# Shannon pipeline plug-in:
## For human mRNA splicing mutations
CLC bio Genomics Workbench plug-in
CLC bio Genomics Server plug-in
### *Features and Benefits*

*Cytognomix* introduces a line of Shannon pipeline plug-ins for prediction of functionally-significant, non-coding variants in genome or exome sequences. Comprehensive genome-scale analysis is now possible for mutations which completely or partially inactivate mRNA splice sites or activate cryptic splicing. The Shannon pipeline plug-in uses our patented and proven information theory-based binding site analysis. Its algorithm has been validated in hundreds of peer-reviewed research studies of splicing mutations, and has been recognized by the American College of Medical Genetics and Genomics in their published guidelines and standards (*Genet. Med.* 7, 571−583).

The CLC Bio Genomics Workbench is used to visualize, export, and further analyze the results of the Shannon pipeline suite.

## Cytognomix Tools

The 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, then annotated against standard databases, and resulting mutations are filtered. Results are displayed as exportable Manhattan-style plots, sortable tables, or browser tracks.

## User Benefits

- The Shannon Human Splicing Pipeline starts with hundreds of thousands of variants, then hones in on the very limited number that potentially alter mRNA splicing.
- Variants are categorized by whether they fully or partially inactivate natural sites or activate cryptic sites proximate to exons or within them.
- Predicts variants missed by many other techniques
- Splicing-related changes are displayed graphically either as "Manhattan-like" plots or as BED tracks.
- Mutations can be sorted according to:
    - the change in information content,
    - the proximity to a natural splice site,
    - the relative strength of cryptic *vs* adjacent natural site,
    - gene and chromosome coordinate.
- Variants that affect known single nucleotide polymorphisms (SNPs) are identified and can be filtered according to allele frequency.
- Results for genome-wide high throughput sequence data obtained in ≤2 hours.

- Fully integrated with CLC-Bio Genomics Workbench and Genomics Server

## Features

- Identifies likely mutations with industry-leading sensitivity and specificity. Trusted legacy of experimentally validated mutation predictions.
- Input accepts variants in VCF, in *Cytognomix*'s simple indexed format, or as CLC Bio variant objects.
- Can accept variants and output results with either hg18 (NCBI36) or hg19 (GRCh37) coordinates.
- Intuitive exportable results based on sensible defaults.
- Run as a standalone application on the Genomics Workbench or configured as a client-server with both the Workbench and Server.
- Fully compatible output for other CLC Bio Workbench Tools.

## Innovative platform for rapid, non-coding mutation analysis

*Cytognomix* uses information theory-based models of mRNA splicing to analyze mutations that alter transcript structure and abundance[1-3]. Information models rank sequences according to their individual information content ($R_i$ in bits). Functional binding sites have $R_i > 0$, corresponding to $\Delta G < 0$ kcal/mol. Strong binding sites have $R_i > R_{sequence}$ while weak sites have $R_i < R_{sequence}$. Variations which alter the affinity of a protein to bind there modify the $R_i$ of the site. A 1 bit change in information content ($\Delta R_i$) corresponds to a $\geq$ 2 fold change in binding affinity. This approach is applicable to any type of nucleic acid binding site, including transcription factors and other conserved non-coding sequences.

Predictions from these models are accurate[2-4], as differences in individual information contents ($\Delta R_i$ in bits) are related to the splicesomal affinities of natural and variant sequences[1,5]. Pathogenicity is related to $\Delta R_i$, which is decreased at natural splice sites and/or increased at cryptic sites [6-7]. Sites with negative $R_i$ values are not recognized. Leaky mutations have modestly reduced $R_i$ values. Cryptic splice sites with $R_i$ values exceeding adjacent natural splice sites are activated[2]. The plug-in reports minimally detectable expression changes of >2 fold [8], or $\Delta R_i > 1$ bit, as significant.

The Shannon pipeline was initially created to address the vexing problem of assessing the many variants of unknown significance that are detected in cancer genetic testing. The pipeline has been used to reanalyze the Breast Cancer Information Core identifying many splicing mutations, most of which were previously unrecognized[9].

The Shannon software pipeline was developed and implemented in C and Perl to perform information analysis fast on a genome-wide scale. Determining $R_i$ values of donor and acceptor sites along a nucleotide sequence is carried out using a convolution-style, sliding-window computation on chromosomes or subsets of chromosomes. $R_i$ values are computed with $R_i(b,l)$ information weight matrix (based on a genome-wide set of verified donor and acceptor splice sites). Variants with significant $\Delta R_i$ values are then filtered, and annotated based on the gene they reside within. For

cryptic splice sites, the distance and relative location of the adjacent natural splice site of the same polarity is reported. Inactivating and leaky natural splice sites and cryptic splicing variants are categorized. Variants with known SNP designations and respective allele frequencies are also reported.

The CLC Bio implementation allows for additional filtering of input, and produces a graphical display of both $\Delta R_i$ and final $R_i$ values for each variant on each chromosome, a table output which can be dynamically sorted that is categorized as separate tab for each type of mutation consequence, and BEDGRAPH output suitable for genome browsing. Results may be exported to a spreadsheet format for further data exploration.

References:
1.Schneider TD. J. Theor. Biol. 189: 427-41, 1997; 2. Rogan PK et al. Hum Mutat 12:153-171, 1998; 3. Rogan PK et al. Pharmacogenetics. 13:207-18, 2003; 4. Rogan PK, Schneider TD.. Hum Mutat 6:74-76, 1995; 5. Gadiraju S, et al. BMC Bioinformatics 4:38, 2003; 6. von Kodolitsch et al Circulation 100:693-9, 1999; 7. von Kodolitsch et al. 12: 258-262, 2006; 8. Nalla VK, Rogan PK. Hum Mut. 25(4):334-42, 2012.; 9. Mucaki et al. Hum. Mut. 32:735-742, 2012.

## Benchmarks

The *Shannon Human Splicing Pipeline* analyzes an average of 3198 variants/min on an I7-based server:

**Performance of Shannon pipeline plug-in on complete genome sequence data**[*]

| Number of variants | Complete analysis time |
| --- | --- |
| 100,000 | 37 m |
| 211,049 | 1h 12 m |
| 290,589 | 1h 22 m |
| 314,637 | 1h 27 m |

**\*hg18/NCBI36**

## Requirements and validation

The Cytognomix Shannon human mRNA splicing plug-in runs in standalone mode on the CLC Genomics Workbench V5.5 or with both the Workbench and CLC Genomics Server V.4.5 (as a standalone server or running Gridworks). Released for **Linux** and **MacOSX** Operating systems supporting Perl and gcc. 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 (Lion release version 10.7.4; gcc v.4.2.1 and Perl 5.12.3) on hardware equipped with an Intel I7 processor and at least 4Gb RAM.

## Support
CLC Bio customer support (primary)
*Cytognomix* Inc. (secondary)
700 Collip Circle #150, London ON N6G 4X8 Canada
info@cytognomix.com ♦ www.cytognomix.com

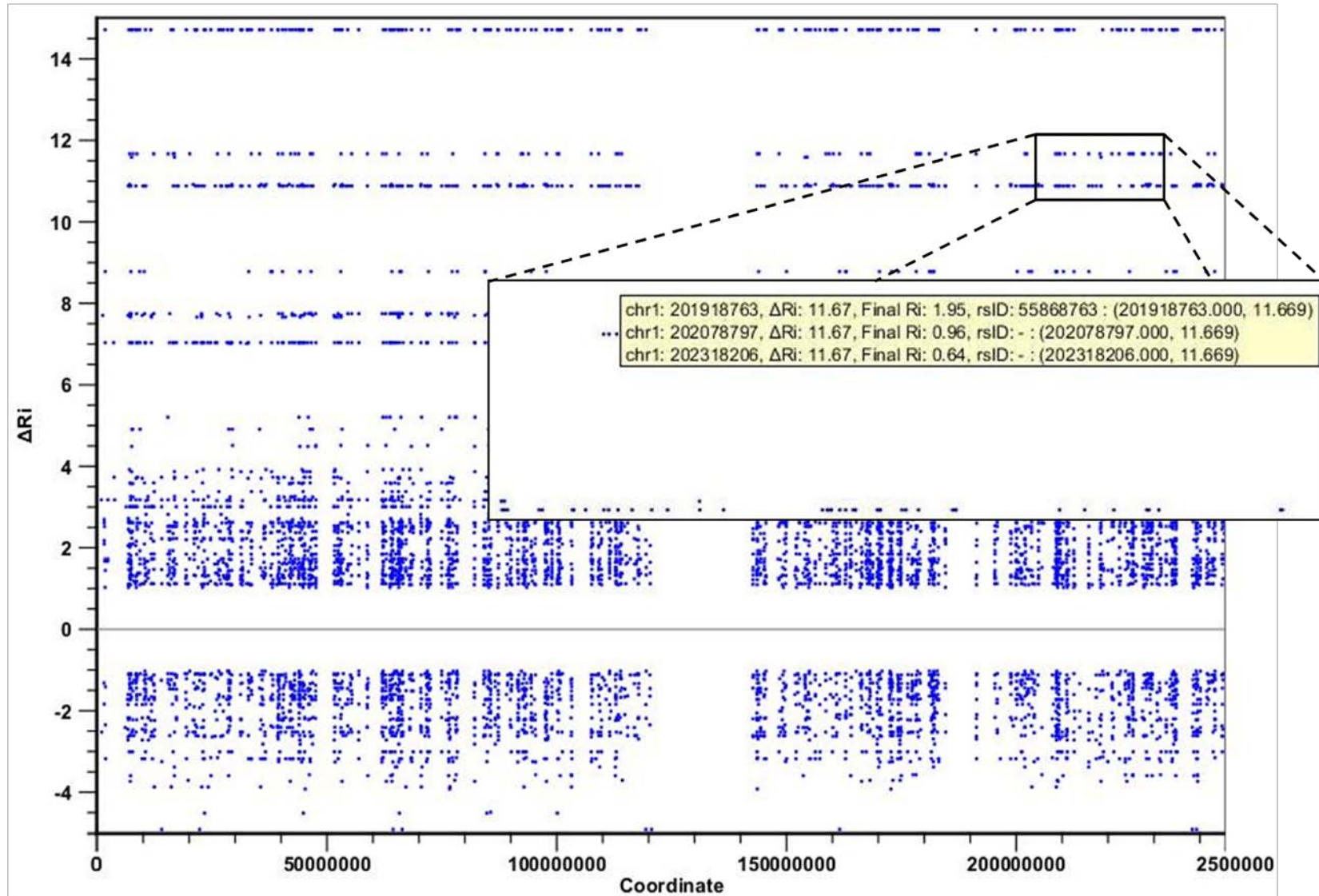Notice: Exclusive marketing rights for products based on **US Pat. No. #5,867,402** have been granted to *CytognomiX*.

Tabular results from genome-wide mutation analysis produced by *CytognomiX*'s Shannon mutation pipeline:

``Manhattan``-like plot output of potential splicing mutations (unfiltered) on chromosome 1 generated by *CytognomiX*'s Shannon mutation pipeline. Inset shows context-dependent mutation detail generated by mouse-over:



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: 202318206, ΔRi: 11.67, Final Ri: 0.64, rsID: - : (202318206.000, 11.669)

Graphical custom genome browser track output produced by *Cytognomix*'s Shannon mutation pipeline indicating information changes of a series of mutations occurring in the *BRCA1* gene (Mucaki et al. Hum. Mut. 32:735-742, 2012):