

Tutorial MAJIQ/Voila (v2.0)

[TOC]

Introduction

What are *MAJIQ* and *Voila* ? {#Soft_Def}

MAJIQ and Voila are two software packages that together define, quantify, and visualize local splicing variations ([LSV](#)) from RNA-Seq data. Conceptually, MAJIQ/Voila can be divided into three modules:

- **MAJIQ Builder:** Uses RNA-Seq (BAM files) and a transcriptome annotation file (GFF/GTF) to define splice graphs and known/novel Local Splice Variations ([LSV](#)).
- **MAJIQ Quantifier:** [Quantifies](#) relative abundance (PSI) of LSVs and changes in relative LSV abundance (delta PSI) between conditions with or without replicates.
- **Voila:** A visualization package that combines the output of MAJIQ Builder and MAJIQ Quantifier using interactive [D3](#) components and HTML5. Voila creates interactive summary files with gene splice graphs, LSVs, and their quantification.

The above three modules are designed to be executed in sequence with one module's output feeding into the other. In most usage cases, the Builder will be executed only once for a given set of RNA-Seq experiments, and then the Quantifier and Voila may be executed on top of it multiple times for different analysis tasks. Note that for samples to be analyzed by the Quantifier **ALL samples analyzed must come from the same execution** of the Builder.

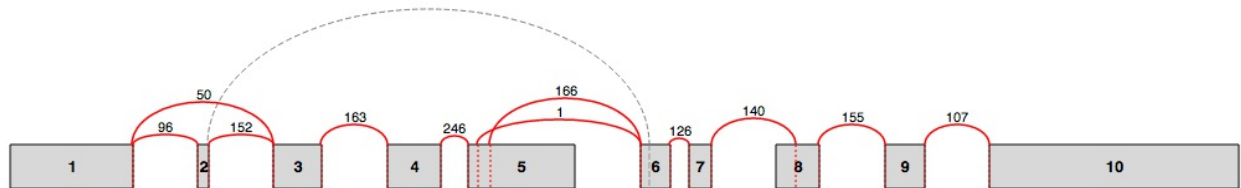
There are two main modes of executing the quantifier: Quantifying the relative inclusion levels of LSVs in a given experimental condition (also known as "*percent selected index*", PSI, or Ψ), and quantifying changes of LSVs inclusion levels between two experimental conditions (aka delta PSI or $\Delta\Psi$).

Voila has two main modes of visualizing the Quantifier's results, whether these are PSI or delta PSI quantifications. The first is a (possibly long) table view of LSVs that can be filtered and ordered by different attributes (columns). The second mode is gene based, in which case each gene's splice graph and matching LSVs are grouped together. In both cases, the experimental condition can be either a single experiment or a set of replicates. In all cases the output is an interactive HTML5 that can be opened in a web browser. There is also an option to dump the output as a tab delimited text file for further analysis with other tools/scripts.

Below there is more information describing what are LSVs, how are they quantified and visualized, and what MAJIQ can (and equally important - cannot) do. You can either go through those or jump directly to the [Quick Start guide](#).

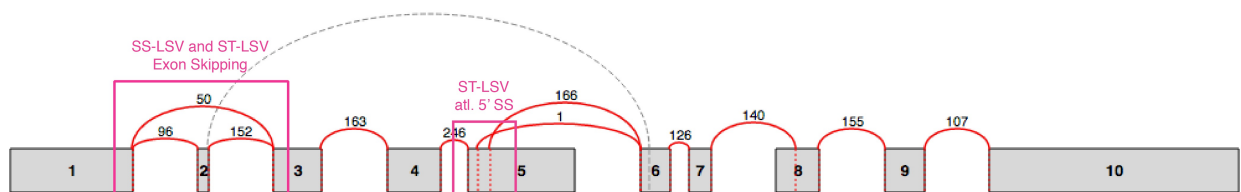
What is an LSV? {#LSV_Def}

LSV stands for "local splicing variation". Briefly, exons that are spliced together can be represented using a *splice graph* [Heber et al. 2002](#) such as this:



In Voila Splice graphs, exons are represented by rectangles and junctions (or edges) by arcs. The raw number of reads spanning a junction is also displayed. For a more detailed description see [splice graphs description](#).

LSVs involve an exon (or node in the splice graph) from which splits in the graph originate (single source LSV, or SS-LSV) or an exon into which several graph edges converge (single target LSV, or ST-LSV). An illustration of a splice graph is shown below with several SS-LSV and ST-LSV marked. The "local" aspect of LSVs definition stems from the fact they involve only a single source or single target exon. For a more formal definition please see [Vaquero-Garcia et al., 2016](#).



The terminology of LSVs generalizes that of alternative splicing (AS) "events" and AS events "types". The most common AS types in mammals are skipped exons, alternative 3' splice site, and alternative 5' splice site ([Wang et al., Nature 2008](#)). These types can all be seen as specific cases of simple or binary LSVs, i.e. LSVs that involve only two way graph splits (see figure above). However, we find that the transcriptome contains many other types of LSVs that involve different combinations of 3' and 5' splice site choices in different exons (see figures above/below, or just run your data through MAJIQ/VOILA...). Consequently, the LSV terminology helps us define and quantitate more accurately the spectrum of local splice variations observed in the transcriptome, many of which are complex.

Conceptually, LSVs are aimed to fill the gap between previously defined AS "types" described above, and full transcripts/isoforms. Ideally, we would like to identify and quantify all existing

isoforms of each gene in a given RNA-Seq experiment. However, the complexity of gene isoforms combined with the shortness of current sequencing reads (typically ~100 nt long) makes isoform quantification from RNA-Seq reads a challenging problem. In contrast, LSVs can arguably still capture a lot of useful information about transcriptome variability while being deduced directly from RNA-Seq reads that span across splice junctions.

What is LSV quantification? {#LSV_Quant}

MAJIQ's LSV quantification is based on estimating the relative inclusion level of each junction in the LSV. For simple, binary, cases such as skipped exons, LSV quantification is equivalent to estimating the exon's percent spliced in (PSI, or Ψ). For more complex LSVs that involve three or more splice graph edges (i.e., exon joining options), *MAJIQ* computes the marginal inclusion level, or PSI, per junction. Computing only these marginals allows *MAJIQ* to handle complex LSVs, keeping computational cost linear with the number of edges while still delivering estimates for the interesting biological question of "how much is each junction used?".

When estimating PSI for a LSV's junctions, *MAJIQ* produces a complete posterior distribution over possible PSI values. This distribution takes into account the number of reads observed at each junction, their distribution across genomic positions, GC content bias, and some possible mapper or technical artifacts. Intuitively, the deeper and smoother the coverage of an LSV, the more concentrated the PSI posterior would be (i.e. the more "sure" *MAJIQ* is about the "true" PSI value), while lower and less even coverage would result in higher variance of the PSI estimate.

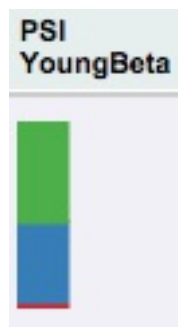
Similarly, *MAJIQ*'s quantification of LSV's differential inclusion when comparing two conditions is based on estimating a posterior distribution for the change in each junction's relative inclusion level, termed delta PSI ($\Delta\Psi$). Naturally, this distribution lies in the range of -100% to +100% (or -1 to +1 when using fractions instead of percentages).

**Note:* For a thorough description of *MAJIQ*'s quantification algorithm for Ψ and $\Delta\Psi$ and the various parameters that control it see [Vaquero-Garcia et al., 2016](#).

Voila's visualization of LSV Quantification uses several different techniques. In all cases, colors are used to represent the different junctions in the LSV. *Voila* uses *violin plots* to represent the posterior probability distribution for Ψ or $\Delta\Psi$. Examples of the violin plots for single Ψ are shown below.

PSI

Compact view



When displaying lists of LSVs *Voila* uses a compact stacked bar chart representation. The **height** of each bar represent the **expected PSI** ($E[\Psi]$) which naturally add to 100% over all the LSV's junctions. Clicking over the bars will open a more detailed representation of the PSI distribution.

Violin plots (binary and multi-way LSV)



Violin plots are *boxplots* plotted over the original distributions. The *box* goes from the 25th to the 75th percentile, with a white horizontal line indicating the 50th percentile (median). The tails represents the 10th and 90th percentile. Additionally, the expected PSI or $E[\Psi]$ is marked with a white circle.

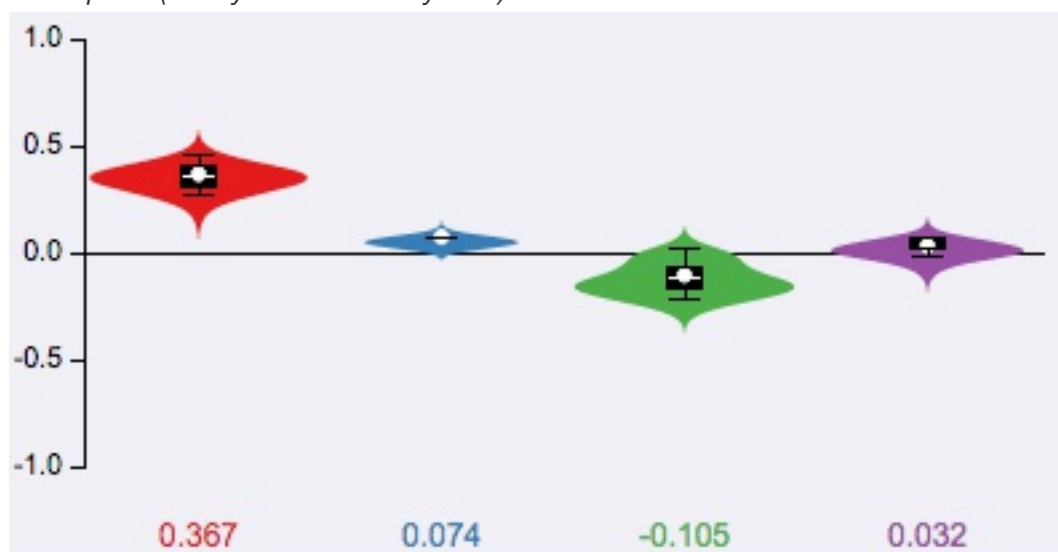
Delta PSI

Compact view



For compact visualization of ΔPsi quantification, each colored bar represents the percentage of the expected differential inclusion for the matching edge in the splice graph. The arrows indicate the preference for one condition versus another. In this example, the blue junction is shown to have 40% more inclusion in *YoungBeta* compared with *OldBeta*. In contrast the green junction is expected to be more included in *OldBeta* by a 35% difference. Note that in this case the numbers do not generally sum to 1 or 100% as they reflect MAJIQ's expected change for each junction separately. In addition, users can use the more detailed violin plots to gain other information/statistics such as the confidence and probability distribution per junction.

Violin plots (binary and multi-way LSV)



Each violin correspond to the posterior distribution of a junction in the LSV (not shown here) for delta PSI analysis of condition1 Vs condition2. For each junction, the expected delta PSI or $E(\Delta\text{Psi})$ is shown at the bottom, where **negative** values correspond to **increased differential inclusion in condition1** compared with condition2 whereas a positive $E(\Delta\text{Psi})$ denotes preference for condition2 Vs condition1.

. In this example, the purple junction (the actual LSV is not shown here) is shown to have a 75% confidence, based on the observed reads, that its inclusion in *condition 1* is increased by

at least 20% compared with *condition 2*. In contrast the red and blue junctions have confidence of 28% and 38% respectively of dropping by at least 20% (i.e. increased exclusion) in *condition 2*. Note that in this case the numbers do not generally sum to 1 or 100% as they reflect MAJIQ's confidence for each junction separately that its inclusion was increased/decreased by more than the given threshold. In addition, users can use the more detailed violin plots to gain other information/statistics such as the expected $\Delta \Psi$ per junction.

-->

What is MAJIQ?

MAJIQ is a software package that allows researchers to define and quantify both known and novel Local Splice Variations ([LSVs](#)) in genes from RNA-Seq data.

MAJIQ's main features

MAJIQ takes as input a set of RNA-Seq experiments (sorted, indexed BAM files) and previous genome annotation ([GFF3 files](#)) and produces the following:

- Splice graph for each gene based on both known transcripts annotation and de-novo junctions detected.
- All detected (known + novel) single source and single targets [LSVs](#) per gene.
- [Quantification of LSVs](#) from a given RNA-Seq experiment (w/wo replicates).

New Majiq changes (v1.1.x)

The new version of Majiq v1.1.x is a reimplementaion of the software in order to achieve a faster and memory efficient execution. The structure has been refactored be able to work with huge datasets with a lower memory imprint.

This new version uses python ≥ 3.5 implementation and cython modules.

The output has been highly reduced with smaller and faster output files and removing the creation of temporary files.

All this has been implementing keeping all the Majiq functionalities like complexity quantification, denovo junctions/exon detection and visualization, that made MAJIQ stands out from other RNASeq differential splicing tools.

We added the possibility to store back in disc the annotation DB complementing it with the denovo elements found in the data, this enriched DB can be feed as input to MAJIQ on future runs.

What MAJIQ is not {#MAJIQ_Not}

There are many RNA-Seq analysis tasks for which *MAJIQ* was *not* designed or is currently not structured to address. Some examples include:

- Gene/isoform expression estimation: *MAJIQ* uses expression levels when it quantifies LSVs. For example, LSVs that are not present in the data will not be quantified and those with lower coverage will result in lower confidence for the relative inclusion of matching RNA segments. However, *MAJIQ* only computes *relative* inclusion of junctions in an LSV (e.g., 80% inclusion of an alternatively skipped exon). Consequently, computing the expression of genes or isoforms and comparing those in the same experiment or between experiments is not supported.
- Relative isoform abundance: *MAJIQ* only operates at the level of local splice variations (LSVs). It does not assume the full spectrum of gene isoforms is known and does not quantify those.
- Novel gene/non coding RNA detection: *MAJIQ* requires a transcriptome annotation file ([GFF3](#)). It supplements by identifying both known and novel splice junctions within the bounds of existing loci in the provided annotation. Putative isoforms of new loci are not inferred during this process.
- Alternative transcription start/end
- Alternative polyadenylation (APA) identification/quantification.

What is Voila?

Voila is a package to interactively visualize splice variations in RNA-Seq data. It is written in Python and produces summary files in HTML5 that can be opened and interactively explored with any modern browser*. It has been conceived as the visual component of *MAJIQ* for analysis of Local Splice Variants (LSVs).

**Voila has been tested on Google Chrome [recommended], Firefox, and Safari.*

How to cite us?

Primary Publication

If you use *MAJIQ* in your published work, please cite this publication:

Vaquero-Garcia J, Barrera A, Gazzara MR, et al. A new view of transcriptome complexity and regulation through the lens of local splicing variations. eLife, 5, e11752.

<http://doi.org/10.7554/eLife.11>

Quick start {#QStart}

Pre MAJIQ

Select a GFF3 annotation file

The general feature format (gene-finding format, generic feature format, GFF) is a file format used for describing genes and other features of DNA, RNA, and protein sequences. The format specification for the gff version 3 can be found at [GFF3 format](#).

In our case we use some of these features in order to define genes, transcripts and exons. An example of this format is shown below

```
chr1 protein_coding gene 107399655 107452689 . + .
Name=Serpib7;ID=ENSMUSG000000067001;Name=ENSMUSG000000067001
chr1 protein_coding mRNA 107399655 107435399 . + .
Parent=ENSMUSG000000067001;Name=Serpib7-002;ID=ENSMUST00000154538
chr1 protein_coding exon 107399655 107399724 . + .
Parent=ENSMUST00000154538;ID=exon:ENSMUST00000154538:1
chr1 protein_coding exon 107428231 107428416 . + .
Parent=ENSMUST00000154538;ID=exon:ENSMUST00000154538:2
chr1 protein_coding exon 107434736 107434786 . + .
Parent=ENSMUST00000154538;ID=exon:ENSMUST00000154538:3
chr1 protein_coding exon 107435327 107435399 . + .
ID=exon:ENSMUST00000154538:4;Parent=ENSMUST00000154538
chr1 protein_coding five_prime_UTR 107399655 107399724 . + .
Parent=ENSMUST00000154538;ID=five_prime_UTR:ENSMUST00000154538:1
chr1 protein_coding five_prime_UTR 107428231 107428248 . + .
ID=five_prime_UTR:ENSMUST00000154538:2;Parent=ENSMUST00000154538
chr1 protein_coding start_codon 107428249 107428251 . + 0
Parent=ENSMUST00000154538;ID=start_codon:ENSMUST00000154538:1
chr1 protein_coding CDS 107428249 107428416 . + 0
ID=CDS:ENSMUST00000154538:1;Parent=ENSMUST00000154538
chr1 protein_coding CDS 107434736 107434786 . + 0
Parent=ENSMUST00000154538;ID=CDS:ENSMUST00000154538:2
...
```

It is important to note that MAJIQ makes some assumptions when parsing the hierarchical GFF3 file and currently has some specific requirements:

- We only consider sequence features with the type (column 3) “gene”
- For every gene, we only consider isoforms of a gene with a type of “mRNA” or “transcript”

- All entries (except for genes) should have a parent attribute
- All genes should have a unique ID attribute
- Within a gene all entries should have a unique ID attribute.
- A gene can have a Name attribute, otherwise the ID will be used instead in the output.

Keeping these in mind will be important for analyzing the types of transcripts you care about and modifying your GFF3 annotation file may be necessary.

In order to obtain this format, we recommend the use of some of the most well known online DB. They provide the annotation files in some format like *GTF*, and you can transform this file to *GFF3* using a script, like this [script](#)

You can also download the annotation files used in [Vaquero-Garcia et al., 2016](#) for the Ensembl hg19 or mm10 genome builds [here](#).

Study configuration file

MAJIQ has a set of parameters needed for its execution. Several of them depend of the RNA-Seq study. This configuration file should include this information in order to be able to pass it the the MAJIQ Builder. Secondly, it is useful to keep the info of the study ready and accessible.

This is an example of the configuration file, divided in two blocks, *info* and *experiments*:

```
[info]
readlen=76
samdir=/data/MGP/ERP000591/bam
genome= mm10 strandness=forward[reverse|None]
[experiments]
Hippocampus=Hippocampus1,Hippocampus2
Liver=Liver1,Liver2
[optional]
Hippocampus1=strandness:None,
Liver2=strandness:reverse,
```

Info

This is the study global information needed for the analysis. The mandatory fields are:

- *readlen*: Length of the RNA-seq reads. MAJIQ *can* handle experiments with multiple read lengths, just indicating the longest read length
- *samdir*: Path where the bam files are located
- *genome*: Genome assembly
- *strandness=forward[reverse|none]*: Some of the RNASeq preparations are strand specific.

This preparations can be reverse strand specific[**reverse**], forward strand specific [**forward**], or non-strand specific [**none**]. This parameter is optional, in which case **None** is used as default.

Experiments

This section defines the experiments and replicates that are to be analyzed. Each line defines a condition and its name can be customized in the following way:

```
<group name>=<experiment file1>[,<experiment file2>]
```

where the experiment file is the sorted bam filename inside the samdir directory (excluding extension *.bam*). MAJIQ expects to find within the same directory bam index files for each experiment with the format `<experiment file>.bam.bai`.

Multiple replicates within an experiment should be comma separated with no spaces.

Optional

This section, newly introduced in **v1.1.x**, allows the user to specify changes in the global parameters specific to single experiments. The syntax goes as follows,

```
<experiment file1>=<option1>:<value1>,...,<optionN>:<valueN>
```

The user only need to add the experiments that have different parameter value than the global parameters, using this section only when is needed.

Currently only **strandness** has been adapted to be overwrite using this optional section, but new options can be added in the future.

MAJIQ Builder

MAJIQ Builder is the part of MAJIQ tool where RNA-Seq data is analyzed in order to detect LSV candidates.

All conditions and replicates that will be analyzed with MAJIQ PSI or delta PSI should be executed TOGETHER in a single Builder execution.

```
maji build <transcript list> -c <configuration file> -js NT -o <build outdir>
```

- **Transcriptome annotation:** This is the file with the annotation database. Currently, we accept only [GFF3 format](#). For a better description, see the annotation file section.
- **Configuration file:** This is the configuration file for the study. This file should define the files and the paths for the bam files, the read length, the genome version, and some other information needed for the builder. For a more detailed information, please check the

[configuration file](#) section.

- **NT**: Number of threads to use.
- **Build outdir**: Directory where the output will be placed. MAJIQ builder has a set of output files including one *.maji* for each bam file and one *splicegraph.sql*. These files will be the input files in the next steps of the analysis.

MAJIQ Builder has several arguments in order to tweak its analysis and performance. Please check the [MAJIQ parameters section](#) for a more detailed explanation.

PSI Analysis

PSI quantification

MAJIQ PSI quantifies the LSV candidates given by the Builder. In order to improve its accuracy and reproducibility, it allows the use of biological replicates.

```
maji psi <build outdir>/<replicate1>.maji [<build outdir>/<replicate2>.maji ...] -j NT -o  
<psi outdir> -n <cond_id>
```

- ***.maji file[s]**: the path to the *.maji* file(s) that were created by the MAJIQ Builder execution.
- **cond_id**: group identifier that you want to use for this execution
- **NT**: Number of threads to use.

Please check the [MAJIQ parameters section](#) for a more detailed explanation of all the arguments.

Visualize results with VOILA

The package VOILA allows the user to generate interactive summaries to display MAJIQ computations and quantifications in the browser. All the information is also provided in TAB-delimited files that can be easily parsed for further analysis.

```
voila psi <psi outdir>/<cond_id>.psi.voila --splice-graph <build outdir>/splicegraph.sql -o  
<voila outdir>
```

- **<cond_id>.psi.voila** is the output file from MAJIQ PSI computation
- **splicegraph.hdf5** contains information about the genes and splice variants identified in the MAJIQ Builder

In the output directory **<voila outdir>** you will find:

- **index.html**: HTML file with a table containing all genes and LSVs identified and analyzed.
- **summaries/xx_<cond_id>.psi.html** files: interactive HTML5 summaries with MAJIQ quantifications, where **xx** is the page counter. Ten genes are displayed per page
- **<cond_id>.psi.tsv**: A tab-delimited file with all LSV information (expected PSI value, variance, exon coordinates, junction coordinates, etc.) and genomic information (chromosome, strand and coordinates).
- **static** folder: needed for the correct visualization of the `index.html` file.
- **voila.log**: log file with the execution information of Voila.

For more information see [VOILA](#) section.

For additional command line arguments use `majiq psi -h` or `voila psi -h`

Delta PSI Analysis

Delta PSI quantification

Majiq Delta PSI quantifies the differential splicing between two different groups (or conditions). Like PSI, Delta PSI is able to use replicates for each group in order to improve its accuracy and reproducibility.

```
majiq deltapsi -grp1 <build outdir>/<cond1_rep1>.majiq [<build outdir>/<cond1_rep2>.majiq
...] -grp2 <build outdir>/<cond2_rep1>.majiq [<build outdir>/<cond2_rep2>.majiq ...] -j NT -o
<dpsi outdir> -n <cond1_id> <cond2_id>
```

- **-grp1 .majiq file[s]**: Set of *.majiq* file[s] for the first condition,
- **-grp2 .majiq file[s]**: Set of *.majiq* file[s] for the second condition,
- **--name cond_id1 cond_id2**: group identifiers for *grp1* and *grp2*, respectively, used for naming output files
- **NT**: Number of threads to use.

Please check the [MAJIQ parameters section](#) for a more detailed explanation of all the arguments.

Visualize results with VOILA

To visualize deltapsi quantification with Voila execute:

```
voila deltapsi <dpsi outdir>/<cond1_id>_<cond2_id>.deltapsi.voila --splice-graphs <build
outdir>/splicegraph.sql -o <voila outdir>
```

- **<dpsi_outdir>/<cond1_id>_<cond2_id>.deltapsi.voila** is the output file from delta PSI computation,

- **splicegraph.hdf5** contains information about the genes and splice variants identified in the MAJIQ Builder

In the output directory `<voila outdir>` you will find:

- **index.html**: HTML file with a table containing all genes and LSVs identified and analyzed and links to more detailed summaries.
- **summaries/xx_<cond1_id>_<cond2_id>.deltapsi.html** files: interactive HTML5 summaries with MAJIQ quantifications.
- **<cond1_id>_<cond2_id>.deltapsi.tsv**: A tab-delimited file with all the genes and LSV quantifications and genomic information.

By default VOILA uses a threshold of a change of $|dPSI| \geq 0.2$ (20%) between conditions. To change this threshold you can use the option `--threshold` and specify a fraction from 0 to 1. To show all LSVs `--show-all` can be used instead.

For additional command line arguments use `majiq deltapsi -h` or `voila deltapsi -h`

MAJIQ

Majiq Parameters {#majiq_params}

In the [quick start](#) section above we described a general execution pipeline for MAJIQ, but those three commands have many other parameters that can be adjusted to modify MAJIQ's behavior.

Builder

```
majiq build [-h] transcripts -c CONF --output OUTDIR [-j NTHREADS] [--silent] [--debug] [--min-experiments MIN_EXP] [--k K_samples] [--m M_samples] [--min-denovo MIN_DENOVO] [--minreads MINREADS] [--min-intronic_cov MIN_INTRONIC_COV] [--minpos MINPOS] [--disable-denovo] [--disable-ir] [--markstacks PVALUE_LIMIT]
```

Mandatory arguments:

- `transcripts` : Transcriptome file with the annotation database. Currently, we accept [GFF3 format](#). For a better description, see the annotation file section. transcripts can be a majiq DB file generated with `DB.npz` name from a previous majiq build execution.

- `-c/--conf CONFIG_FILE` : This is the configuration file for the study. This file should define the files and the paths for the bam files, the read length, the genome version, and some other information needed for the Builder. For a more detailed information, please check the [configuration file](#) section.
- `-o/--output OUTDIR` : Directory where the output will be placed. MAJIQ Builder has a set of output files *.majiQ* per each bam file and one *splicegraph.sql*. These files will be the input files in the next steps of the analysis.

Optional arguments:

- `-h, --help` : Show help message and exit
- `-j/--nproc NTHREADS` : Number of threads to use.
- `--prebam` : prebam option will assume that the bam analysis was done before in the specified output folder. Causes the Builder to skip redoing this step and look for the temporary files that should be already generated.
- `--disable-denovo` : Avoid *de novo* detection of junction, splice-sites, and exons. This will speedup the execution but reduce the number of LSVs detected.
- `--disable-ir` : Avoid *intron retention* detection. This will speedup the execution but reduce the number of LSVs detected.
- `--k K` : Number of positions to sample per iteration. [Default: 50]
- `--m M` : Number of bootstrapping samples. [Default: 100]
- `--minreads MINREADS` : Minimum number of reads threshold combining all positions in a LSV to consider that the LSV "exists in the data". [Default: 3]
- `--minpos MINPOS` : Minimum number of start positions with at least 1 read in a LSV to consider that the LSV "exists in the data" [Default: 2]
- `--min-intronic_cov MIN_INTRONIC_COV` : Minimum number of reads on average in intronic sites, only for intron retention. [Default: 1.5]
- `--min-experiments MIN_EXP` : Used to lower the threshold for group filters. min_experiments is the minimum number of experiments where the different filter checks in order to pass an lsv or junction.
- `--min-denovo MIN_DENOVO` : Minimum number of reads threshold combining all positions in a

LSV to consider that denovo junction is real". [Default: 2]

- `--markstacks MARKSTACKS` : Mark stack positions. Expects a p-value. Use a negative value in order to disable it. [Default: 1e-07]

Logger arguments:

- `--logger LOGGER_PATH` : Path for the logger. Default is output directory
- `--silent` : Boolean argument used to silence the logger.
- `--debug` : Activate this flag to activate debug messages.

PSI

```
`majiQ psi [-h] [-j NTHREADS] -o OUTDIR -n NAME  
[--logger LOGGER] [--silent] [--debug]  
[--min-experiments MIN_EXP] [--minreads MINREADS] [--minpos MINPOS]
```

files [files ...]

Mandatory arguments:

- `files` : *majiQ* file[s] that were created by the MAJIQ Builder execution
- `-n/--name NAME` : The name that identifies the quantification group.

- `-o/--output OUTDIR` : PSI output directory. It will contain the *psi.voila* file once the execution is finished.

Optional arguments:

- `-h, --help` : Show help message and exit
- `-j/--nprocs NTHREADS` : Number of threads to use.
- `--minreads MINREADS` : Minimum number of reads to pass the quantifiable threshold combining all positions in a LSV to considered. [Default: 10]
- `--minpos MINPOS` : Minimum number of start positions with at least 1 read in a LSV to considered. [Default: 3]
- `--min-experiments MIN_EXP` : Use to alter the threshold for group filters. `min_experiments` is the minimum number of experiments where the different filter checks must be met in order to consider LSV or junction quantifiable.

Logger arguments:

- `--logger LOGGER_PATH` : Path for the logger. Default is output directory
- `--silent` : Boolean argument used to silence the logger.
- `--debug` : Activate this flag to activate debug messages.

DeltaPSI

```
`majiq deltapsi -grp1 FILES1 [FILES1 ...] -grp2 FILES2 [FILES2 ...]
-n NAMES [NAMES ...] -o OUTDIR [-h] [-j NTHREADS]
[--logger LOGGER] [--silent] [--debug]
[--min-experiments MIN_EXP] [--minpos MINPOS] [--minreads MINREADS]
[--binsize BINSIZE] [--default-prior] [--prior-minreads PRIORMINREADS]
[--prior-minnonzero PRIORMINNONZERO] [--prior-iter ITER]
```

Mandatory arguments:

- `-grp1 FILES1 [FILES1 ...]` : Set of *.majiq* file[s] for the first condition
- `-grp2 FILES2 [FILES2 ...]` : Set of *.majiq* file[s] for the second condition
- `-n/--names NAMES [NAMES ...]` : *cond_id1 cond_id2*: group identifiers for grp1 and grp2 respectively.
- `-o/--output OUTDIR` : PSI output directory. It will contain the *deltapsi.voila* file once the execution is finished.

Optional arguments:

- `-h, --help` : Show help message and exit
- `-j/--nprocs NTHREADS` : Number of threads to use [Default: 4].
- `--minreads MINREADS` : Minimum number of reads to pass the quantifiable threshold combining all positions in a LSV to considered. [Default: 10]
- `--minpos MINPOS` : Minimum number of start positions with at least 1 read in a LSV to considered. [Default: 3]
- `--min-experiments MIN_EXP` : Use to alter the threshold for group filters. *min_experiments* is the minimum number of experiments where the different filter checks must be met in order to consider LSV or junction quantifiable.
- `--binsize BINSIZE` : The bins for PSI values. With a BINSIZE of 0.025 (default), we have 40 bins
- `--default-prior` : Use a default prior instead of computing it using the empirical data
- `--prior-minreads PRIORMINREADS` : Minimum number of reads combining all positions in a junction to be considered (for the 'best set' calculation). [Default: 20]
- `--prior-minnonzero PRIORMINNONZERO` : Minimum number of positions for the best set.
- `--prior-iter ITER` : Max number of iterations of the EM

Logger arguments:

- `--logger` `LOGGER_PATH` : Path for the logger. Default is output directory
- `--silent` : Boolean argument used to silence the logger.
- `--debug` : Activate this flag to activate debug messages.

VOILA

HTML5 Summaries

View

Here is the command line usage statment output from `voila view --help`.

```
usage: voila view [-h] [-p PORT] [--force-index] [-j NPROC] [--debug]
                  [-l LOGGER] [--silent]
                  files [files ...]

positional arguments:
  files                List of files or directories which contains the splice
                       graph and voila files.

optional arguments:
  -h, --help            show this help message and exit
  -p PORT, --port PORT  Set service port. Default is a random.
  --force-index          Create index even if already exists.
  -j NPROC, --nproc NPROC
                       Number of processes used to produce output. Default is
                       half of system processes.
  --debug
  -l LOGGER, --logger LOGGER
                       Set log file and location. There will be no log file
                       if not set.
  --silent              Do not write logs to standard out.
```

Here is a sample usage of the voila view command. The *files* positional argument just need the directories or files

location for the splice graph and voila files. Voila will detect if the voila file is for psi or delta psi.

```
$ voila view splicegraph.sql deltapsi.voila
```

TSV

The usage statement for `voila tsv --help` is the following:

```
usage: voila tsv [-h] -f FILE_NAME [--threshold THRESHOLD]
               [--non-changing-threshold NON_CHANGING_THRESHOLD]
               [--probability-threshold PROBABILITY_THRESHOLD] [--show-all]
               [--lsv-types-file LSV_TYPES]
               [--lsv-types [LSV_TYPES [LSV_TYPES ...]]]
               [--lsv-ids-file LSV_IDS] [--lsv-ids [LSV_IDS [LSV_IDS ...]]]
               [--gene-names-file GENE_NAMES]
               [--gene-names [GENE_NAMES [GENE_NAMES ...]]]
               [--gene-ids-file GENE_IDS]
               [--gene-ids [GENE_IDS [GENE_IDS ...]]] [-j NPROC] [--debug]
               [-l LOGGER] [--silent]
               files [files ...]
```

positional arguments:

files List of **files** or **directories** which contains the splice graph and voila **files**.

optional arguments:

-h, --help *show this help message and exit*

--threshold THRESHOLD Filter out LSVs with no junctions predicted to change over a certain **value**. Even when show-all is used this **value** is still used to calculate the probability in the TSV. The default is **"0.2"**.

--non-changing-threshold NON_CHANGING_THRESHOLD The default is **"0.05"**.

--probability-threshold PROBABILITY_THRESHOLD This is off by default.

--show-all *Show all LSVs including those with no junction with significant change predicted.*

--lsv-types-file LSV_TYPES Location of **file** that contains a list of LSV types which should remain in the results. One type per **line**

--lsv-types [LSV_TYPES [LSV_TYPES ...]] LSV types which should remain in the results

--lsv-ids-file LSV_IDS Location of **file** that contains a list of LSV IDs which should remain in the results. One ID per **line**.

--lsv-ids [LSV_IDS [LSV_IDS ...]] LSV IDs, separated by spaces, which should remain in the results. e.g LSV_ID1 LSV_ID2 ...

--gene-names-file GENE_NAMES Location of **file** that contains a list of common gene

names which should remain in the results. One name per line.

`--gene-names [GENE_NAMES [GENE_NAMES ...]]`
Common gene names, separated by spaces, which should remain in the results. e.g. GENE1 GENE2 ...

`--gene-ids-file GENE_IDS`
Location of file that contains a list of gene IDs which should remain in the results. One name per line.

`--gene-ids [GENE_IDS [GENE_IDS ...]]`
Gene IDs, separated by spaces, which should remain in the results. e.g. GENE_ID1 GENE_ID2 ...

`-j NPROC, --nproc NPROC`
Number of processes used to produce output. Default is half of system processes.

`--debug`

`-l LOGGER, --logger LOGGER`
Set log file and location. There will be no log file if not set.

`--silent`
Do not write logs to standard out.

required named arguments:

`-f FILE_NAME, --file-name FILE_NAME`
Set the TSV file's name and location.

For the voila command you will have to supply a filename for the tsv. Again, psi or deltappsi will be detected from

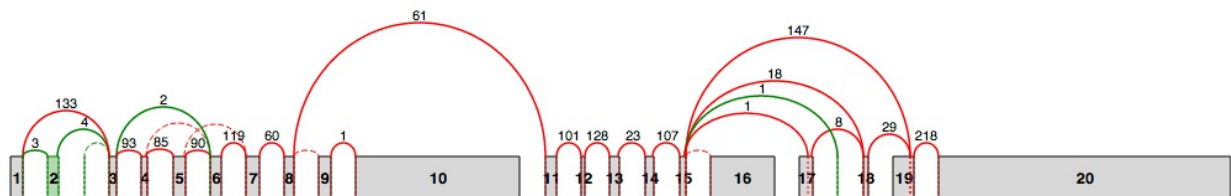
the voila file supplied. A sample command for voila tsv would be something similar to:

```
$ voila tsv splicegraph.sql psi.voila -f psi.tsv
```

Splice Graphs {#Sgraph}

Gene name: Tom1l2; chr11:-:60226714-60352905; 
Gene ID: ENSMUSG00000000538;


Cer_CT28.mm10.sorted.splicegraph  **Coordinates: 60252831-60254816; Length: 1986**



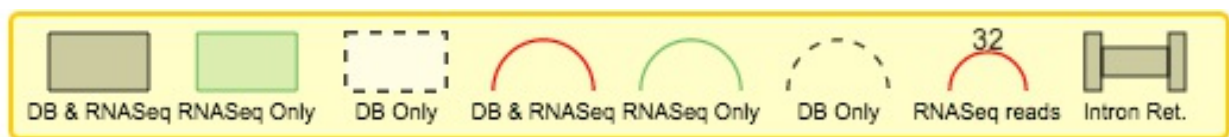
The splice graph gadget included in VOILA summarizes all the splice variants found in a gene by MAJIQ. Splice Graphs include the following features:

- Easy differentiation between exons and junctions annotated (red) and *de novo* detected (green) in RNA-Seq data,
- Contextual information about the coordinates for each exon and intron (hovering over

exons/junctions).

- Raw reads counts for each junction. (Note MAJIQ applies several normalization factors to these raw values)
- Scaled view of the gene and the splice graph . By default, introns are trimmed to obtain a more compact representation of the splice graph. Switching between scaled and default view is accomplished by clicking on the wand.
- Zoom in/out to explore complex splice graphs.
- Possibility to switch between replicates (condition members) when more than one splice graph is available via the dropdown box.




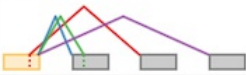





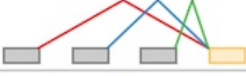


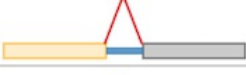


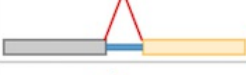





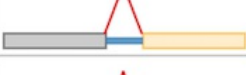


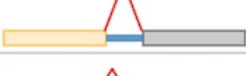





All *Gene Summaries* include a descriptive legend of what you might find in the Splice Graphs:



DB refers to exons and junctions annotated in the GFF3 file and *RNASeq* to exons and junctions found in RNA-Seq data. *RNASeq reads* alludes to the raw reads found in RNA-Seq data. Please, note that retained introns (narrow rectangles connecting two exons) do not appear currently in the legend.

Note that when MAJIQ creates the splice graphs that VOILA visualizes, it considers the bounds of each individual exon in all transcripts containing that exon and the longest version is represented in the splice graph. So in the above example, the annotation database had longer versions of exons 17 and 19 corresponding to alternative transcription start sites for this gene and thus the starting positions are extended to reflect this.

PSI Summary

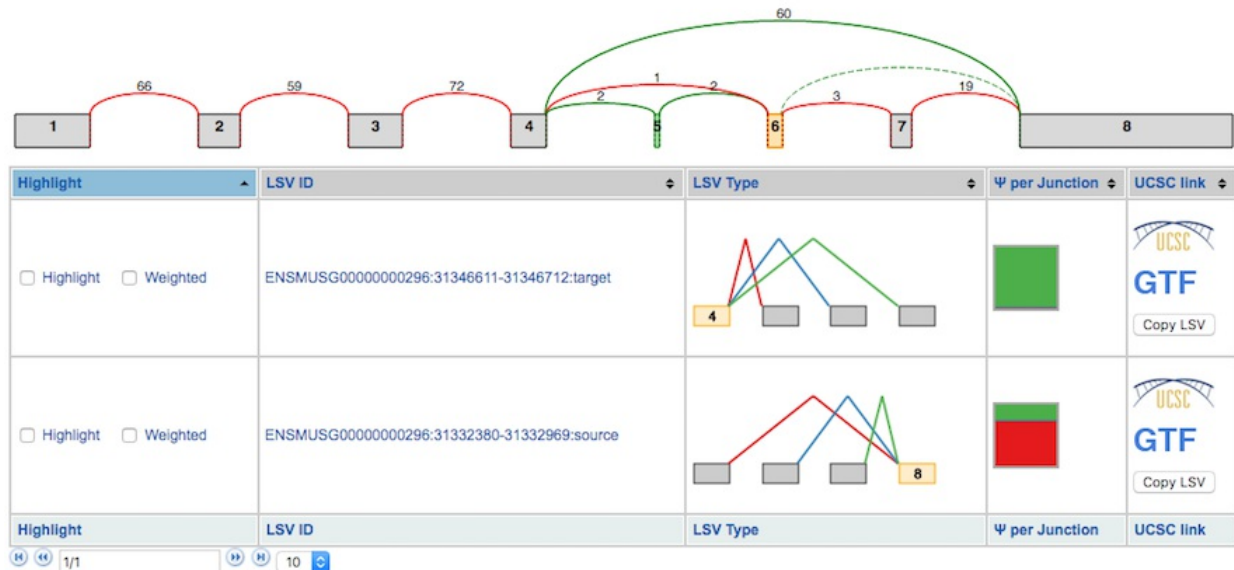
#	Gene	LSV ID	LSV Type	Ψ per Junction	Coordinates Links
1	Zranb2	ENSMUSG00000028180:157544964-157545122:source			
2	Zranb2	ENSMUSG00000028180:157536334-157536386:source			
3	Zranb2	ENSMUSG00000028180:157546675-157548390:target			
4	Zranb2	ENSMUSG00000028180:157537587-157537669:target			
5	Ap3d1	ENSMUSG00000020198:80708747-80708868:target			
6	Ap3d1	ENSMUSG00000020198:80708458-80708537:source			
7	Ap3d1	ENSMUSG00000020198:80710089-80710467:target			
8	Ap3d1	ENSMUSG00000020198:80709201-80709286:source			
9	Ap3d1	ENSMUSG00000020198:80709400-80709488:target			
10	Idh3g	ENSMUSG0000002010:73780213-73780315:source			
#	Gene	LSV ID	LSV Type	Ψ per Junction	Coordinates Links

1/1687 10

VOILA PSI Index file offers an overview of all the LSVs detected in a table, providing links to detailed summaries of LSVs and genes. Clicking over a gene or LSV ID opens up a new tab with a summary containing interactive splice graphs, distributions of PSIs per junction and links to the UCSC. Below is an example of PSI summaries for *Tpd52l1*:

Gene name: **Tpd52l1**; chr10:-:31332380-31445921; 
Gene ID: **ENSMUSG00000000296**;

Combined  Coordinates: 31341140 - 31341178; Length: 39



Tip: PSI Summaries can be navigated through the *Previous* and *Next* links, without having to go back to the *index* file.

The information is broken into genes (10 per page), each of them with an interactive splice graph and an associated table with LSV quantification data. The table has the following information about the LSV:





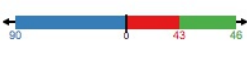





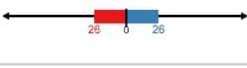




- **Highlight:**
 - **Highlight.** To move the visualization of the LSV to the splice graph, select the “highlight” checkbox in the LSV table. All junctions will be hidden except the ones included in the LSV. The visible junctions will be colored to match the highlighted LSV.
 - **Weighted.** To weight the junctions of the highlighted LSV, check the “weighted” checkbox. The junctions for this LSV will now be weighted proportionally to the LSV’s expected PSI value for this experiment.
- LSV ID: a unique identifier for the LSV.
- LSV type: a thumbnail representing the splicing event. Each junction has a different color.
- PSI per junction: the expected Percent Selected Index (PSI) per junction. It has two *views* that can be switched between by clicking on them:
 - **Compact view.** Initially, the PSI per junction is represented by the height of a colored box. Each color refers to the expected inclusion of a particular junction. The taller the rectangle, the more included the junction is expected to be. Clicking on the rectangle (zoom in pointer) will open the *expanded view*.
 - **Expanded view.** Distributions of probabilities of PSI per LSV junction (violin

boxplots). The *white* dot represents the expected PSI, whereas the box plot indicates the 10, 25, 50, 75 and 90 percentile of the distribution. Consistently with the Compact view, each color refers to a particular junctions of the LSV.

- LSV links:
 - **GTF**. Link to the GTF file associated with the LSV.
 - **UCSC**. Coordinates link to explore the LSV on UCSC Genome Browser (when available).
 - **Copy LSV**. The new VOILA supports copy-paste functionality into a new web-tool, MAJIQ-SPEL. SPEL allows users to connect the LSV quantification to the various gene isoforms associated with those, automatically design RT-PCR primers for validating the LSV (with optional control over primer design parameters), and map the LSV to protein domains on the genome browser for functional integration and downstream analysis. See more details:
<http://biorxiv.org/content/early/2017/05/09/136077>

Lastly, the panel *LSV filters* allows the user to screen out LSVs with certain properties like having alternative 5-prime splice sites, involve Exon Skipping or contain a certain amount of exons and junctions.

Delta PSI Summary

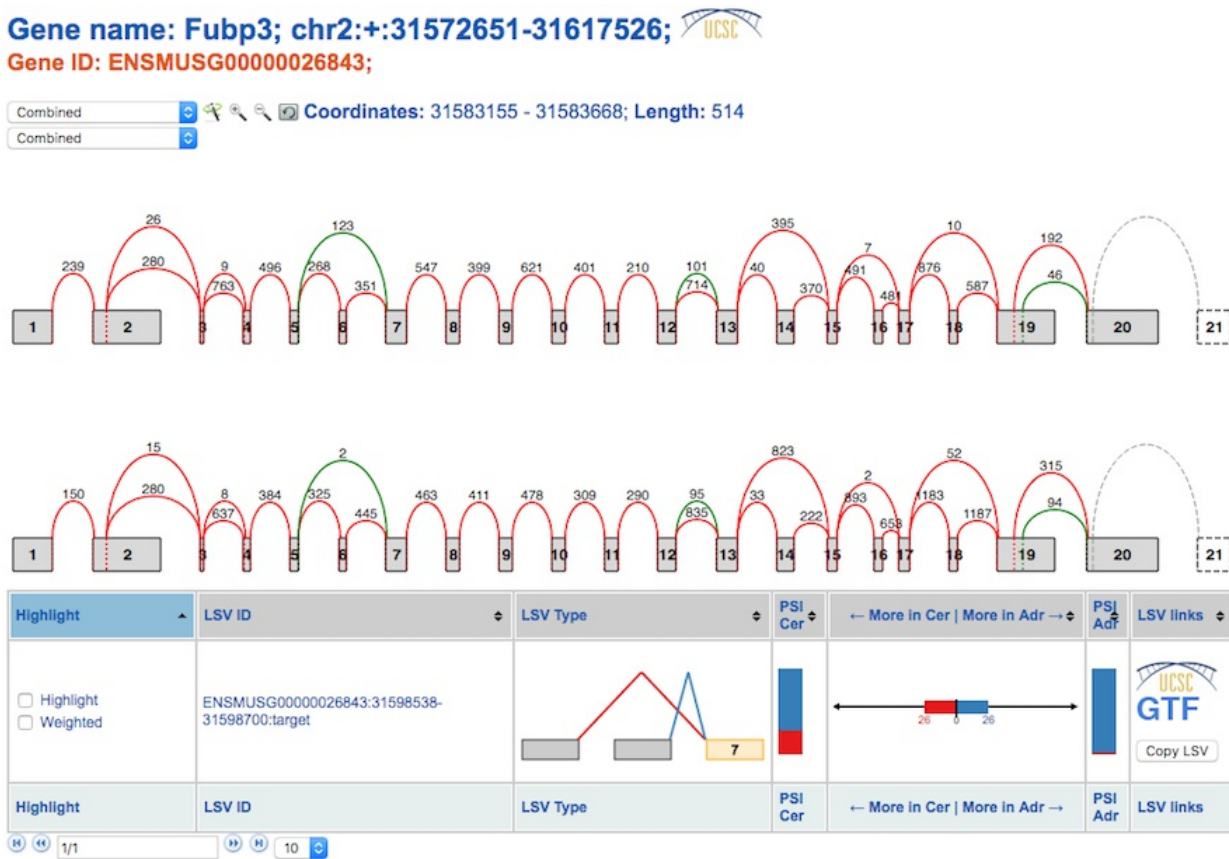
#	Gene	LSV ID	LSV Type	← More in Cer More in Adr →	Coordinates Links
1	Rapgef1	ENSMUSG00000039844:29722224-29722309:target			
2	Clip1	ENSMUSG00000049550:123630225-123631169:target			
3	Dlgap4	ENSMUSG00000061689:156745835-156746196:target			
4	Fubp3	ENSMUSG00000026843:31598538-31598700:target			
5	Cita	ENSMUSG00000028478:44025449-44025560:source			
#	Gene	LSV ID	LSV Type	← More in Cer More in Adr →	Coordinates Links

1/1 All

Tip1: index tables can be sorted by *gene* name, *LSV Type* and/or *most changing junction* clicking on the column header respectively. It is possible to sort the table by multiple columns holding the *shift* key.

Tip2: You can display all LSVs that changed above the specified dPSI threshold (default 0.2) by clicking the dropdown box at the bottom left and select "All" (default: display 10 LSVs). This allows for quick searching (Ctrl+F) for genes of interest.

Delta PSI index file is identical to PSI index except for the *PSI per junction* column that represents the estimated differential inclusion levels as bars leaning towards condition1 or condition2 (see *Compact view* below). Clicking over LSV IDs and genes opens up individual summaries.



The results are broken into 10-genes per page summaries, with LSV quantifications grouped by gene. Unlike in single PSI summaries, there are 2 splice graphs (one per condition) which facilitate quick visual inspection of possible differences. If multiple replicates were used in a condition, the splice graphs for each replicate can be switched between through the drop down box located next to the wand. The LSV information is displayed as follows:

- Highlight:
 - Highlight.** To move the visualization of the LSV to the splice graph, select the “highlight” checkbox in the LSV table. All junctions will be hidden except the ones included in the LSV. The visible junctions will be colored to match the highlighted LSV.
 - Weighted.** To weight the junctions of the highlighted LSV, check the “weighted” checkbox. The junctions for this LSV will now be weighted proportionally to the LSV’s expected PSI value for this experiment.
- LSV ID: a unique identifier for the LSV.
- LSV type: a thumbnail representing the splicing event. Each junction has a different color.
- PSI *condition*[1|2]: individual \$Psi\$ in *condition*[1|2], similar to the \$Psi\$ *per junction*

column in PSI summary, with compact and expanded views.

- More in *condition 1* | More in *condition 2*: the expected delta PSI for all junctions. It has two *views* which can be switched between by clicking on them:
 - **Compact view**. Initially, the observed delta PSI preference for a certain condition (more included) per junction. A bar going towards *condition 1* (left) with a value of 26 means that the junction is 26% more included in condition 1 than in condition 2. Each color refers to a particular junction of the LSV. In the above example, the red junction is 26% more included in cerebellum compared to adrenal gland.
 - **Expanded view**. Distributions of probabilities of Delta PSI per LSV junction (violin boxplots). The *white* dot represents the expected PSI, whereas the box plot indicates the 10, 25, 50, 75 and 90 percentile of the distribution. Consistent with the Compact view, each color refers to a particular junction of the LSV.
- LSV links:
 - **UCSC**. Coordinates link to explore the LSV on UCSC Genome Browser (when available).
 - **Copy LSV**. The new VOILA supports copy-paste functionality into a new web-tool, MAJIQ-SPEL. SPEL allows users to connect the LSV quantification to the various gene isoforms associated with those, automatically design RT-PCR primers for validating the LSV (with optional control over primer design parameters), and map the LSV to protein domains on the genome browser for functional integration and downstream analysis. See more details:
<http://biorxiv.org/content/early/2017/05/09/136077>

Tab-delimited file

VOILA provides a tab-delimited text file to allow users to parse MAJIQ results and further analyze particular LSVs or genes of interest. Most fields are shared between single PSI and delta PSI computations for the expected values and the confidence measures (variance in the case of single PSI and the probability of $|\text{delta psi}| > 0.2$ (or your specified alternative threshold) in delta PSI analysis). The common fields are: Gene name; LSV ID; LSV Type; LSV attributes (A5SS, A3SS, ES, Num. Junctions and Num. Exons); chromosome; strand; LSV coordinates (junctions and exons coordinates); and finally, if additional evidence is required to determine what is the start/end of an LSV, a list with all possible alternative starts and ends is provided.

FAQ

In VOILA gene summaries, what is that *number at*

the beginning of the HTML file?

To achieve a better performance, VOILA creates HTML files of up to 10 genes. Therefore, if MAJIQ detected and quantified LSVs from N genes, there will be $N/10$ pages (always rounded to the upper integer limit).

For example, let say that we executed `voila psi data/Liver1.majiq_psi.pickle --genes-exp1 data/Liver1.splicegraph -o psi_gene_out/` and 182 genes were detected. There will be $182/10=19$ pages (starting with 0): `0_Liver1.majiq_psi_lsv_single_gene.html`, `1_Liver1.majiq_psi_lsv_single_gene.html`, ..., `18_Liver1.majiq_psi_lsv_single_gene.html`.