

# geneviz-java: Design Planning Documentation

## Table of Contents

<b>1. Software Requirements.....</b>	<b>2</b>
1.1 User Needs:.....	2
1.2 Design Inputs .....	3
1.3 Software Requirement Specifications.....	5

## 1. Software Requirements

### 1.1 User Needs:

ID	Description	Downstream
UN-001	As an end user I need a GUI-based tool to visualize gene models based on gene structure data.	DI-001 DI-002
UN-002	As an end user I need a tool that I can pass FASTA and GTF files.	DI-003 DI-004 DI-005 DI-006
UN-003	As an end user I need the gene features to be displayed in a list or table.	DI-007
UN-004	As an end user I need the sequence data from the FASTA file displayed as text.	DI-008
UN-005	As an end user I need the interface to be simple in design and easy to navigate.	DI-009
UN-006	As an end user I need the tool to show me basic statistics on data from both files [FASTA, GFF3].	DI-010
<b>UN-007</b>	As an end user I need the tool to load multiple input files. e.g. multiple FASTA files.	DI-003 DI-010 DI-011 DI-016
UN-008	As an end user I need ability to <emphasize> exons in a displayed sequence.	DI-012
UN-009	As an end user I need ability to display a visual representation of the gene model.	DI-013 DI-014 DI-015

*Note: Entries reflect user requirements only; they do not guarantee deliverables. Bold elements indicate outstanding implementation.*

## 1.2 Design Inputs

ID	Description	Downstream
DI-001	The tool shall provide the user a front end GUI.	SRS-001
DI-002	The tool shall allow the user to visualize a gene model.	SRS-012
DI-003	The tool shall allow the user means to provide input data — eg a file.	SRS-006 SRS-009
DI-004	The tool shall have a parsing module to parse a FASTA file.  Notes: Consider BioJava.	SRS-005
DI-005	The tool shall have a parsing module to parse a GTF file.  Notes: GFF3 support should be prioritized, the formats are similar enough so consider implementing both GTF and GFF3. Main difference seems to be the last column.  Consider BioJava.	SRS-003 SRS-004
DI-006	The tool shall have a system to store data.  Notes: This could be dedicated class objects or an SQL database. If SQL database, this is likely to be an ephemeral database that exists while the user has the application open. Avoid mySQL/postgres as these require hosting on a server. SQLite database is an option to consider.	SRS-004 SRS-005 SRS-006
DI-007	The tool shall allow the user means visualize the gene features (from the GFF3) in a tabular visualization.	SRS-002
DI-008	The tool shall allow the user means visualize the FASTA file as display text.	SRS-007
DI-009	The tool shall allow consider multiple tabs to separate out visual features.  Notes: Consider separating out gene-feature table from sequence annotations. Consider where the exon sequence and stats displays fit into the GUI design	SRS-002 SRS-007
DI-010	The tool shall allow the user means visualize basic statistics from the FASTA file.	SRS-008
DI-011	The tool shall allow the user means visualize basic statistics from the GFF3 file.	SRS-010
DI-012	The tool shall allow the FASTA display text to be intractable — either by bulk highlighting exons on import (MVP) or by allowing the user the ability to select a gene and emphasize its exon.  Note:	SRS-007 SRS-011

	Could be highlighted, bold, underlined — anything that makes sense (or is easy!?)	
<b>DI-013</b>	The tool shall allow the user a method to select a specific gene.	SRS-011
DI-014	The tool shall allow the user an interface with a visual representation of a gene.	SRS-012
<b>DI-015</b>	Context menu for export (PNG/SVG) of gene model.	
<b>DI-016</b>	The tool shall be able to handle zipped or unzipped files.	SRS-003 SRS-004

*Note: Entries reflect design inputs only; they do not guarantee deliverables. Bold elements indicate outstanding implementation.*

### 1.3 Software Requirement Specifications

ID	Description
SRS-001	The tool shall be written in Java since this is an easy to work with framework for GUI applications.
SRS-005	<p>The tool shall be able to parse FASTA formats.</p> <p>Notes:</p> <ul style="list-style-type: none"> <li>— date accessed: 15-Nov-25: <a href="https://www.ncbi.nlm.nih.gov/genbank/fastaformat/">https://www.ncbi.nlm.nih.gov/genbank/fastaformat/</a></li> <li>— date accessed: 15-Nov-25: <a href="https://www.ncbi.nlm.nih.gov/genbank/mods_fastadefline/">https://www.ncbi.nlm.nih.gov/genbank/mods_fastadefline/</a></li> </ul> <p>In FASTA format the line before the nucleotide sequence, called the FASTA definition line, must begin with a carat ("&gt;"), followed by a unique SeqID (sequence identifier).</p> <p>The SeqID must be unique for each nucleotide sequence and should not contain any spaces.</p> <p>Please limit the SeqID to 25 characters or less.</p> <p>The SeqID can only include letters, digits, hyphens (-), underscores (_), periods (.), colons (:), asterisks (*), and number signs (#).</p> <p>Information about the source organism from which the sequence was obtained follows the SeqID and must be in the format [modifier=text]. Do not put spaces around the "=".</p> <p>At minimum, the scientific name of the organism should be included.</p> <p>Addition modifiers exist.</p> <p>e.g: "SeqABCD [organism=Mus musculus] [strain=C57BL/6] ....."</p> <p>The FASTA definition line must not contain any hard returns. All information must be on a single line of text.</p> <p>Note: The format is not fixed, so don't over engineer this part of the code base.</p> <p>Examples from <a href="https://www.bioinformatics.nl/tools/crab_fasta.html">https://www.bioinformatics.nl/tools/crab_fasta.html</a> have different characters in the header, some FASTA files seem to have " " character.</p>

SRS-003	<p>The tool shall be able to parse GTF formats.</p> <p>Notes:</p> <ul style="list-style-type: none"> <li>— date accessed: 15-Nov-25</li> <li>— date accessed: 15-Nov-25</li> </ul> <p>The GTF (General Transfer Format) is identical to GFF version 2. GFF2 can only represent 2 level feature hierarchies, while GFF3 can support arbitrary levels. GFF2 also does not require that column 3, the feature type, be part of the sequence ontology. The GFF format is a flat tab-delimited file, each line of which corresponds to an annotation, or feature. Each line has nine columns and looks like this:</p> <pre>##### Chr1 curated CDS 365647 365963 . + 1 Transcript "R119.7" #####</pre> <p>reference sequence: This is the ID of the sequence that is used to establish the coordinate system of the annotation. In the example above, the reference sequence is “Chr1”.</p> <p>source: The source of the annotation. This field describes how the annotation was derived. In the example above, the source is “curated” to indicate that the feature is the result of human curation. The names and versions of software programs are often used for the source field, as in “tRNAscan-SE/1.2”.</p> <p>method: The annotation method, also known as type. This field describes the type of the annotation, such as “CDS”. Together the method and source describe the annotation type.</p> <p>start: The start of the annotation relative to the reference sequence.</p> <p>stop: The stop of the annotation relative to the reference sequence. Start is always less than or equal to stop.</p> <p>score: For annotations that are associated with a numeric score (for example, a sequence similarity), this field describes the score. The score units are completely unspecified, but for sequence similarities, it is typically percent identity. Annotations that do not have a score can use “.”</p> <p>strand: For those annotations which are strand-specific, this field is the strand on which the annotation resides. It is “+” for the forward strand, “-“ for the reverse strand, or “.” for annotations that are not stranded.</p> <p>phase: For annotations that are linked to proteins, this field describes the phase of the annotation on the codons. It is a number from 0 to 2, or “.” for features that have no phase.</p> <p>group: GFF provides a simple way of generating annotation hierarchies (“is composed of” relationships) by providing a group field. The group field contains the class and ID of an annotation which is the logical parent of the current one. In the example given above, the group is the Transcript named “R119.7”. The group field is also used to store information about the target of sequence similarity hits, and miscellaneous notes.</p>
---------	--

SRS-004	<p>The tool shall be able to parse GFF3 formats.</p> <p>Notes:</p> <ul style="list-style-type: none"> <li>— date accessed: 15-Nov-25: <a href="https://gmod.org/wiki/GFF2">https://gmod.org/wiki/GFF2</a></li> <li>— date accessed: 15-Nov-25: <a href="https://www.ensembl.org/info/website/upload/gff.html">https://www.ensembl.org/info/website/upload/gff.html</a></li> </ul> <p>The GFF (General Feature Format) format consists of one line per feature, each containing 9 columns of data, plus optional track definition lines. Fields must be tab-separated.</p> <p>All but the final field in each feature line must contain a value; "empty" columns should be denoted with a '.'.</p> <p>seqname: name of the chromosome or scaffold; chromosome names can be given with or without the 'chr' prefix. Important note: the seqname must be one used within Ensembl, i.e. a standard chromosome name or an Ensembl identifier such as a scaffold ID, without any additional content such as species or assembly. See the example GFF output below.</p> <p>source: name of the program that generated this feature, or the data source (database or project name).</p> <p>feature: feature type name, e.g. Gene, Variation, Similarity</p> <p>start: Start position* of the feature, with sequence numbering starting at 1.</p> <p>end: End position* of the feature, with sequence numbering starting at 1.</p> <p>score: A floating point value.</p> <p>strand: defined as + (forward) or - (reverse).</p> <p>frame: One of '0', '1' or '2'. '0' indicates that the first base of the feature is the first base of a codon, '1' that the second base is the first base of a codon, and so on..</p> <p>attribute: A semicolon-separated list of tag-value pairs, providing additional information about each feature.</p>
SRS-006	<p>User Input module to get file from user. Following Friday's practical I could implement a menu bar with option to load a file — seems generally conventional. I could also have a button somewhere on the GUI to say load file. The requirement is to have the user be able to add their file.</p> <p>The SRS should also restrict the user to opening only permitted file types — eg the user can select a GFF file but cannot select a PDF file.</p> <p>Notes:</p> <p>Check out jMenuBar.</p> <p>Check out jMenu.</p> <p>Check out jMenuItem.</p> <p>Check out FileDialog   JFileChooser</p>
SRS-009	<p>User Input from multiple FASTA / GFF3 files.</p> <p>Note:</p> <p>How shall the tool manage this. Does every FASTA file require a GFF3 file? What happens if the user only provides one.</p>

SRS-002	<p>Visualize tables in GUI.</p> <p>Notes: Check out jTable Component.</p>
SRS-007	<p>Visualize sequences in GUI.</p> <p>Notes: Check out JTextArea Component.</p>
SRS-011	<p>The tool shall have a way to determine what is an exon (from GFF3 file) and what is the start and stop point of this feature (GFF3 file).</p> <p>Notes: Genes can have more than one exon, so you likely need a method that can return the start:stop points for all the exons in a gene.</p> <p>Again, database operations could help.</p> <pre>" SELECT exon_id, start, stop, FROM gff_table WHERE gene_id IN ("x1", "x2") "</pre> <p>Then select these locations from the sequence and highlight the corresponding values.</p> <p>In addition, I think this approach makes it easy to have a drop down box and select individual genes. Consider an option to have all genes selected, or a selection of genes.</p> <p>Simple way uses JTextPane and HTML tags — explore other options. See HighlightPainter</p>
SRS-008	<p>The tool shall calculate:</p> <p>Length of sequence from FASTA file.</p> <p>Average sequence length (multi-file).</p> <p>GC content from FASTA file. [BioJava: DNASequence.getGCCCount()]:</p> <p>—date accessed 15-Nov-25: <a href="https://biojava.org/docs/api7.1.4/org/biojava/nbio/core/sequence/DNASequence.html">https://biojava.org/docs/api7.1.4/org/biojava/nbio/core/sequence/DNASequence.html</a></p> <p>Note: Database operation could make these summaries quick and easy.</p>

SRS-010	<p>The tool shall calculate:</p> <ul style="list-style-type: none"> <li>Average number of exons per gene from GFF3 file.</li> <li>Longest and shortest gene models from GFF3 file.</li> <li>Average gene length from GFF3 file.</li> </ul> <p>Note: Database operation could make these summaries quick and easy</p>
SRS-012	<p>The tool shall render a gene image and superimpose exons (from GFF3).</p> <p>— date accessed: 15-Nov-25:  <a href="https://biojava.org/docs/api7.1.4/org/biojava/nbio/genome/parsers/gff/FeatureList.html">https://biojava.org/docs/api7.1.4/org/biojava/nbio/genome/parsers/gff/FeatureList.html</a></p> <p>Notes:</p> <p>Consider drawing shapes on a JPanel.</p> <p>Reuse methods from SRS-011 to get start and stop points of exons from a gene.</p> <p>Implement a Utils method to draw a generic rectangle of size y-axis size fixed to gene shape and x-axis determined by the proportionate start/stop from the gene to the length of the gene-shape rectangle.</p> <p>Consider external packages — they look a bit heavy!</p> <p>Add tooltip OnHover to get information on the exon.</p> <p>How does it handle overlapping exons? [omitOverlapping method on the FeatureList]</p> <p>How does it handle multiple FASTA files?</p>

*Note: Entries reflect design inputs only; they do not guarantee deliverables. Bold elements indicate outstanding implementation.*