# *Sequence Updater* User Guide
Wilson Leung

## Table of Contents

# Introduction

The genome assemblies used by the GEP scientific projects have not undergone manual sequence improvement and they might contain consensus errors (i.e. substitutions, insertions, and deletions) that impact the annotation of protein-coding genes. For example, the Baylor College of Medicine Human Genome Sequencing Center (BCM-HGSC) has previously sequenced eight *Drosophila* species using the 454 and Illumina sequencing platforms. The 454 sequencing technology exhibits low accuracy in long (> 6bp) mononucleotide runs. Similarly, sequencing reads produced by the Pacific Biosciences (PacBio) and Nanopore sequencing platforms have lower read accuracy (~90%) and higher rates of insertions and deletions compared to Sanger sequencing data.

The GEP developed the *Sequence Updater* tool to provide a standardized way to document consensus errors in the project sequence using the [Variant Call Format](#) (VCF). The VCF files generated by the *Sequence Updater* can be used in conjunction with the [*Gene Model Checker*](#) to verify gene models that contain consensus sequence errors. Multiple VCF files can be combined into a single project VCF file using the [*Annotation Files Merger*](#) in preparation for project submission.

This user guide provides an overview of the strategies that could be used to identify consensus errors in the project sequence, and it describes how to document these consensus errors using the *Sequence Updater*. The [*Gene Model Checker* User Guide](#) describes how the VCF file produced by the *Sequence Updater* can be used to verify a gene model with errors in the consensus sequence.

# Acknowledgements

The *Sequence Updater* is developed by Wilson Leung at Washington University in St. Louis for the Genomics Education Partnership (GEP).

# Questions about the *Sequence Updater*

Please contact Wilson ([wleung@wustl.edu](mailto:wleung@wustl.edu)) if you have any questions or encounter any problems with the *Sequence Updater*.

# Availability

The *Sequence Updater* is available under the "Resources & Tools" section of the [F Element project page](#) and the [Pathways project page](#) on the GEP website.

# Overview of the *Sequence Updater*

We will use a consensus error within the *tgo* gene in the *Drosophila sechellia* May 2011 (Broad dsec_caf1/DsecCAF1) assembly to illustrate the key functionalities of the *Sequence Updater*. This document will discuss how you can use the *D. sechellia* RNA-Seq data to identify potential discrepancies between the RNA-Seq reads and the project consensus sequence, how to document these errors using the *Sequence Updater*.

## Using *tblastn* searches to identify potential consensus sequence errors

When nucleotide insertions or deletions (indels) that are not multiples of three are introduced into the coding regions of a project sequence, the *blastx* or *tblastn* alignment of the coding exons (CDS) against the project sequence will often be split into multiple alignment blocks. This occurs because the indel will introduce a frame-shift within the CDS, which results in the nucleotides before the indel being translated in one reading frame, and the nucleotides after the indel being translated in a different reading frame. Hence the split alignment blocks are typically located adjacent to each other in the project sequence but they are in different reading frames.

We will use the annotation of the *tgo* gene in *D. sechellia* to illustrate the issue with split alignments. The putative ortholog of the *tgo* gene in *D. sechellia* is located at approximately 17,034,485-17,036,408 on scaffold super_0 (GenBank accession number CH480815.1). According to the *Gene Record Finder*, the *tgo* gene has a single CDS (1_13223_0) in *D. melanogaster* (Figure 1).



**Figure 1. The *Gene Record Finder* record shows that there are two isoforms of *tgo* in *D. melanogaster* (A and B). Both isoforms have only one CDS (1_13223_0).**

A *tblastn* alignment of CDS 1_13223_0 from the *D. melanogaster tgo* gene against the 17,000,000-17,100,000 region of the *D. sechellia* scaffold super_0 shows two highly significant alignment blocks. The first alignment block (with an E-value of 0.0) covers the first 565 residues of CDS 1_13223_0, which spans from 17,034,485-17,036,176 on the *D. sechellia* scaffold super_0 in frame +2 (Figure 2).



**Figure 2. The first alignment block from the *tblastn* search of the *D. melanogaster tgo* CDS 1_13223_0 (query) against the *D. sechellia* scaffold super_0 (CH480815.1; subject) placed the beginning of the CDS at 17,034,485-17,036,176 in frame +2.**

The second alignment block covers the residues 550-643 of CDS 1_13223_0, which spans from 17,036,130-17,036,411 on the *D. sechellia* scaffold super_0 in frame +3 (Figure 3).



**Figure 3. The second alignment block from the *tblastn* search of the *D. melanogaster tgo* CDS 1_13223_0 (query) against the *D. sechellia* scaffold super_0 (CH480815.1; subject) placed the end of the CDS at 17,036,130-17,036,411 in frame +3.**

These two alignment blocks cover the entire length (643aa) of CDS 1_13223_0. However, 16 residues (from 550-565) of CDS 1_13223_0 appear in both alignment blocks. In addition, since the first alignment block ends at 17,036,176 and the second alignment block begins at 17,036,130 on the *D. sechellia* scaffold super_0, there is a 47bp overlap between the two alignment blocks. The reading frame has also changed from +2 in the first alignment block to +3 in the second alignment block.

Collectively, the available evidence suggests that CDS 1_13223_0 of the *tgo* gene in *D. sechellia* has either acquired a novel intron or there is an error in the super_0 consensus sequence which leads to a frame shift in the region around 17,036,130 (i.e. the start of the second alignment block).

## Using RNA-Seq data to identify potential consensus errors

To determine if the frame shift was caused by a genuine difference in the *D. sechellia* super_0 genomic sequence or an error in the consensus sequence, we will examine the RNA-Seq reads that have been mapped to the region surrounding position 17,036,130 in scaffold super_0. The RNA-Seq reads were generated by the Illumina sequencing platform, which has an average base accuracy greater than 99%. If a genomic region shows multiple high quality discrepancies between the RNA-Seq reads that aligned to the region and the consensus sequence, then either the RNA-Seq reads have been placed in the wrong part of the assembly or the consensus sequence is incorrect.

Based on the start of the second alignment block from the *tblastn* search, the potential issue with the consensus sequence likely occurs at approximately 17,036,130. Open a new web browser window and navigate to the *GEP UCSC Genome Browser*. Click on the "Genome Browser" link on the left sidebar. Select "*D. sechellia*" under the "Browse/Select Species" section, select "May 2011 (Broad dsec_caf1/DsecCAF1)" under the "D. sechellia Assembly" field, and then enter "super_0:17036130" into the "Position/Search Term" field (Figure 4). Click on the "GO" button.



**Figure 4. Navigate to position super_0:17036130 in the *D. sechellia* May 2011 (Broad dsec_caf1/DsecCAF1) assembly using the *GEP UCSC Genome Browser*.**

To get a broader view of the genomic region surrounding this position, zoom out 100x. Click on the "hide all" button beneath the Genome Browser image to hide all the evidence tracks. Scroll down to the track configuration section and then change the display settings for the following evidence tracks:

Under "Mapping and Sequencing Tracks":
- Base Position: **dense**
- Gap: **dense**
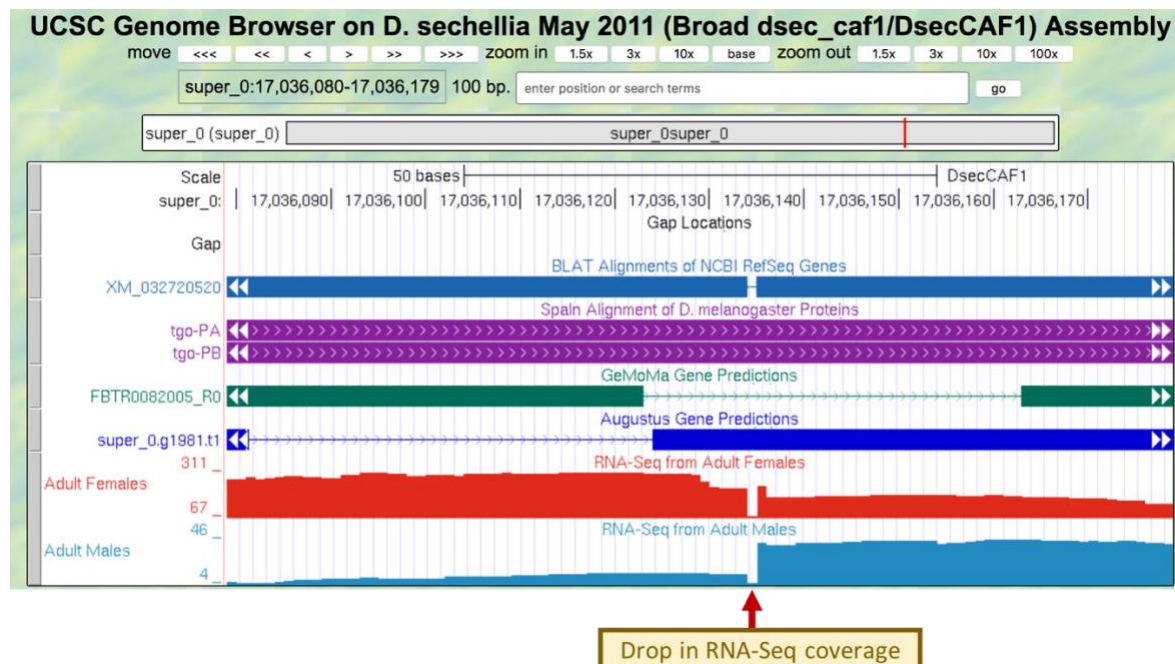
Under "Genes and Gene Prediction Tracks":
- RefSeq Genes: **pack**
- D. mel Proteins: **pack**
- *GeMoMa* Genes: **pack**
- *Augustus*: **pack**

Under "RNA-Seq Tracks":
- RNA-Seq Coverage: **full**

Click on one of the "refresh" buttons to update the Genome Browser display.

The *GeMoMa* prediction FBTR0082005_R0 and the *Augustus* gene prediction super_0.g1981.t1 both introduced an intron in this region. The *BLAT* alignment of XM_032720520 in the NCBI "RefSeq Genes" track shows a single base gap. The "RNA-Seq Coverage" tracks show a substantial drop in read coverage in the adult females and adult males RNA-Seq samples at the same position as the gap in the "RefSeq Genes" track (Figure 5).



**Figure 5. Introns predicted by multiple gene predictors and the decrease in RNA-Seq read coverage demarcate the position within the *D. sechellia* super_0 scaffold which contains the potential consensus error.**

To inspect the position which shows a substantial decrease in RNA-Seq coverage more closely, enter "super_0:17,036,125-17,036,145" into the "enter position or search terms" text box and then click on the "go" button. We find that G nucleotide at 17,036,135 corresponds to the position with the substantial drop in RNA-Seq read coverage (Figure 6). This drop in RNA-Seq read coverage typically indicates a discrepancy between the RNA-Seq reads aligned to this region and the consensus sequence.
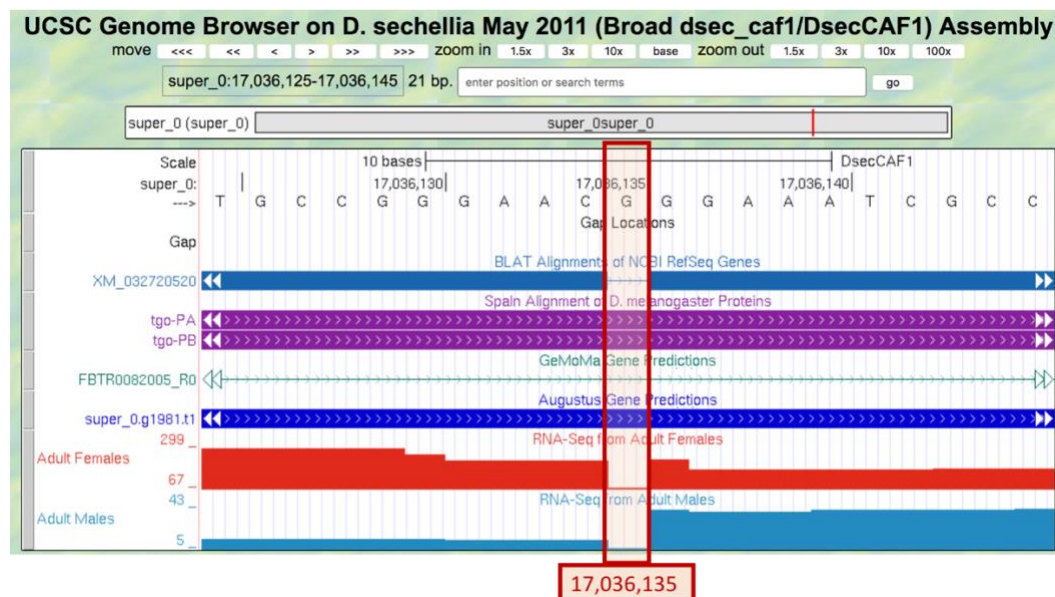


**Figure 6. The RNA-Seq Read coverage tracks indicate that the discrepant position between the RNA-Seq reads and the consensus sequence is located at 17,036,135 in scaffold super_0.**

## Use *SRA-BLAST* to gather additional evidence for the consensus error

To assess the potential consensus sequence error at position 17,036,135 in scaffold super_0, we will compare this genomic region against the *D. sechellia* RNA-Seq reads using NCBI *BLAST*.

In order to perform this analysis, we will need to first determine the source of the RNA-Seq data used to generate the RNA-Seq read coverage track. Scroll down to the "RNA-Seq Tracks" section and then click on the "RNA-Seq Coverage" link. The "Description" section indicates that the adult females RNA-Seq data was obtained from the BioProject with the accession number PRJNA205470, and the adult males RNA-Seq data was obtained from the BioProject PRJNA414017.
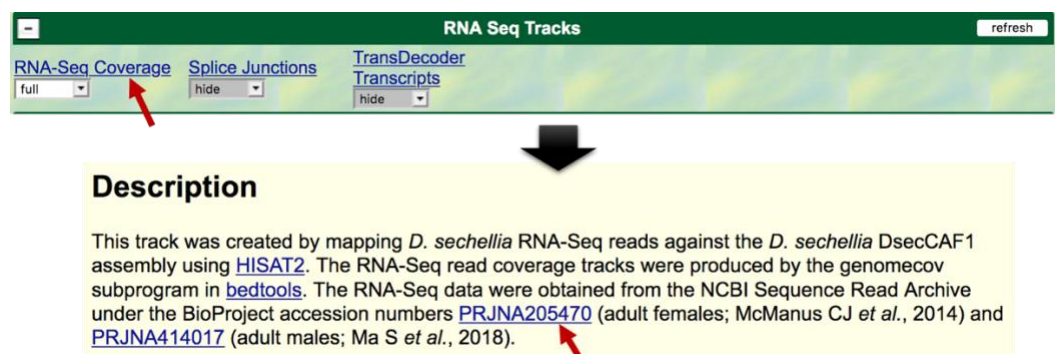


**Figure 7. The "Description" section of the details page for the "RNA-Seq Coverage" track indicates that the *D. sechellia* adult females RNA-Seq data was obtained from the NCBI BioProject PRJNA205470.**

Click on the "PRJNA205470" link in the "Description" section to navigate to the BioProject record for the *D. sechellia* RNA-Seq sample. Scroll down to the "SRA Experiment" row under the "Project Data" section, and then click on the "9" link (Figure 8).



**Figure 8. Click on the link under the "SRA Experiments" row in the "Project Data" table to access the sequencing data associated with the BioProject PRJNA205470.**

The NCBI Sequence Read Archive (SRA) search results page shows that the *D. sechellia* adult females RNA-Seq data has the SRA Experiment (SRX) accession number SRX287399 (Figure 9).



**Figure 9. Use the SRA search results page to determine the accession number for the *D. sechellia* adult females RNA-Seq experiment (i.e. SRX287399).**

In order to compare the *D. sechellia* RNA-Seq reads from this experiment against the consensus sequence from the super_0 scaffold, open a new web browser window and navigate to the NCBI SRA home page. Click on the "SRA-BLAST" link under the "Tools and Software" section (Figure 10).
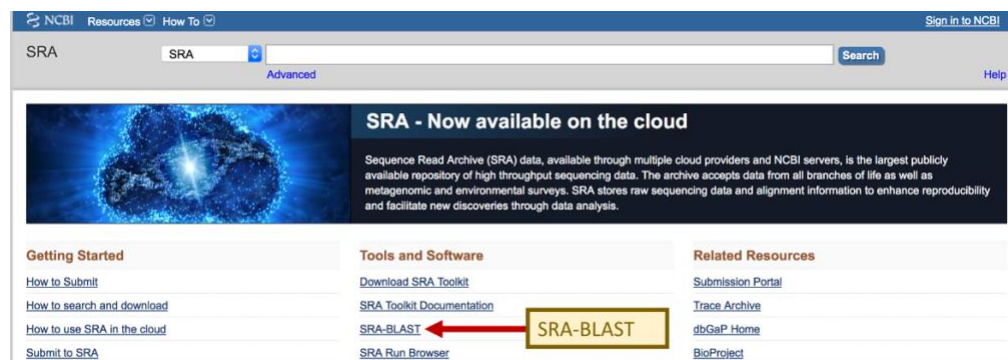


**Figure 10. Access the *SRA-BLAST* service from the NCBI SRA home page.**

Enter the GenBank accession number for the *D. sechellia* scaffold super_0 (i.e. CH480815.1) into the "Enter Query Sequence" text box. Since the potential discrepant position is located at 17,036,135, we will limit the search region to 17,036,100-17,036,200 of the scaffold. Under the "Query subrange" section, enter "17036100" into the "From" field and "17036200" into the "To" field.

Under the "SRA Experiment set (SRX)" field, enter the accession number "SRX287399" to search the reads from the *D. sechellia* adult females RNA-Seq experiment. As you enter the accession number, the field will show a suggestion with the taxonomy and SRA Run accession number for the experiment [i.e. SRX287399 (taxid:7238; run:SRR869601)] (Figure 11). Click on the "BLAST" button to run the search.



**Figure 11. Configure *SRA-BLAST* to search the 17,036,100-17,036,200 region of the *D. sechellia* scaffold super_0 (CH480815.1; query) against the sequencing reads in the SRA Experiment SRX287399 (database).**

Once the *SRA-BLAST* search is complete, click on the "Alignments" tab to see the alignments between the *D. sechellia* RNA-Seq reads and the consensus sequence. These alignments consistently show an extra G at 17,036,135 of the *D. sechellia* scaffold super_0 compared to the RNA-Seq reads (Figure 12).
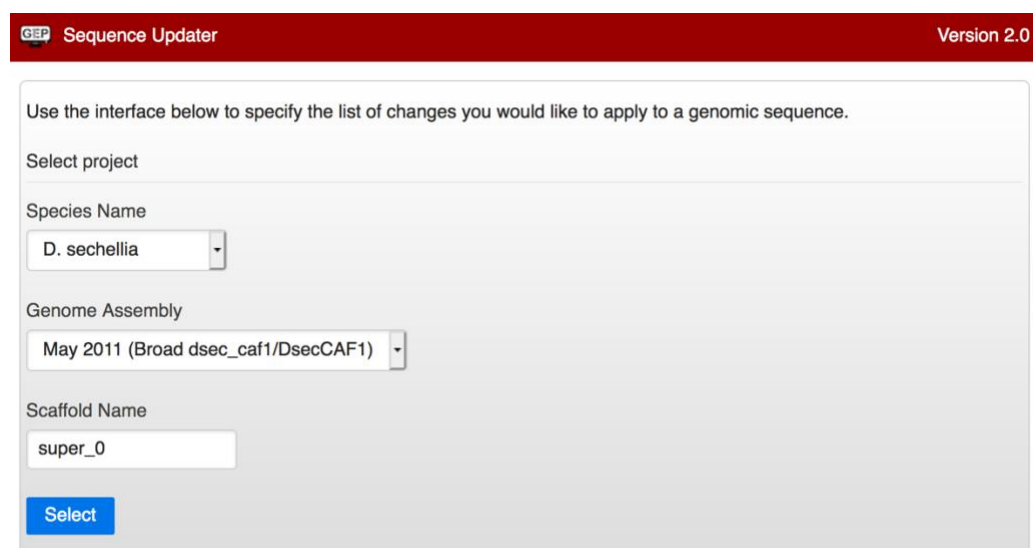


Figure 12. The *SRA-BLAST* alignment of the *D. sechellia* scaffold super_0 (CH480815.1; query) against the RNA-Seq read SRR869601.2846847.2 (top) and SRR869601.3144243.1 (bottom). Both alignments show an extra G at 17,036,135 in the *D. sechellia* scaffold super_0.

Based on the analysis of the *tblastn* search result for CDS 1_13223_0, the gene predictions and RNA-Seq read coverage tracks on the *GEP UCSC Genome Browser*, and the *SRA-BLAST* search result, the G at 17,036,135 should be removed from the *D. sechellia* super_0 sequence.

# Document the consensus sequence error with the *Sequence Updater*

Now that we have identified a consensus error at position 17,036,135 of scaffold super_0, we will use the *Sequence Updater* to document this error. Open a new web browser tab and navigate to the [Pathways project page](#) on the GEP website. Click on the "*Sequence Updater*" link under the "Resources & Tools" section.

Select "D. sechellia" under the "Species Name" field, select "May 2011 (Broad dsec_caf1/DsecCAF1) under the "Genome Assembly" field, and then enter "super_0" into the "Scaffold Name" field (Figure 13). Click on the "Select" button.



**Figure 13. Specify the species, genome assembly, and the scaffold that will be modified by the *Sequence Updater*.**

The *Sequence Updater* uses the [Variant Call Format](#) (VCF) to describe the changes to the original project sequence. In order to document a change in the sequence, we must specify the start coordinate of the change (relative to the original sequence), the original sequence, and the new sequence. The start coordinate corresponds to the first position that changed between the original and the new sequence. In the case of base substitutions, this will correspond to the position where the first base substitution occurs. However, in the case of base insertions or deletions, the start position will correspond to the base just before the indels.

Since we would like to remove the G nucleotide at 17,036,135 from the consensus sequence, we will use 17,036,134 as the start coordinate, and then omit the G at 17,036,135 as part of the change.

Enter "17036134" into the "Start Position" field. A tooltip will appear which shows the nucleotide at this position (C) and the nucleotides surrounding this position. To remove the G at 17,036,135, enter "CG" into the "Original Sequence" field and "C" into the "New Sequence" field (Figure 14). Click on the "Add" button to add the proposed modification to the "List of sequence changes" section (Figure 15).

**Figure 14. Use the *Sequence Updater* interface to describe the changes that should be applied to the original sequence. When you specify a "Start Position", a tooltip will appear on the right which shows the nucleotide at that position (e.g., the C in red) and the 15 nucleotides before and after that position (nucleotides in blue).**
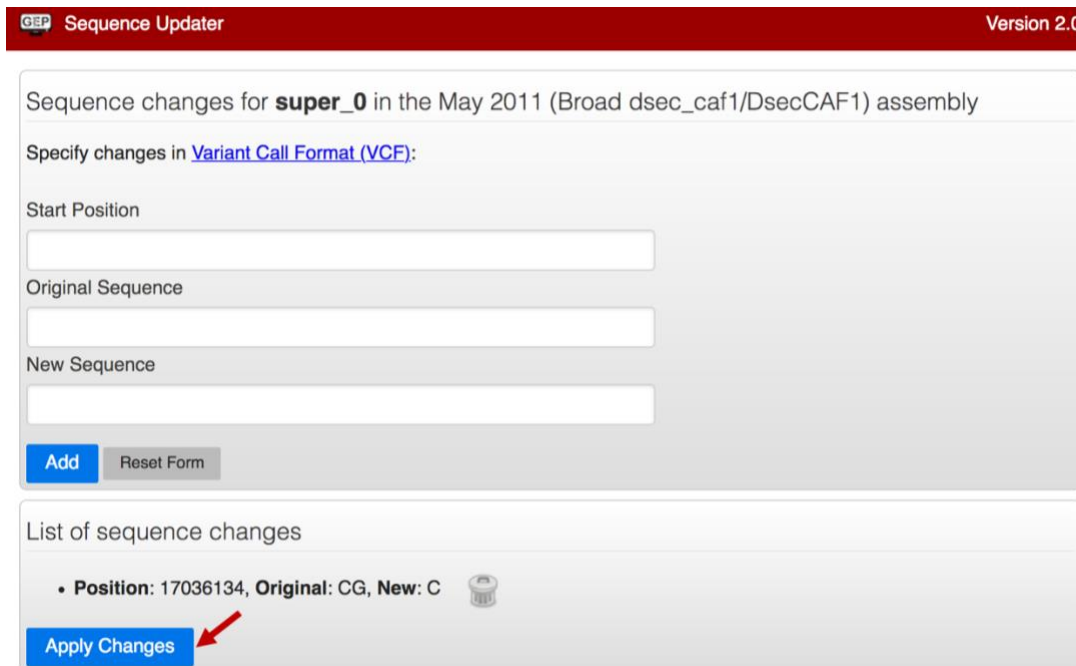


**Figure 15. The changes that will be applied to the original sequence are listed under the "List of sequence changes" section. Click on the trash can icon to delete the modification from the list of sequence changes.**

If the region contains multiple consensus errors, you can use the same procedure to add the additional changes to the "List of sequence changes" section.

> **Note**: The "Start Position" and "Original Sequence" are based on the sequence in the *GEP UCSC Genome Browser*. In cases where there are multiple consensus errors, **all the modifications should be relative to the original sequence** in the *GEP UCSC Genome Browser*. The *Sequence Updater* will automatically transforms the provided start positions as it iteratively applies the modifications to the original sequence.

Click on the "Apply Changes" button to create the VCF file which describes the list of changes to the original sequence (Figure 16).

**Figure 16. After specifying the sequence changes using the *Sequence Updater* interface, click on the "Apply Changes" button to generate the VCF and the revised sequence file.**

Once the analysis is complete, A "Download Results" panel will appear with a link to the VCF file (Figure 17). Right click (control-click on macOS) on the "VCF file" link and then select "Save Link As…" or "Download Linked File As…" to save the VCF file on your computer.

For projects with errors in the consensus sequence, you will need to submit the VCF file in conjunction with the annotation report, the GFF file, the transcript sequence file, and the peptide sequence file. For the purpose of this tutorial, we will name the VCF file "**dsec_tgo.vcf.txt**".



**Figure 17. Right-click (control-click on macOS) on the "VCF file" link and select "Save Link As…" or "Download Linked File As…" to save the VCF file generated by the *Sequence Updater*.**

## Use the VCF file to verify gene models with errors in the consensus sequence

After documenting the error at 17,036,135 of the *D. sechellia* scaffold super_0, we can apply the changes in the VCF file to the project sequence when we verify a gene model using the *Gene Model Checker*. The *Gene Model Checker* User Guide includes a walkthrough which illustrates how to use the VCF file with the *Gene Model Checker* to verify the gene model for *D. sechellia* tgo-PA (Figure 18).



**Figure 18. The VCF file produced by the *Sequence Updater* can be provided to the *Gene Model Checker* through the "File with Changes to the Consensus Sequence" field.**

> **Note**: When you use the *Gene Model Checker* to verify a gene model with consensus errors, all the exon coordinates should be **relative to the original project sequence**.

# Conclusion

This user guide demonstrates how you can use the *Sequence Updater* to document potential consensus errors in the project sequence. It shows how *tblastn* CDS searches and RNA-Seq data on the *GEP UCSC Genome Browser* can be used to identify potential consensus errors. It then uses the NCBI SRA and the *SRA-BLAST* service to compare the RNA-Seq reads against the project sequence to confirm a consensus error. Finally, the user guide provides an overview of the *Sequence Updater* interface and how to use it to create a VCF file that is suitable for use with the *Gene Model Checker*.