

Gene Annotation: Finding and characterizing genes in genomes

Malcolm Arnott '22, Blake Tellinghusen '23, Vy Lam '22, Jack Allen '21, Sam Alper '19, Juliana Choza '20, and Sami Zimmerman '19, *Department of Biology¹, BCMB², Psychology³, Lewis & Clark College.*

Introduction

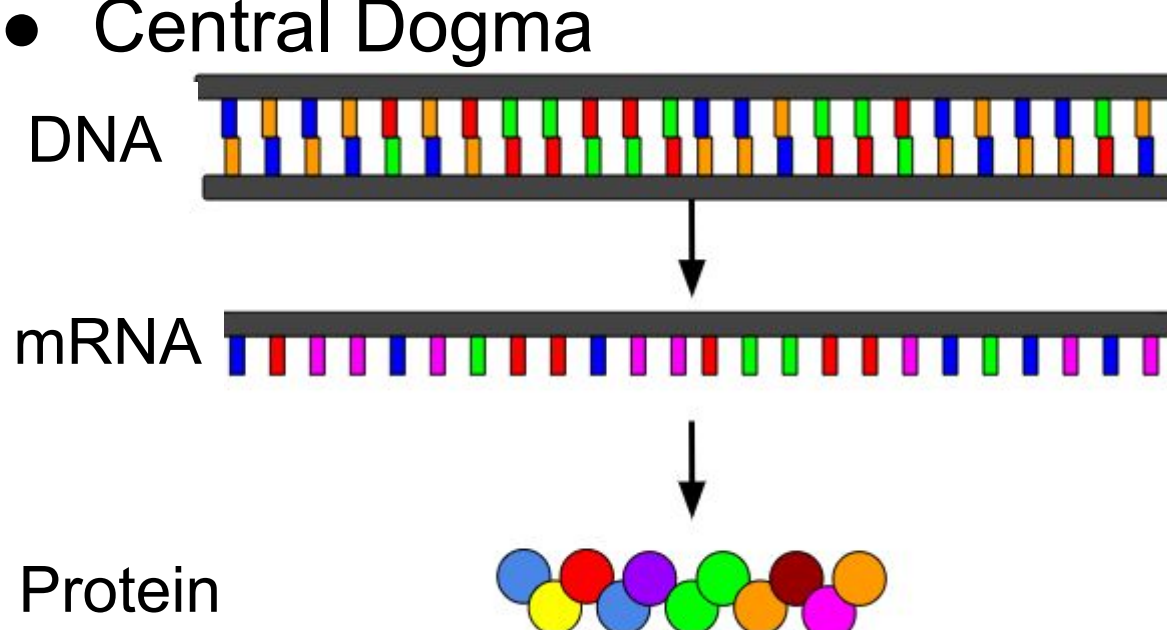
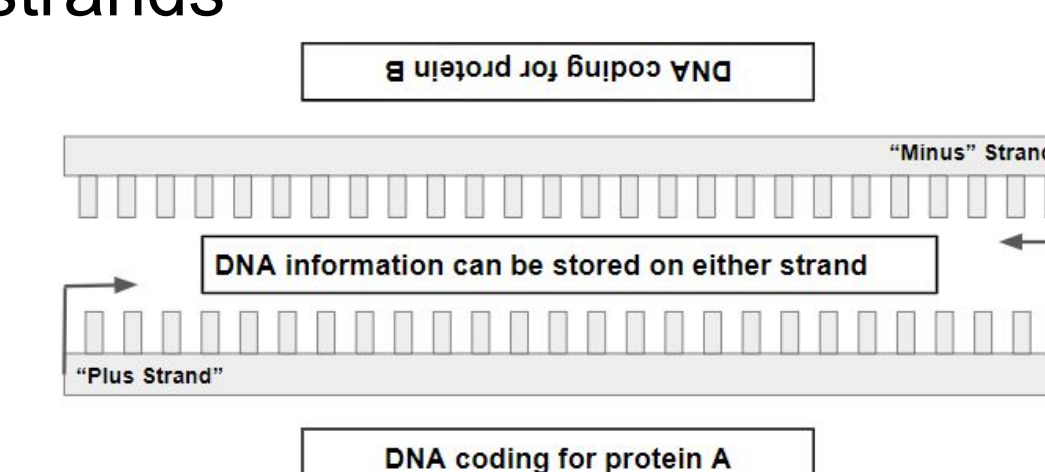
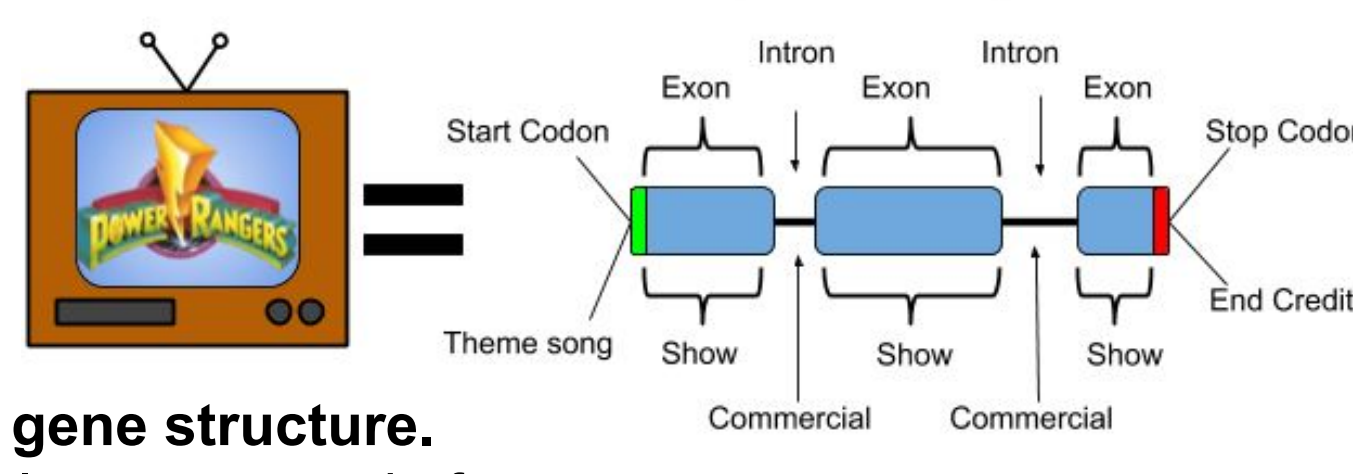
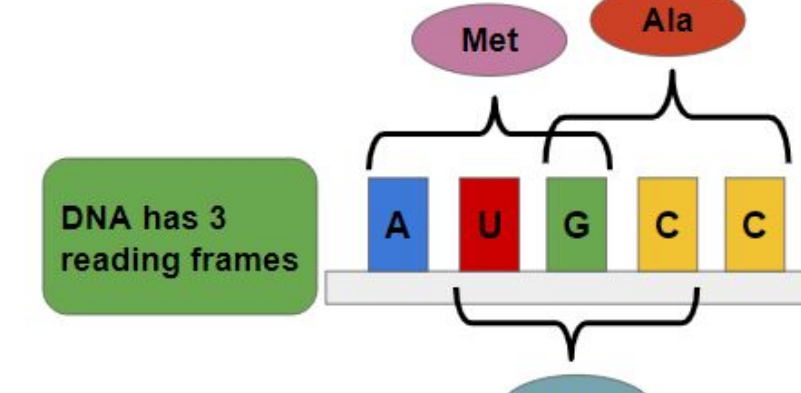
- Central Dogma
 
- DNA strands
 
- Gene structure
 - Both have:
 - Beginning, middle, and end
 - Irrelevant components.
 - Introns are like commercials. They can be removed without altering the "story"
- Reading Frames & codons
 

Figure 1. Schematic of the central dogma in biology and gene structure. Genomic information in an organism is stored in DNA molecules, composed of a sequence of nucleotides. Genes have a defined structure.

Terminology

- Open Reading Frame (ORF):** Sequence of DNA without a stop codon
- Exon:** Segment of a gene expressed as a protein
- Intron:** Gene segment removed during splicing
- Splicing:** Removal of introns
- Ortholog:** Homologous (similar) gene sequences

Important Sequences

- Start codon: ATG
- Stop codons: TAA, TGA, TAG
- Splice sites: DNA sequences prompting intron removal
 - Donor: GT
 - Acceptor: AG

Methods

Gene annotation process:

- Select the gene prediction
- Identify the ortholog in *D. melanogaster*
- Review structure of ortholog
- Determine approximate locations of exons using BLAST/RNA-Seq
- Identify splice sites (GT/AG) and determine phases
- Verify splice sites with tophat where possible
- Verify annotation with Gene Model Checker

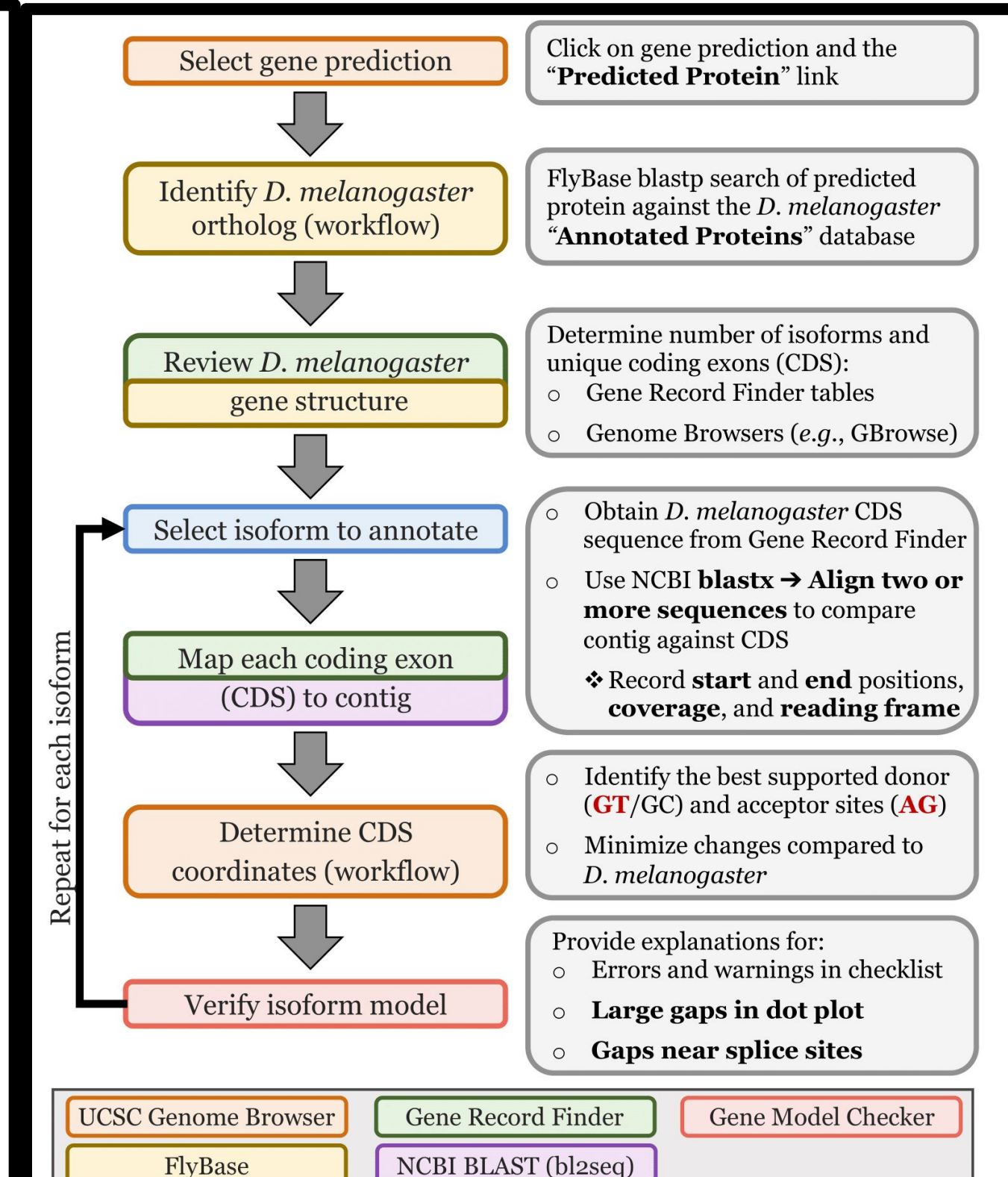
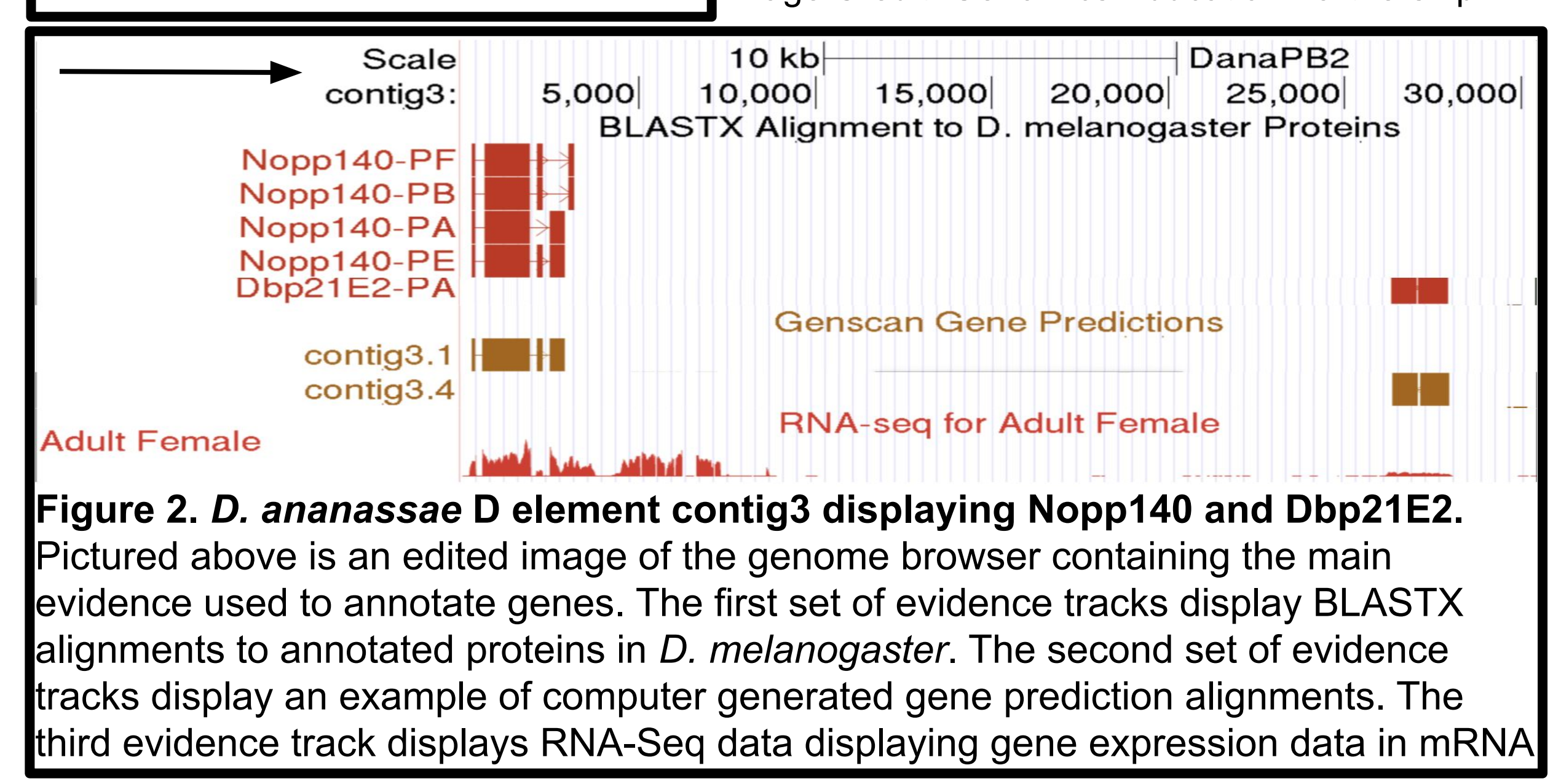


Image Credit: Genomics Education Partnership



Results

A) Nopp140-PA gene model: Exon 1 (267-365), Exon 2 (568-1923), Exon 3 (2483-2908). Splice sites: GT-AG. Start codon (green), stop codon (red).

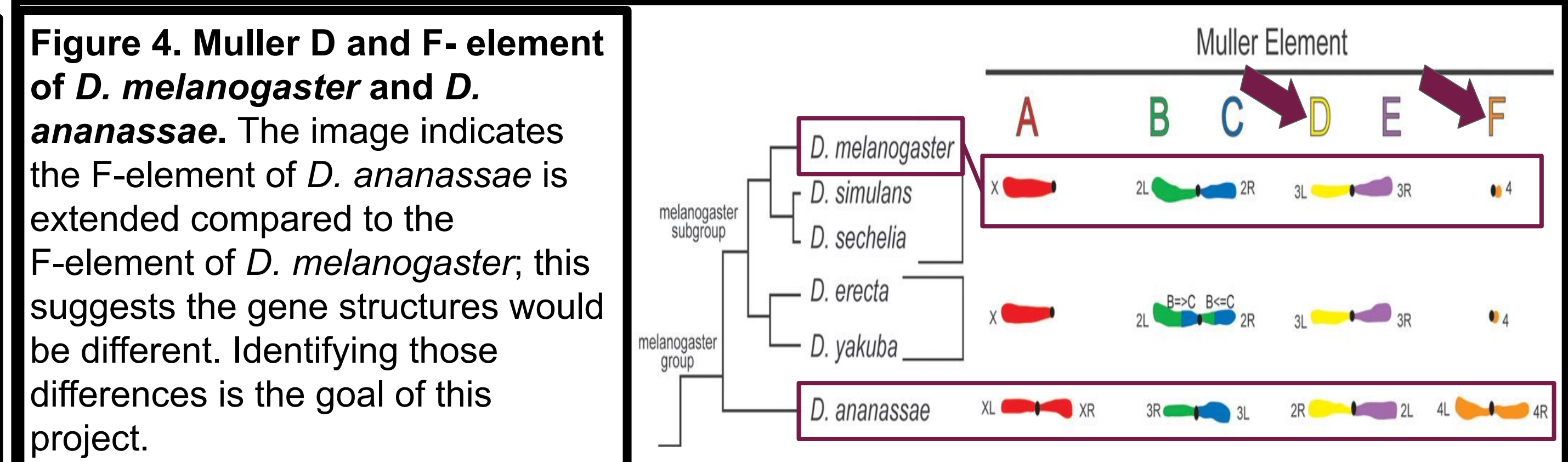
B) Nopp140-PB gene model: Exon 1 (267-365), Exon 2 (568-1923), Exon 3 (2106-2267), Exon 4 (3015-3164). Splice sites: GT-AG. Start codon (green), stop codon (red).

C) Dbp21E2-PA gene model: Exon 2 (26346-27095), Exon 1 (27147-27992). Splice site: GA TG. Start codon (green), stop codon (red).

Figure 3. *D. ananassae*'s Nopp140 and Dbp21 Gene Models. Gene models including labeled introns (standard) and exons (bold), splice donor (GT) and acceptor (AG) sites, start (green box) and stop codons (red box), and coordinates for each exon. Arrows denote beginning and direction of transcription. A) Nopp140-PA gene model. B) Nopp140-PB gene model. C) Dbp21E2-PA gene model.

Future directions

- Dot plots of our annotations (comparison of sequence similarity to *D. melanogaster*) did not come out as expected. This highlights the need for independent corroboration.
- The genes P5CDh1 and CG14565 still require annotation, as well as other contigs of *D. ananassae*. Contig 12 of Muller D element and contig 59 of Muller F element will be annotated next.



Conclusions: Why does this matter?

- The process of Gene Annotation can be used to find and characterize genes in previously uncharacterized genomes.
- Applications:
 - These data allow us to answer questions about evolutionary conservation and protein function across species (including in humans).
 - For example: If there is a human protein of unknown function involved in a disease we can search the genomes of other species to find a similar protein (an ortholog) to make an educated guess as to what the human function is. Gene annotation allows for these educated guesses.

Acknowledgements

- Genomic Education Partnership. Prof. Norma Velázquez-Ulloa.