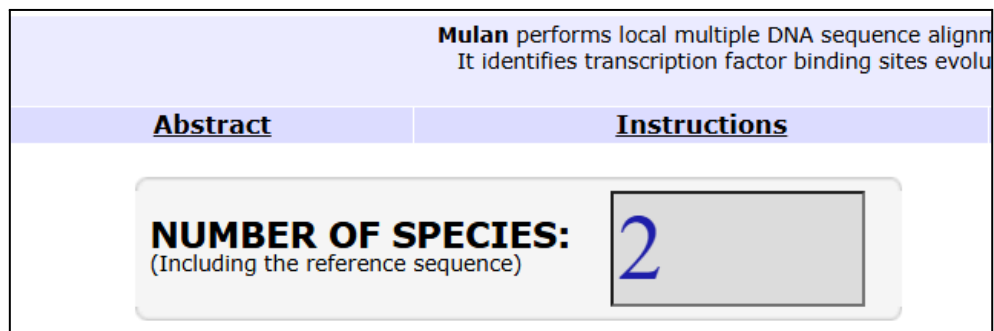**Supplementary Document 1.** Detailed workflow for predicting Transcription Factor Binding Sites (TFBS) using the MultiTF tool within the Mulan platform (MUltiple sequence Local AligNment and conservation visualization).

---

1. Visit Mulan website (https://mulan.dcode.org/).
2. Select the desired number of species for the alignment (see Figure 1).



Figure 1

3. In the "ALL FINISHED SEQUENCES :: TBA alignment" section, click on the "**SELECT**" option to access the MultiTF tool (see Figure 2).
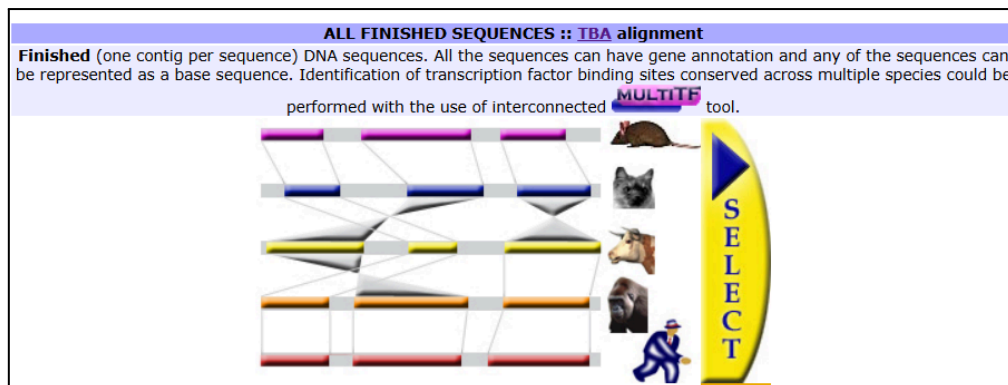


Figure 2

4. Upload the sequences in FASTA format into the designated fields (see Figure 3), then click "**SUBMIT**". For regulatory element analysis, providing an "Annotation" file is not required.

Figure 3

5. Once the sequences are submitted, the system generates a "pitstop" page (Figure 4) that presents the phylogenetic relationships among the compared sequences. Press "**Continue**" on this page (Figure 4).



Figure 4

6. Select the "**MultiTF**" option to proceed with TFBS prediction (Figure 5).



Figure 5

7. On the following page, select "TRANSFAC professional V10.2 library" > "vertebrates" > "Optimized for function" (Figure 6), then click "**SUBMIT**".



Figure 6

8. On the "SELECT TRANSCRIPTION FACTORS" page, click "**SELECT ALL**" and then "**SUBMIT**" (Figure 7).



Figure 7

9. Press "**CHECK IT**" (Figure 8).



Figure 8

10. Once the analysis is complete, the system generates the results. To view the list of predicted TFBS, click the link in the "**Summary**" section (highlighted in orange, Figure 9). The summary includes the transcription factor name, binding site sequence, DNA strand orientation, and genomic position (Figure 10).
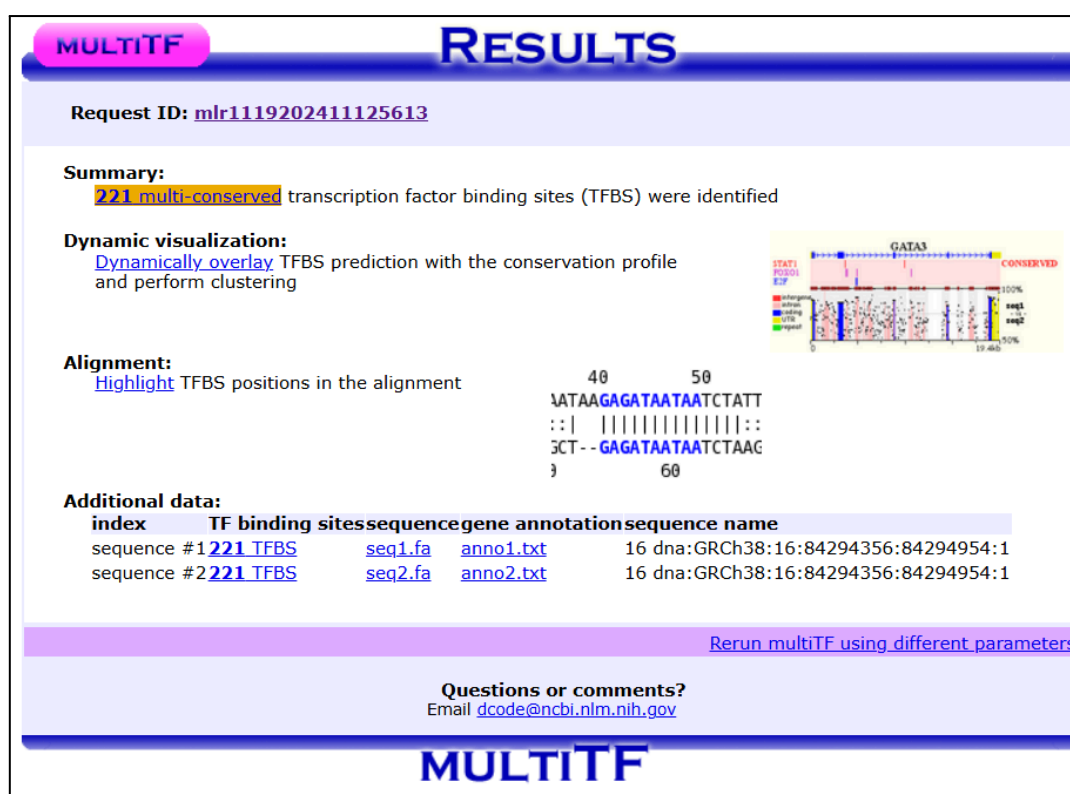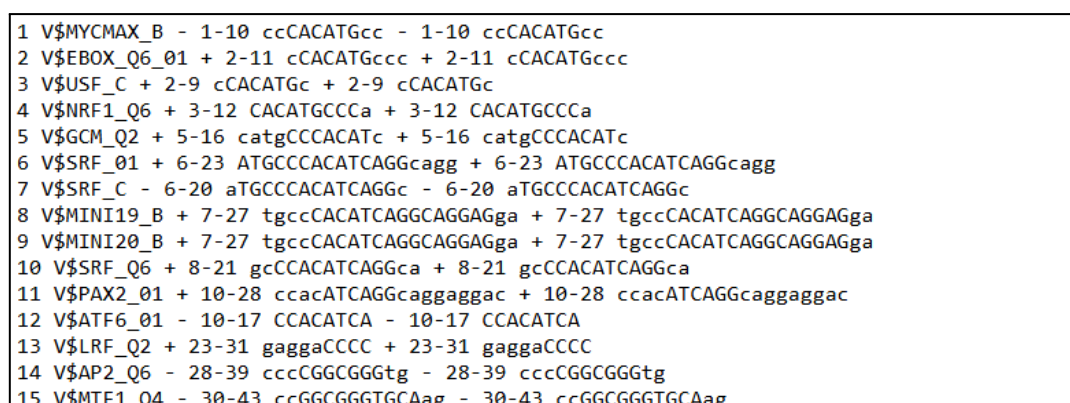
Figure 9



Figure 10

11. Additional information, such as the conservation profile of TFBS across species (if multiple species were analyzed), can be accessed on the Results page by using the "**Dynamically overlay**" feature (Figure 11).
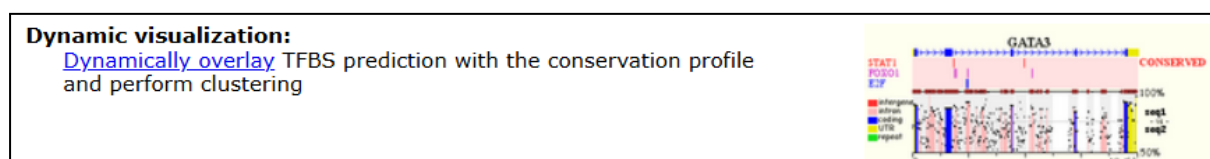


Figure 11