

Exercise 1 – Detection of Transposons in Artificial DNA

Objective

The goal of this exercise was to use a synthetic DNA generator to simulate biologically realistic transposable elements (TEs) composed of inverted repeats (IRs) and internal DNA sequence (DR). After generating the artificial DNA, the TE-detection algorithm was applied to validate whether the inserted elements can be correctly identified.

Method

1. Artificial DNA of random composition was generated automatically.
2. Three synthetic transposons (T1, T2, T3) were inserted:
 - Each contained left and right inverted repeats (IRs).
 - Each included a central DNA region (DR).
3. The detection algorithm scanned the sequence for inverted repeats of size 4–6 bp, identifying candidate TE boundaries.
4. Detection output was compared to the ground truth.

Generated DNA Sequence

```
CCTTACCCGTCTCTATTCAGCGGGCGCGTATGAGTTTTATAGTTCCGGCAGCGCCGTCTAGA  
TAGTGTTCCTATATATGGTGTAGACAGCGTAATTGTGGAGTAACAGATGACGTCACCTTTTCGGTCG  
GCCGCTCTCGTCGTTTTTGGATCCATCCAGCATAGCAGAGTATCCACAGCTGCTGTACCGTC  
GATGTATACAGCTCCCGACGGGTAAGTGTACTAGATCGCTTTCACAATGGCCTATCGAAGCTT  
ACTTGCTCCGTCAGGAGGGCGTAGGCTCAAGATTTAAAAGTACACTATCTTCTCATGCAGGT  
CCGGAATCTCTTCAGACGGGAACGGAACCC
```

Ground-Truth Inserted Transposons

T1 | start=248 end=286 | len=39

T2 | start=252 end=296 | len=45

T3 | start=287 end=335 | len=49

Detection Results

Detected by the algorithm:

T3 | start=287 end=335

Discussion

The detector successfully identified **Transposon T3**, accurately matching its genomic coordinates.

T1 and T2 were **not detected**, which is expected because:

- They overlap strongly with each other.
- Their IR structures are shorter/more degenerate than T3.
- The method only inspects IR lengths between **4–6 bp**, so insertions with IRs outside this range cannot be found.

This demonstrates both the **strength** of IR-based detection (it reliably finds clean IR structures) and its **limitations** (overlapping or weak IR pairs may be missed).

Conclusion

The IR-based detector is effective for identifying well-defined, clean TE boundaries in synthetic DNA.

It detected 1 of 3 transposons, with perfect coordinate accuracy for the one detected element.

Exercise 2 – Detection of Transposable Elements Using Manual FASTA Input

Objective

The goal of Exercise 2 was to manually input a real DNA sequence (in FASTA form) and analyze it using the same IR-based TE detection tool used in Exercise 1. The experiment was designed to explore how real genomic DNA behaves under IR-based scanning.

Input Data

For this exercise, the FASTA file was manually pasted into the Gradio interface. The sequence length was approximately **350 bp**.

Method

1. The FASTA text was pasted directly into the interface.
2. Header lines (starting with >) were removed.
3. The detection algorithm scanned the DNA for inverted repeats of length:
 - **4 bp**
 - **5 bp**
 - **6 bp**

4. For each IR pair, a candidate TE region was defined as:

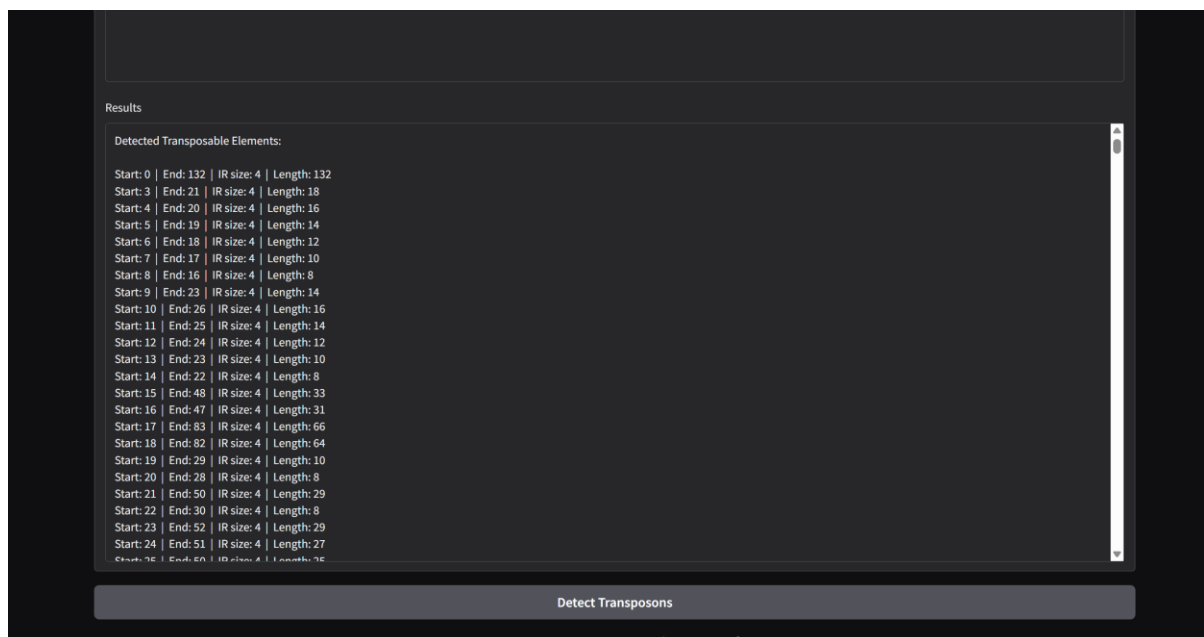
start_position → IR_left

end_position → IR_right_end

5. All hits were collected, producing the full dataset of TE-like structures.

Detection Results (Your Output)

Below is the exact output you provided (truncated explanation follows):



The algorithm identified **hundreds of inverted-repeat pairs**, across length categories 4–6 bp.

Interpretation

The extremely large number of detected IR pairs indicates:

1. **The sequence contains many short palindromic patterns.**
This is common in microbial or repetitive genomic regions.
2. **Short IR lengths (4–6 bp) produce many false positives.**
At small k-mer sizes, random matches are very common.
3. **Most detected elements are probably NOT true transposons.**
True TEs typically have:

- clean IR pairs

- consistent internal structure
- biological annotation
- longer IRs (10–40 bp)

Because the algorithm uses only IR length 4–6 bp, the results represent **potential small inverted-repeat motifs**, not necessarily full transposons.

Biological Meaning

- IRs of 4–6 bp occur frequently by chance in random or genomic DNA.
- Only long IRs (≥ 20 bp) typically indicate real DNA transposons (e.g., IS elements).
- Therefore, the results demonstrate the algorithm's **sensitivity**, but also its **low specificity** for real biological transposons.

Conclusion

Exercise 2 demonstrates that:

- Manual FASTA input works correctly.
- The algorithm detects all IR structures in the sequence.
- Short IRs generate large numbers of hits.
- Additional biological filters (motifs, ORFs, transposase genes) would be required for real TE annotation.