

# Project: Discover and Characterize MicroDNA in a Cancer Cell Line

## Computational Genomics

**Objectives:** Apply your computational genomics knowledge to solve a real-world problem.

**Due Date:** May 7, 11:59PM

**Introduction:** MicroDNA are small, circular, non-coding DNA molecules that were first described in 2012 (<https://www.science.org/doi/10.1126/science.1213307>), and their role remains not fully understood. It is currently believed to influence cellular homeostasis by binding to transcription factors and has been used as a biomarker for cancer. The detection of microDNA from sequencing data is essential for understanding their role in cancer biology.

**What MicroDNA looks like:** When standard linear DNA is sequenced and those reads are aligned to the genome, all of the reads fully align to the genome and are uniformly distributed. When we sequence a molecule of circularized linear DNA, some of the reads will span the circle *junction*, or the point where the two ends of the linear DNA joined. When these reads are aligned to the genome, only a portion will match and the remaining sequence will be *clipped*.

### Linear DNA

...CCCTCA**CCCT**TGGAGAGTCCACAGGTACCAGGGGTTGGTCTGAACCC**CC**AGCACAG...

### Reads

CCCTCA**CCCT** ACCAGGGGTT  
A**CCCT**TGGAGA GGGTTGGTCT  
GGAGAGTCCA GGTCTGAACC  
GTCCACAGGT GAACCC**CC**AG  
CAGGTACCAG CCCAGCACAG

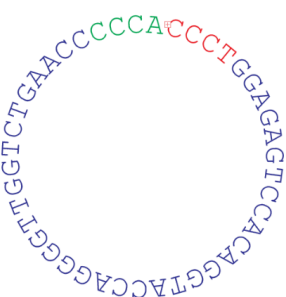
### Reference Genome

AGAAAACCAATCTCGCAGCCCTCACCCCTGGAGAGTCCACAGGTACCAGGGGTTGGTCTGAACCC**CC**AGCACAGAGCACCT

### Alignments

CCCTCA**CCCT** GTCCACAGGT GGGTTGGTCT CCCAGCACAG  
A**CCCT**TGGAGA CAGGTACCAG GGTCTGAACC  
GGAGAGTCCA ACCAGGGGTT GAACCC**CC**AG

### microDNA



### Reads

CCCTTGGAGAG TTGGTCTGAA  
TCCACAGGTA CCC**CC**CA**CCCT**  
CCAGGGGTTG GGAGAGTCCA  
GTCTGAACCC CAGGTACCAG  
CC**CCCT**TGGA GGGTTGGTCT  
GAGTCCACAG GAACCC**CC**AC  
GTACCAGGGG CCTGGAGAGT

### Reference Genome

AGAAAACCAATCTCGCAGCCCTCACCCCTGGAGAGTCCACAGGTACCAGGGGTTGGTCTGAACCC**CC**AGCACAGAGCACCT

### Alignments

CCCTTGGAGAG TCCACAGGTA CCAGGGGTTG CC**CCCT**TGGA  
CC**CCCT**TGGA GAGTCCACAG GTACCAGGGG GTCTGAACCC CCC**CC**CA**CCCT**  
TTGGTCTGAA

↑  
Clipped  
alignment

- M - Match or mismatch (aligned to the reference)
- S - Soft clipping (clipped sequence not included in alignment)

<u>Chromosome</u>	<u>Start</u>	<u>Score</u>	<u>CIGAR</u>	<u>Sequence</u>
NC_000001.10	2383746	60	<b>42M</b>	CCTGCCTGGCAGGTAGCAGCCCCGTGGAAGTATTTTCATCTTG

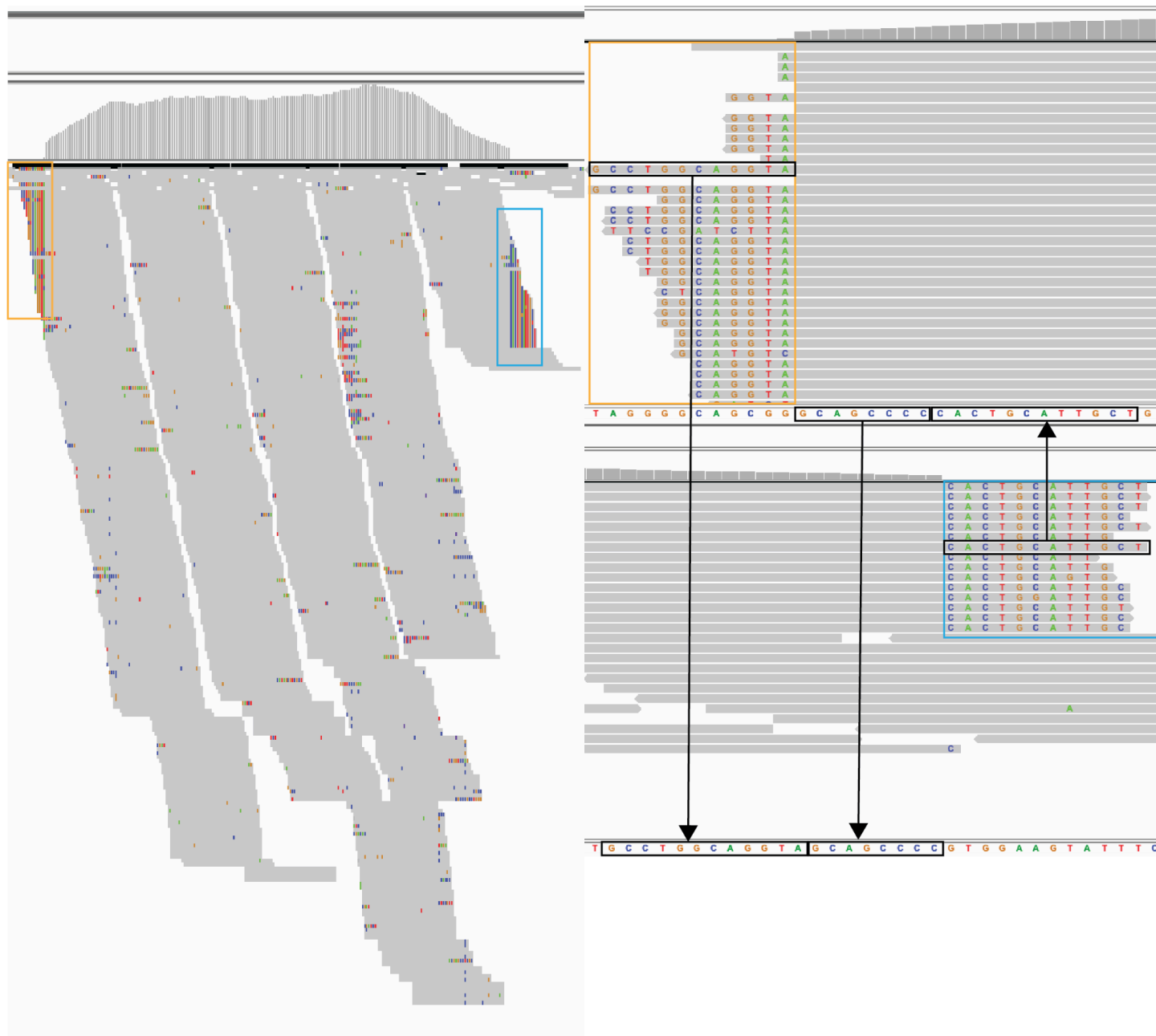
NC_000001.10	2383556	60	<b>12S30M</b>	GCCTGGCAGGTAGCAGCCCCACTGCATTGCTGAGCCTGGAA
NC_000001.10	2383739	60	<b>30M12S</b>	CCCCGGCCCTGCCTGGCAGGTAGCAGCCCCACTGCATTGCT

Ref: TAGGGGCAGCGGGCAGCCCCACTGCATTGCTGAGCCTGGAA  
xxxx| ||||xxx| ||||||||| ||||||||| |||||||||  
Read: GCCTGGCAGGTAGCAGCCCCACTGCATTGCTGAGCCTGGAA

Ref: CCCC GGCCCTGCCTGGCAGGTAGCAGCCCCGTTGAAGTATTT  
| | | | | | | | | | | | | | | | | | | | | x x x x x x x | x x x |  
Read: CCCC GGCCCTGCCTGGCAGGTAGCAGCCCCACTGCAATTGCT

GCCTGGCAGGTA**GCAGCCCC**ACTGCATTGCTGAGCCTGGAA  
 |||||  
 CCCCggccctGCCTGGCAGGTA**GCAGCCCC**ACTGCATTGCT

Using IGV (<https://igv.org/>), we can visualize all of the evidence for a circle:



**Assignment:** Develop an algorithm to detect and quantify microDNAs from alignment data. In addition to your code (in a GitHub repo that follows the standard for course homeworks), write a short (~2 pages) report that details your algorithm and results. Your report should include how you validated a subset of your circles, and a link to your GitHub repository.

Your algorithm should take as input a BAM file (provided). Reading BAM files in Python is best accomplished with the pysam library (<https://github.com/pysam-developers/pysam>) which has extensive documentation (<https://pysam.readthedocs.io/en/latest/api.html>). For example, you can find alignments with softclips in the CIGAR string with:

```
import pysam

# Path to your BAM file
bam_file = "your_file.bam"

# Open the BAM file
with pysam.AlignmentFile(bam_file, "rb") as bam:
    for read in bam:
        # Check if the read has soft clips in the CIGAR string
        if any(cigar_op[0] == 4 for cigar_op in read.cigartuples):
            print(f"Read {read.query_name} has soft clips: {read.cigarstring}")
```

Your algorithm will scan the bam file for alignments that indicate the existence of a circle and aggregate the alignments that seem to be evidence for the same circle. For each suspected circle, define a metric that indicates how likely the circle is real and not an unrelated alignment artifact (e.g. how well the start and end tags agree, the number of junction tags, the total depth of circle DNA, etc), and report all circles that meet some predefined threshold. This report should give the position and score of each circle.