**Turn on your Laptop**

**Connect to Internet (Check your connection)**

**Launch WSL (Windows) or Terminal (macOS/Linux)**

# Tutorial Outline

1. Downloading data from NCBI
   a. RefSeq vs GenBank
   b. FNA vs GTF vs GFF vs GBFF

2. Brief introduction to k-mers, GC, Clump and Origin

3. Practice Tasks
   a. ORI signal checker: K-mer enrichment and plotting
   b. Clump finder: L,k,t clumps
   c. GC skew calculator
   d. ORI (Origin of Replication) finder

**National Center for Biotechnology Information (NCBI)**

1. **Public bioinformatics resource** maintained by the NIH (US).
2. Provides access to a w**ide range of biological databases, analysis tools, and reference datasets**.
3. Covers **genomics, transcriptomics, proteomics, and biomedical literature** datasets.

**Key NCBI Databases:**
1. **GenBank**: comprehensive, public repository of sequence data
2. **RefSeq**: non-redundant, curated, and standardized reference sequence data.

**Common NCBI File Formats:**
1. FASTA Nucleotide (**FNA**): Genomic Sequence
2. Gene Transfer Format (**GTF**): gene and transcript annotations
3. General Feature Format (**GFF**): genomic features such as genes, CDS, exons, etc.
4. GenBank Flat File (**GBFF**): human-readable annotated genome file

# Downloading Data

# Downloading Data *Continued*

# Brief introduction to k-mers

1. Substring of length k extracted from a biological sequence data.

2. Have application in genome assembly, sequence comparison and and clustering, etc.



4-mers

# Brief introduction to GC

1. Refers to the percentage of guanine (G) and cytosine (C) nucleotides in a DNA sequences.

2. GC content affects gene density, replication and transcription efficiency, etc.



Nucleobases of DNA

DNA
Deoxyribonucleic acid

*Image Source:*
*https://www.technologynetworks.com/genomics/articles/what-are-the-key-differences-between-dna-and-rna-296719*

# Brief introduction to (L, k, t)-Clumps

1. Identify locally frequent k-mers within a DNA sequence.

2. A k-mer is said to form an (L, k, t)-clump if it appears at least t times within any window of length L in the genome

3. Given:
   - k → length of the k-mer
   - L → length of the sliding window
   - t → minimum number of occurrences

# Brief introduction to Origin (ORI)

1. Origin of Replication (ORI) is a specific genomic region where DNA replication begins.

2. Often AT-rich, making strand separation easier

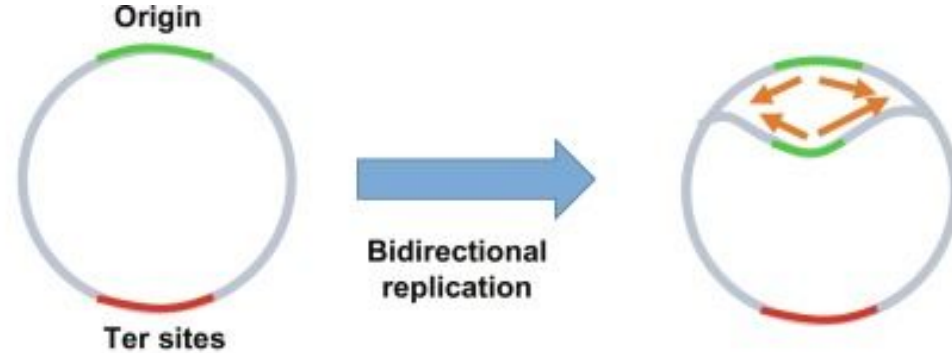3. Usually single ORI in bacteria and multiple ORIs in eukaryotes



*Image Source: https://doi.org/10.1016/B978-0-323-91788-9.00006-5*

# Vibe Coding Session

Practice Tasks
   a. ORI signal checker: K-mer enrichment and plotting
   b. Clump finder: L,k,t clumps
   c. GC skew calculator
   d. ORI (Origin of Replication) finder



Vibe Coders looking at their own code after exhausting their credits

(confused unga bunga)

# Thank You