# Finding patterns in Genome Data

Project Study - Summer Semester 2024

Presentation Date: 19[th] July 2024, Friday

## Masters in Professional IT Business and Digitalization

**htw.**
Hochschule für Technik
und Wirtschaft Berlin

University of Applied Sciences

# Project Guidance

- Prof. Piotr Wojciech Dabrowski

  Hochschule für Technik und Wirtschaft (HTW) Berlin

# Team Members

- Kalyan Shencottah – ██████
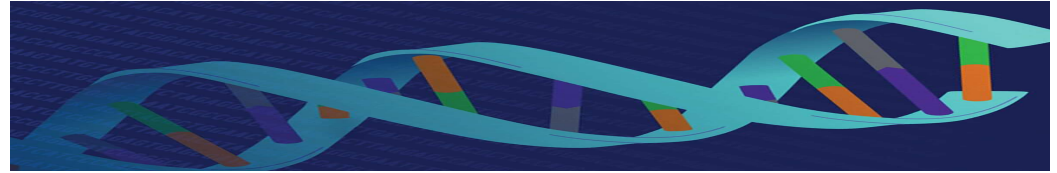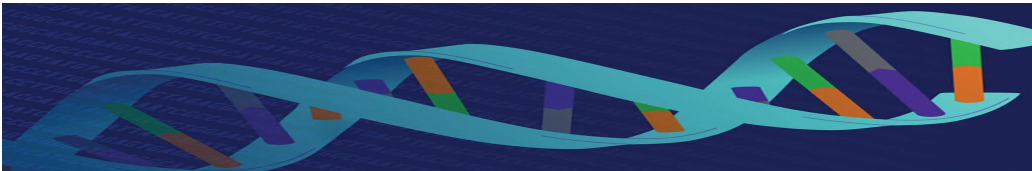
# Agenda



- *Project Study Overview*
  - Motivation
  - Objective

- *The Approach*
  - Conceptual Understanding of Genome Data
  - Data Collection from Genome Databases
  - Genome Data Interpretation
  - Things Accomplished

- *Technical Implementation*
  - Architecture, Technical Components, Algorithms
  - Processing Coding Sequences
  - Processing Non-Coding Sequences
  - Generate clean data

- *Challenges and Resolutions*
  - Example Challenges and Resolutions
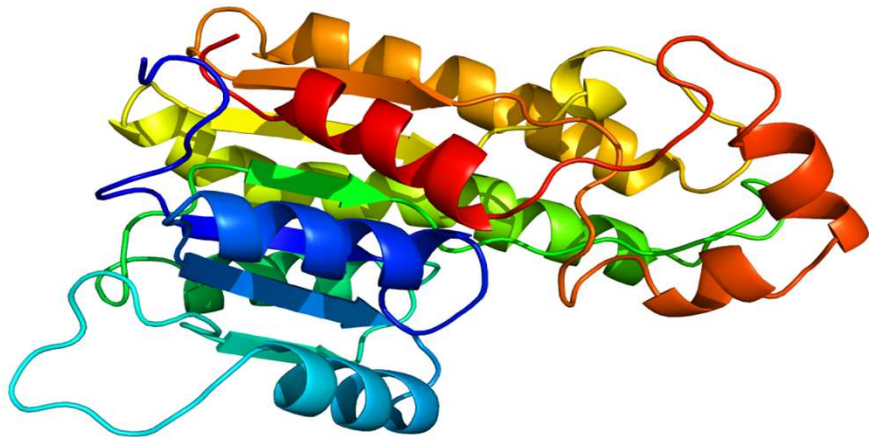
- *Conclusion and Next Steps*

- *Extras – Python Code*



**htw.**
Hochschule für Technik
und Wirtschaft Berlin

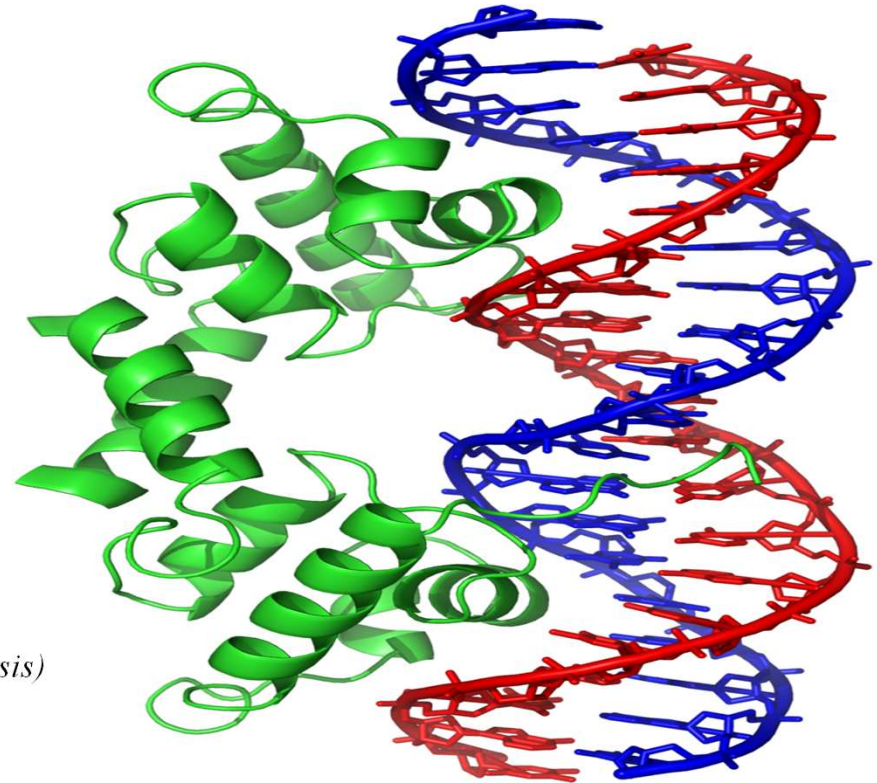University of Applied Sciences

# Global Motivation



- *Genome Data*

  - Evolving field of Bio-Informatics
  - Complex Data – DNA, Proteins
  - Discover Mutations, Cause of Diseases



- *Many ML Algorithm Possibilities*

  - Supervised and Unsupervised
  - Ex. Classification, Clustering, Neural Network



**htw**

**Hochschule für Technik
und Wirtschaft Berlin**

**University of Applied Sciences**

# Personal Motivation

- *To gain deeper understanding on BioInformatics / Computational Biology*

- *Apply my Knowledge gained*

  - *Machine Intelligence Lab, Univ of Cincinnati, USA*

    - *ML, Pattern Recognition, Spatial Data Mining(Thesis)*

  - *HTW Course work*

- *Python and Sequence Analysis*

# Objective

- ***Research on Genome Data / Databases***

  - *NCBI / Ensembl Databases*

  - *Identify Organism (say, Bacteria)*

  - *RefSeq files (Sequence files / FASTA format)*

  - *Feature files (GFF files)*

- ***Build a Pipeline***

  - *Extract Sequence files from Genome Databases*

  - *Retrieve Data pertaining to Features (CDS, Gene)*

  - *Identify Coding and Non-Coding Regions*

  - *Store Coding and Non-Coding regions in FASTA files*

  - *Preparation of Training and Test Data to be fed to NN*



**htw.**

**Hochschule für Technik
und Wirtschaft Berlin**

**University of Applied Sciences**

# Conceptual Understanding of Genome Data
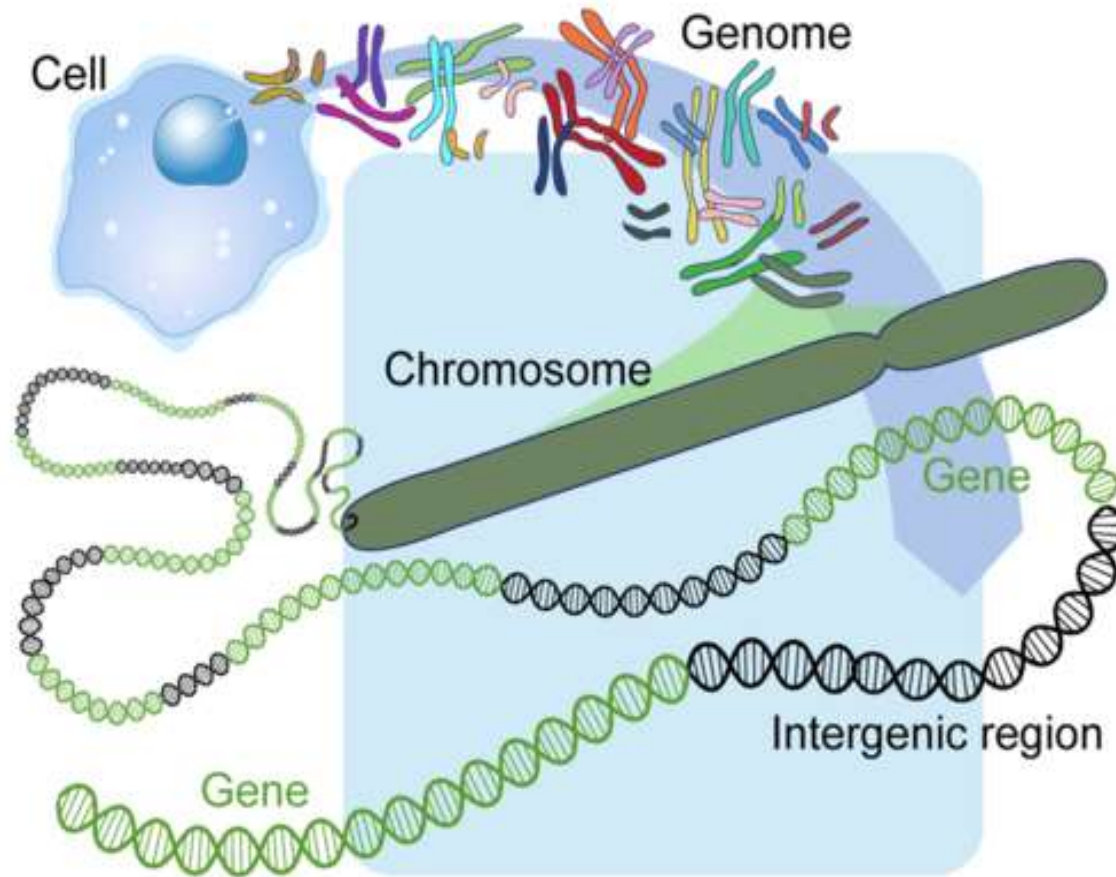


COUPLED



NON-LINEAR



COMPLEXITY



PLASTICITY



NON-EQUILIBRIUM



COMPLEX PHENOTYPES

htw

**Hochschule für Technik und Wirtschaft Berlin**

University of Applied Sciences

# Conceptual understanding of Genome



Two types of Organisms: Prokaryotes(Bacteria), Eukaryotes (Humans, Animals, Plants)

# Data Collection from Genome Databases

# Interpreting Genome Data (Prokaryotes-Bacteria)



UTR – Untranslated Regions

RBS – Ribosome Binding Sites
Codons – 3 Nucleotides – START and STOP
ORF – Open Reading Frame

Source: Wikipedia

# Interpreting Genome Data – Mapping Codon Table



Start Codon – AUG (ATG)
    Initiates translation process

Stop Codon – UAA, UAG, UGA (TAA, TAG, TGA)
    Initiates the termination

| **DNA** | **(m)RNA** |
| --- | --- |
| A – Adenine | A |
| T – Thymine | T |
| C - Cytosine | U - Uracil |
| G - Guanine | G |

Source: Wikipedia
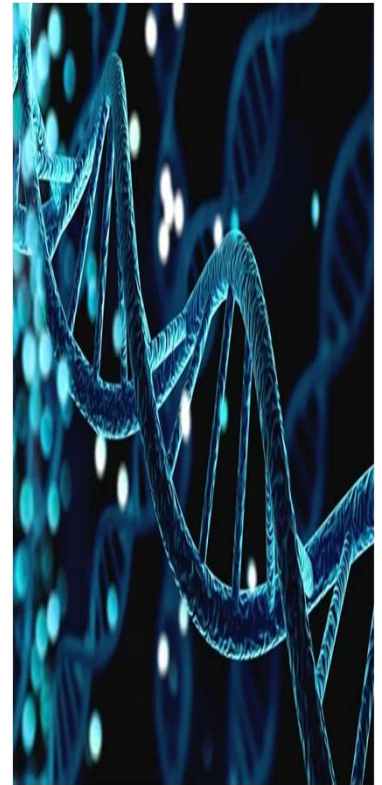
# Things Accomplished

- **Download Data from NCBI Database**

  - *RefSeq/GenBank (FASTA files, GFF files)*

  - *FASTA – Sequence files, GFF – Features*

- **Processing of Fasta Sequences**

  - *Create separate files for Genes (say, AE006468.2, AE006471.2)*

  - *Cleanup unwanted characters(say, N), spaces*

- **Processing of GFF Files (Features)**

  - *Create Features DB (SQLLite)*

  - *Create different GFF files for Genes (say, AE006468.2, AE006471.2)*

- **For the Sequences, for given feature (say, CDS)**

  - *Identify Coding and Non-Coding Sequences*

  - *Process Forward (+) Strands and Reverse (-) Strand. Reverse complement (-) strands*

  - *Generate Coding Sequences in FASTA format, Non-Coding Sequences in FASTA format*

  - *Finalize Training and Test Data to be fed to NN*



**htw**

Hochschule für Technik
und Wirtschaft Berlin
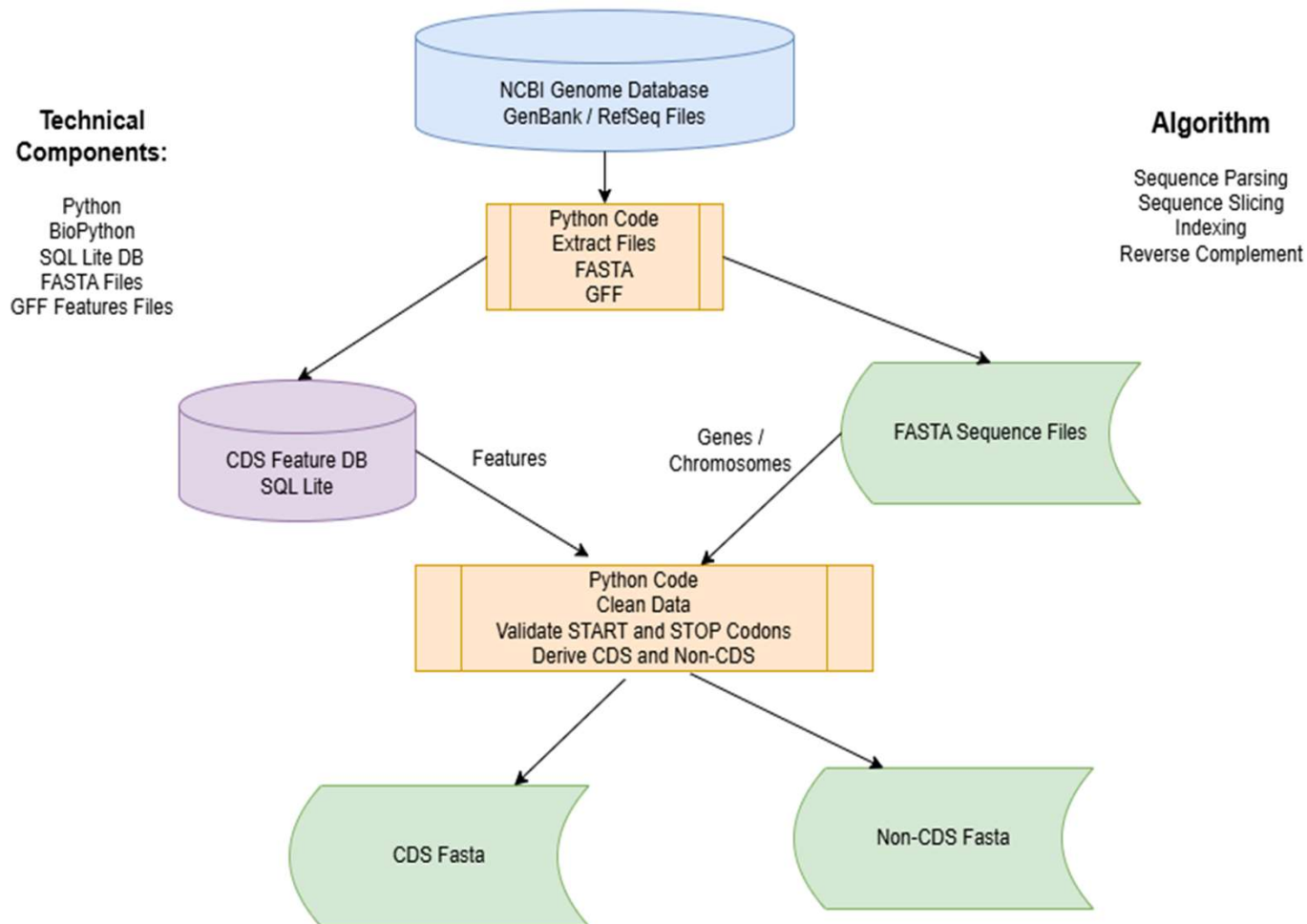
**University of Applied Sciences**

# Technical Implementation

# Architecture For Finding Patterns in Genome Data

**Technical Components:**

Python
BioPython
SQL Lite DB
FASTA Files
GFF Features Files

**Algorithm**

Sequence Parsing
Sequence Slicing
Indexing
Reverse Complement

NCBI Genome Database
GenBank / RefSeq Files

Python Code
Extract Files
FASTA
GFF

CDS Feature DB
SQL Lite

Features

Genes /
Chromosomes

FASTA Sequence Files

Python Code
Clean Data
Validate START and STOP Codons
Derive CDS and Non-CDS

CDS Fasta

Non-CDS Fasta

htw.
**Hochschule für Technik
und Wirtschaft Berlin**

**University of Applied Sciences**

# Generate Output in FASTA – CDS and Non-CDS

## CDS – Coding Sequences FASTA Output



## Non-CDS – Non-Coding Sequences FASTA Output

# Challenges and Resolutions



- **Finding an Organism to work with**

  - *Choose Bacteria.  Understand complexity, clear goal on what data to work with*

- **Memory Issues in handling Datasets**

  - *Reduce the Data to execute Logic.*

- **Understanding Sequences, Features**

  - *Gain knowledge on the Features, its Data Structure*

htw

**Hochschule für Technik
und Wirtschaft Berlin**

**University of Applied Sciences**

# Next Steps

## Bio-Polymer (alphabet)

DNA (A, T, C, G)

mRNA (U, A, C, G)

Proteins (20 amino acids)

## Process(algorithm)

replication

transcription

splicing

translation

folding

Interactions

## Problem -> Algorithms

Sequencing -> Fragment assembly problem -> The shortest superstring problem

Gene Finding -> Hidden Markov Models, Pattern Recognition Methods

Sequence comparison -> pairwise and multiple sequence alignments -> Dynamic programming algorithm, Heuristic Methods

Finding Unknown Patterns -> CNN

Source: folding.chmcc.org, Intro to BioInformatics, Professor Jarek Meller

# References

- Deep Modeling of DNA Sequences with Python and Keras, Martin Preusse, Gokcen Eraslan(Researchers at Helmholtz Munich)

- Intro to BioInformatics , Prof. Jarek Meller, Cincinnati Childrens Hospital Research(www.chmcc.org)

- Demystify DNA Sequencing with Machine Learning and Python, The AI Dream, Oct 26, 2020

# Python Coding

# (Extra, if time permits)

# Processing Sequence Files – Python Code

**Install BioPython, PyTorch, GFFUtils**

**Import BioPython Libraries**

```
## Code is to Open and read the file\n",
BCTGCA20 = open('dataset\\data1\\GCA_000006945.2_ASM694v2_genomic.fna')
readBCTGCA20 = BCTGCA20.read()
```

**Read FASTA File**

```
## Reading Sequences using BioPython
## Using BioPython we can extract ID, Name, Description, Number of Features, Seq of the Genes
## Seq stores the sequence

fBactDataGCA="C:\\Kalyan\\GeneSequence\\dataset\\data1\\GCA_000006945.2_ASM694v2_genomic.fna"
sequences=[i for i in SeqIO.parse(fBactDataGCA, 'fasta')]
```

**Parse FASTA File**

```
SeqIO.write(gene_of_interest,'C:\\Kalyan\\GeneSequence\\dataset\\data1\\GCA_AE006468_2.fna', 'fasta')
```

```
SeqIO.write(gene_of_interest,'C:\\Kalyan\\GeneSequence\\dataset\\data1\\GCA_AE006471_2.fna', 'fasta')
```

**Identify Genes of Interest**

**Store them in Separate Fasta Files**

```
ID: AE006468.2
Name: AE006468.2
Description: AE006468.2 Salmonella enterica subsp. enterica serovar Typhimurium str. LT2, complete genome
Number of features: 0
Seq('AGAGATTACGTCTGGTTGCAAGAGATCATGACAGGGGGAATTGGTTGAAAATAA...ATA')
```

**Understand Sequence Data**

**SeqID, No. of Genes, Length, Extra Characters, Spaces**

# Processing Features GFF – Python Code

```
## Analyzing GFF file
## pip install gffutils
## pip install bcbio.gff
## Reference tutorial https://daler.github.io/gffutils/


import gffutils
import pprint
from BCBio.GFF import GFFExaminer
from Bio import SeqIO
from Bio.Seq import Seq
##from Bio.Alphabet import IUPAC
```

**Import BioPython, GFF Libraries**

```
featTypes = ['gene','CDS','mRNA','exon','intron', 'utr']
for feat in featTypes:
    dbfeatcount = dbgff.count_features_of_type(feat)
    print("DB GFF Count Features:", feat, dbfeatcount)
```

```
DB GFF Count Features: 14390
DB GFF Count Features: gene 4678
DB GFF Count Features: CDS 4555
DB GFF Count Features: mRNA 0
DB GFF Count Features: exon 118
DB GFF Count Features: intron 0
DB GFF Count Features: utr 0
```

**Understand Features**
**CDS**
**Gene**
**Exon**
**Intron**
**mRNA**

**Process CDS and Non-CDS + and - strands**

```
Feature on genome with + strand AE006468.2 from 7665 to 8618, strand: +
> CDS : Gene sequence ID AE006468.2 : start 7665 : end 8618 : strand: +

Feature on genome with - strand AE006468.2 from 10092 to 10805, strand: -
> CDS : Gene sequence ID AE006468.2 : start 10092 : end 10805 : strand: -
AE006468.2

Non-Coding Region Start : 256, Non-Coding Region End: 336
CGCGTACAGGAAACACAGAAAAAAGCCCGCACCTGAACAGTGCGGGCTTTTTTTTCGACCAGAGATCACGAGGTAACAACC
```

```
dbgff = gffutils.FeatureDB('bactGFFDB.db', keep_order=True)

##Prints the schema of the FeatureDB
dbgffschema = dbgff.schema()
print("DB gff schema:", dbgffschema)
```

**Create Database**

![htw Hochschule für Technik und Wirtschaft Berlin — University of Applied Sciences](logo)

www.htw-berlin.de