# Chapter 01. Introduction

**Python Programming for Bioinformatics**

**Robert C. Chi**

# Agenda

- **About This Course**

- **Introduction to Biopython**

- **Installing Biopython**

- **A Quick Tour of Biopython**

# ABOUT THIS COURSE

# Robert C. Chi (紀俊男)

- **Education**
  - Ph.D. Candidate / Bioinformatics
    *Taiwan International Graduate Program (TIGP), 2003-2007*
  - Master / Computer Sciences
    *Queens College, CUNY, 1994-1996*
  - Bachelor / Computer Sciences
    *Fu-Jen Catholic University*

- **Experience**
  - Training Director / AMI (2014-2020)
  - Founder / Hatch Information Co., Ltd. (2007-2013)
  - Research Assistant / Academia Sinica (2000-2007)
  - Manager of Tech Support / Trend Micro Co., Ltd. (1998-2000)
  - Game Developer / CG Animation Co., Ltd. (1997-1998)

- **Expertise**
  - Artificial Intelligence (AI), Embedded System, Computer Security, Game Programming.

# Syllabus

- **Part I. Python (10 Hr)**    *DONE*
  - Python & Environments
  - Literals & Variables
  - Input & Output
  - Branch & Loop
  - String Manipulation
  - Compound Data Types
    - Tuple, List, Dictionary, Set
  - Functions
  - Data Science Packages
    - NumPy, Pandas, MatPlotLib

- **Part II. BioPython (14 Hr + 7Hr)**
  - Introduction
  - Read/Write Bioinformatic Files
    - FASTA, GenBank, SwissProt, ExPASy, KEGG…
  - Sequence Manipulation
    - Transcription, Translation, Alignment
  - Databases Handling
    - BLAST, NCBI Entrez…
  - Working with 3D Structures
  - Machine Learning
    - Data Pre-Processing
    - Classification
    - Clustering

# Schedule & Environment

- **Schedule**
  - Part I (5 Weeks)
    - 2021/08/06 ~ 2021/09/03
    - Fri. 15:00 ~ 17:00
  - Part II (7 Weeks)
    - 2021/10/29 ~ 2021/12/10
    - Fri. 14:00 ~ 17:00 (2Hr + 1Hr)

- **Location: Online**
  - https://www.gotomeet.me/TeacherChi/BioPython

- **Environments**
  - Google Colab

- **Lecturing in**
  - English

- **Teaching Style**
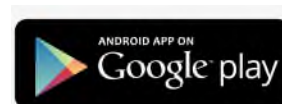  - Part I: Lecturing
  - Part II: Lecturing + Practiceing

# Live Broadcasting

- **URL: https://www.gotomeet.me/TeacherChi/BioPython**



**App Download**

**GoToMeeting**



Meeting ID： 117-684-245

# Pre-requisites



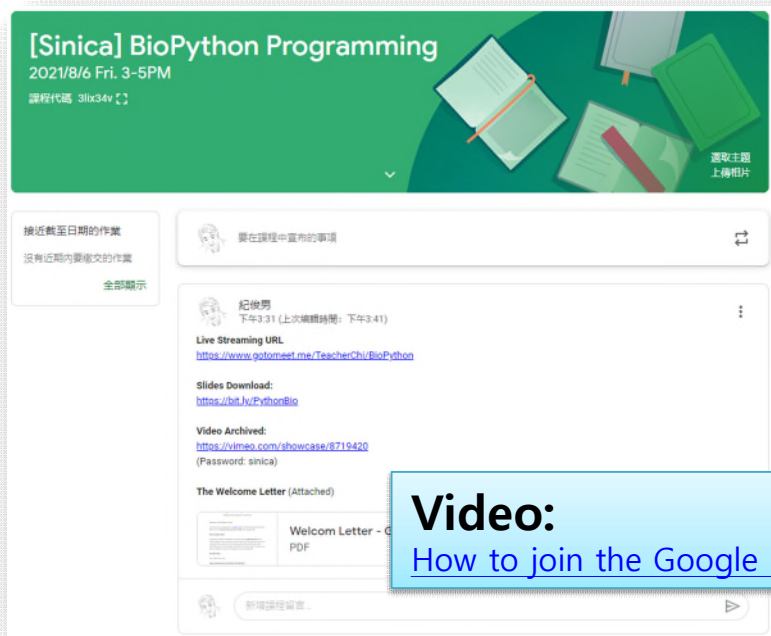**Python Programming**
(Flow Controls, Compound Data...etc.)

**Google Colab**
(Development Environment)

- **Google Classroom：https://bit.ly/BioPy-202108**



**App Download**



**Google Classroom**

**Video:**
How to join the Google Classroom

**Join ID: 3lix34v**

# INTRODUCTION TO BIOPYTHON

- **The most popular Python package for computational molecular biology.**

biopython

**https://biopython.org/**

**Official Name:**

**Biopython** ✅

**BioPython** ❌

# Abilities of Biopython

- **Communicating with Bioinformatic Services**
  - NCBI Services (Blast, Entrez, PubMed)
  - ExPASy Services (Swiss-Prot, Prosite entries, Prosite searches)
- **Parsing Bioinformatic Files**
  - Blast, Clustalw, FASTA, GenBank, PubMed/Medline, ExPASy (Enzyme/Prosite), SCOP ('dom' & 'lin' files), UniGene, SwissProt
- **Performing Sequence Operations**
  - Translation, transcription, weight calculations, alignments…etc.
- **Performing Classification of Data**
  - k Nearest Neighbors, Naive Bayes, Support Vector Machines…etc.
- **Integrating with BioSQL**
  - A sequence database schema also supported by the *BioPerl* and *BioJava* projects.

# How to Cite Biopython?

- **The main Biopython reference**
  - [Cock *et al.*, 2009] ([URL](URL))
- **The official project announcement**
  - [Chapman and Chang, 2000] ([URL](URL))
- **For Bio.PDB**
  - [Hamelryck and Manderick, 2003] ([URL](URL))
- **For Bio.Cluster**
  - [De Hoon *et al.*, 2004] ([URL](URL))
- **For Bio.Graphics.GenomeDiagram**
  - [Pritchard *et al.*, 2006] ([URL](URL))
- **For Bio.Phylo and Bio.Phylo.PAML**
  - [Talevich *et al.*, 2012] ([URL](URL))
- **For the FASTQ file format as supported in Biopython, BioPerl, BioRuby, BioJava, and EMBOSS**
  - [Cock *et al.*, 2010] ([URL](URL))

# INSTALLING BIOPYTHON

# Install Biopython

**pip   install   biopython**

**Python Installation Program**
(Tool for Installing Packages)

**Sub-command of pip**
(to install a package)

**The name of package**

```
1    import Bio
2    print(Bio.__version__)
```

1.79

- **Install Biopython**

    – Open a Colab page called "BiopythonInstall.ipynb".

    – Write and run the following codes:



```
1    !pip install biopython

Requirement already satisfied: biopython in /usr/local/lib/python3.7/dist-packages (1.79)
Requirement already satisfied: numpy in /usr/local/lib/python3.7/dist-packages (from biopython) (1.19.5)

[5]  1    import Bio
     2    print(Bio.__version__)

1.79
```

**(Solution URL of this Practice)**

# A QUICK TOUR OF BIOPYTHON

# Convert Strings to Bio.Seq.Seq

```python
1   # Install Biopython
2   !pip install biopython
3
4   # Import Bio.Seq.Seq
5   from Bio.Seq import Seq
6
7   # Convert String to Bio.Seq.Seq
8   my_seq = Seq("AGTACACTGGT")
9   print(my_seq)
10
11  # Convert Bio.Seq.Seq to String
12  bio_seq = str(my_seq)
13  print(bio_seq)
14  print(type(bio_seq))
```

**Install Biopython onto Colab**

**Import Bio.Seq.Seq**

```
AGTACACTGGT
```

```
AGTACACTGGT
<class 'str'>
```

# Practice

- **Convert Strings to Bio.Seq.Seq**
    - Write and Run the following codes on a Colab page called "QuickTour.ipynb":

```
▼ Install Biopython

✓  [1]   1   !pip install biopython

▼ Convert String to Bio.Seq.Seq

✓  [2]   1   from Bio.Seq import Seq
         2
         3   my_seq = Seq("AGTACACTGGT")
         4   print(my_seq)

▼ Convert Bio.Seq.Seq to String

✓  [3]   1   bio_seq = str(my_seq)
         2   print(bio_seq)
         3   print(type(bio_seq))
```

**(Solution URL of this Practice)**

# Complement of a Sequence

DNA coding strand (aka Crick strand, strand +1)

5' ATGGCCATTGTAATGGGCCGCTGAAAGGGTGCCCGATAG 3'

| | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | |

3' TACCGGTAACATTACCCGGCGACTTTCCCACGGGCTATC 5'

DNA template strand (aka Watson strand, strand −1)

## Complement

A → T          C → G

T → A          G → C

## Reverse Complement

```
5'    ATGGCCATTGTAATGGGCCGCTGAAAGGGTGCCCGATAG 3'
      |||||||||||||||||||||||||||||||||||||||
3'    TACCGGTAACATTACCCGGCGACTTTCCCACGGGCTATC 5'  ← Complement
5'    CTATCGGGCACCCTTTCAGCGGCCCATTACAATGGCCAT 3'  ← Reverse Complement
```

# Complement of a Sequence

```python
from Bio.Seq import Seq
my_seq = Seq("AGTACACTGGT")

# Complement & Reverse Complement of a Sequence
print(my_seq.complement())
print(my_seq.reverse_complement())
```

TCATGTGACCA
ACCAGTGTACT

# **Practice**

- **Complement of a Sequence**
  - Write and Run the following codes on a Colab page called "QuickTour.ipynb":



```
Sequence Complement & Reverse Complement

[4]   1   print(my_seq.complement())
      2   print(my_seq.reverse_complement())
```

**(Solution URL of this Practice)**

# Parsing FASTA Files

- **"ls_orchid.fasta" ( Download Link )**

# Parsing FASTA Files

- **Upload "ls_orchid.fasta" to Colab (by the Local File)**

# Parsing FASTA Files

- **Upload "ls_orchid.fasta" to Colab (by fetching from Internet)**

```
34  # Download Data File: ls_orchid.fasta
35  import os
36  if not os.path.isfile("ls_orchid.fasta"):
37    os.system("wget https://raw.githubusercontent.com/biopython/biopython/master/Doc/examples/ls_orchid.fasta")
```

**isfile(): =True if exist, =False if miss**

**wget: Web GET.  Download file by URL.**

**os.system(): Send the command to OS.**

Files

.. 
sample_data
ls_orchid.fasta

# Parsing FASTA Files

```
1  from Bio import SeqIO
2  for seq_record in SeqIO.parse("ls_orchid.fasta", "fasta"):
3      print(seq_record.id)
4      print(seq_record.seq)
5      print(len(seq_record))
```

**1** Parse the file in FASTA format

**2** Get one record a time

**3** Print ID, Sequence, and Length of the current record

1st SeqRecord
```
gi|2765582|emb|Z78457.1|PCZ78457
CGTAACAAGGTTTCCGTAGGTGAACCTCCGGAAGGATCATTGTTGAGATCACATAATAATTGATC
739
```

2nd SeqRecord
```
gi|2765581|emb|Z78456.1|PTZ78456
CGTAACAAGGTTTCCGTAGGTGAACCTGCGGAAGGATCATTGTTGAGATCACATAATAATTGATC
740
```

3rd SeqRecord
```
gi|2765580|emb|Z78455.1|PJZ78455
CGTAACCAGGTTTCCGTAGGTGGACCTTCGGGAGGATCATTTTTGAGATCACATAAAAATTGATC
745
```

# Practice

- **Parsing FASTA Files**
  - Write and Run the following codes on a Colab page called "QuickTour.ipynb":

```
Parse FASTA File

[5]  1  # Download Data File: ls_orchid.fasta
     2  import os
     3  if not os.path.isfile("ls_orchid.fasta"):
     4    os.system("wget https://raw.githubusercontent.com/biopython/biopython/master/Doc/examples/ls_orchid.fasta")

[6]  1  from Bio import SeqIO
     2  for seq_record in SeqIO.parse("ls_orchid.fasta", "fasta"):
     3    print(seq_record.id)
     4    print(seq_record.seq)
     5    print(len(seq_record))
```

**(Solution URL of this Practice)**

# Parsing GenBank Files



- **"ls_orchid.gbk" ( Download Link )**
  1. **.name** → **LOCUS**
  2. **.description** → **DEFINITION**
  3. **.id** → **VERSION**
  4. **.annotations**[‘references’] → **REFERENCE** xN
  5. **.features** → **FEATURES** xN
  6. **.seq** → **ORIGIN**
  7. **.annotations** → All the other information

# Parsing GenBank Files

- **Upload "ls_orchid.gbk" to Colab (by the Local File)**



- **Upload "ls_orchid.gbk" to Colab (by fetching from Internet)**

```
1  # Download Data File: ls_orchid.gbk
2  import os
3  if not os.path.isfile("ls_orchid.gbk"):
4    os.system("wget https://raw.githubusercontent.com/biopython/biopython/master/Doc/examples/ls_orchid.gbk")
```

# Parsing GenBank Files

**1** **Parse the file in GenBank format**

**2**

```
1  from Bio import SeqIO
2  for seq_record in SeqIO.parse("ls_orchid.gbk", "genbank"):
3      print(seq_record.name)  # LOCUS name
4      print(seq_record.description) # DEFINITION line
5      print(seq_record.id)  # VERSION line
6      print(len(seq_record.features)) # FEATURES part
7      print(seq_record.seq) # ORIGIN part
8      print(seq_record.annotations) # All the other misc info
9      print("---------------")
```

**Get one record a time**

**3**

**Print Information of the current record**

**1st SeqRecord**

```
Z78533
C.irapeanum 5.8S rRNA gene and ITS1 and ITS2 DNA
Z78533.1
5
CGTAACAAGGTTTCCGTAGGTGAACCTGCGGAAGGATCATTGATGAGACCGTGGAATAAACGATCGAG
{'molecule_type': 'DNA', 'topology': 'linear', 'data_file_division':
---------------
```

**2nd SeqRecord**

```
Z78532
C.californicum 5.8S rRNA gene and ITS1 and ITS2 DNA
Z78532.1
5
CGTAACAAGGTTTCCGTAGGTGAACCTGCGGAAGGATCATTGTTGAGACAACAGAATATATGATCGAG
{'molecule_type': 'DNA', 'topology': 'linear', 'data_file_division':
---------------
```

# Practice

- **Parsing GenBank Files**
  - Write and Run the following codes on a Colab page called "QuickTour.ipynb":



```
Parse GenBank File

[7]  1  # Download Data File: ls_orchid.gbk
     2  import os
     3  if not os.path.isfile("ls_orchid.gbk"):
     4    os.system("wget https://raw.githubusercontent.com/biopython/biopython/master/Doc/examples/ls_orchid.gbk")

[13] 1  from Bio import SeqIO
     2  for seq_record in SeqIO.parse("ls_orchid.gbk", "genbank"):
     3    print(seq_record.name)   # LOCUS name
     4    print(seq_record.description) # DEFINITION line
     5    print(seq_record.id)   # VERSION line
     6    print(len(seq_record.features)) # FEATURES part
     7    print(seq_record.seq) # ORIGIN part
     8    print(seq_record.annotations) # All the other misc info
     9    print("---------------")
```

**(Solution URL of this Practice)**

# Summary

- **Install Biopython**
  - pip install biopython
- **Import Biopython**
  - import Bio
- **Check the Version of Biopython**
  - print(Bio.__version__)
- **Convert Strings to Bio.Seq.Seq**
  - my_seq = Seq("AGTACACTGGT")

- **Convert Bio.Seq.Seq to Strings**
  - bio_seq = str(my_seq)
- **Complement & Reverse Complement**
  - my_seq.complement()
  - my_seq.reverse_complement()
- **Parse FASTA / GenBank**
  - SeqIO.parse("ls_orchid.fasta", "fasta")
  - SeqIO.parse("ls_orchid.gbk", "genbank")