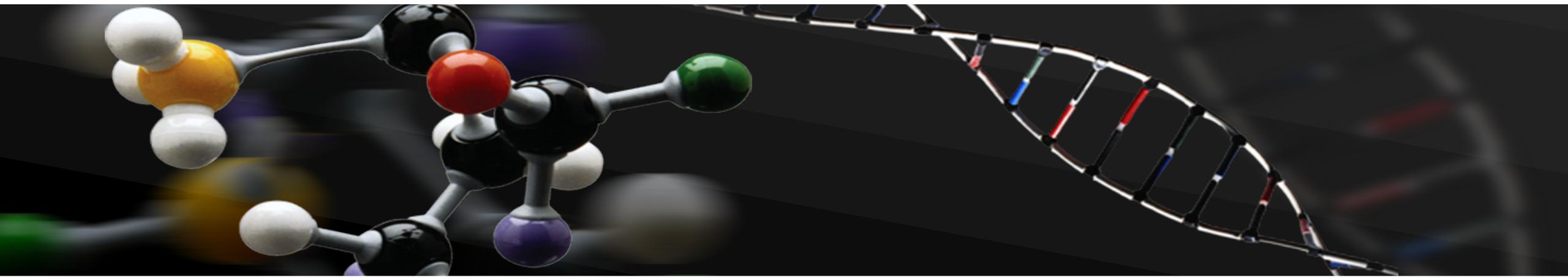


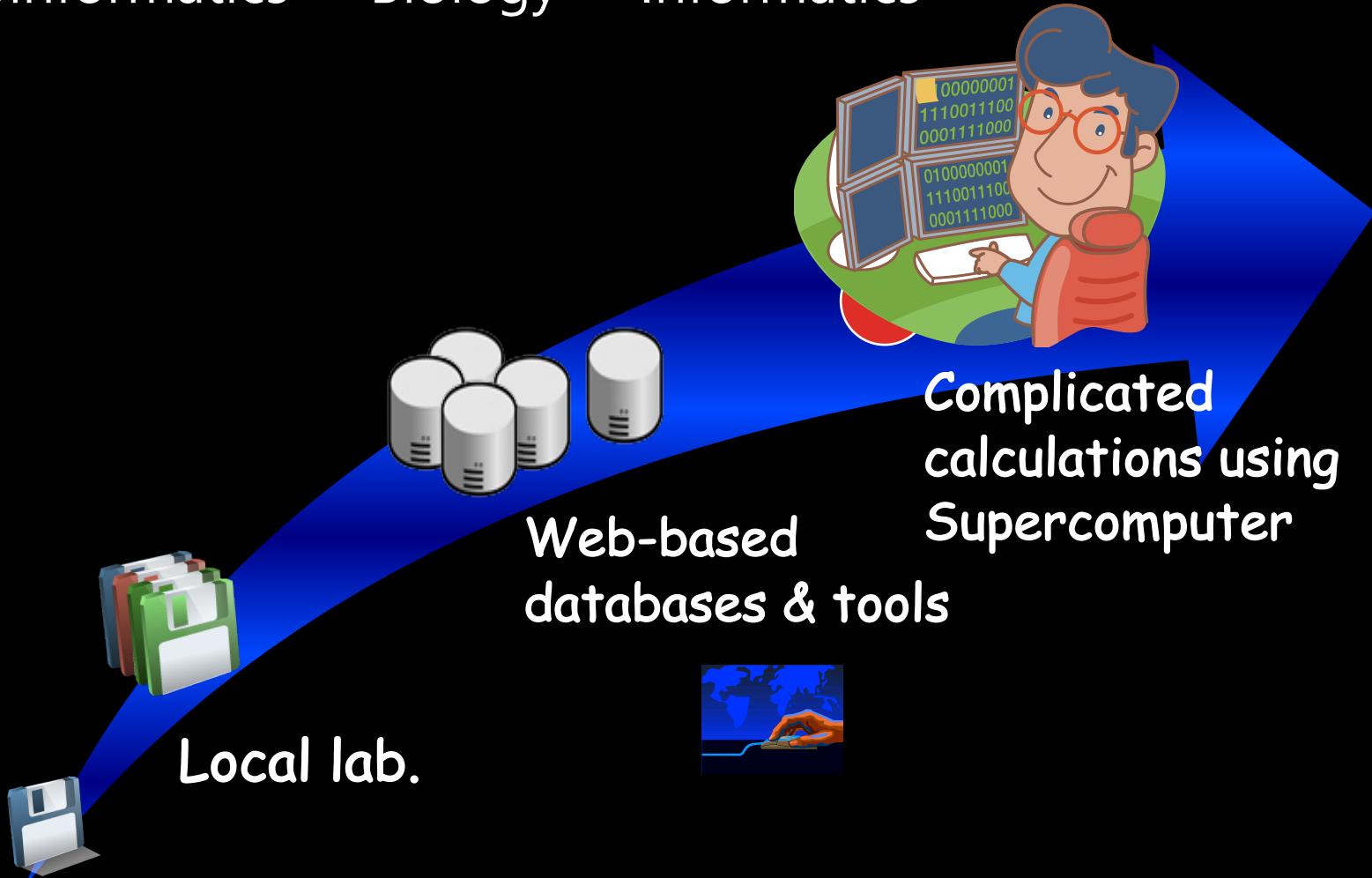
BIOINFORMATICS 2



16th November

WHAT IS BIOINFORMATICS..?

- Bioinformatics = Biology + Informatics



ACCESS TO DATA IN NCBI

- Entrez Nucleotide (<https://www.ncbi.nlm.nih.gov/nucleotide>)
 - ✓ Search GenBank for sequence identifiers and annotations
- BLAST (<https://blast.ncbi.nlm.nih.gov/Blast.cgi>)
 - ✓ Search and align GenBank sequences to a query sequence
- NCBI e-utils (<https://www.ncbi.nlm.nih.gov/books/NBK25501>)
 - ✓ Search, link, and download sequences programmatically
- FTP server (<https://ftp.ncbi.nlm.nih.gov/ncbi-asn1>)
 - ✓ The ASN.1 and flatfile formats are available at NCBI's anonymous FTP server

SARS-CoV-2 Data Hub
[Download](#) [Quick Links](#)
[Betacoronavirus BLAST](#)
[CDC Outbreak Information](#)
[SARS-CoV-2 Articles in PubMed](#)
[SRA Data](#)
[NCBI SARS-CoV-2 Resources](#)
[Datasets command line](#)
[Tabular View](#)
[Dashboard Visualizations](#)
[Mutations in SRA](#) [Complete Tree](#)
[Align](#) [Build Phylogenetic Tree](#)

Selected Results: 0

Refine Results

Reset

Virus

+

Severe acute respiratory syndrome coronavirus 2 (SARS-CoV-2), taxid:2697049

x

Accession

+

Sequence Length

+

Ambiguous Characters

+

Expand Table

	Nucleotide (2,195,773)		Protein (12,570,357)		RefSeq Genome (1)	Select Columns
	Accession	Submitters	Release Date	Pangolin	Isolate	Species
<input type="checkbox"/>	NC_045512 RefSeq	Wu,F., et al.	2020-01-13	B	Wuhan-Hu-1	Severe acute respiratory syndrome coronavirus 2 (SARS-CoV-2), taxid:2697049
<input type="checkbox"/>	OL373892	Moates,D., et al.	2021-11-06		AL-UAB-GX1382	Severe acute respiratory syndrome coronavirus 2 (SARS-CoV-2), taxid:2697049
<input type="checkbox"/>	OL373893	Moates,D., et al.	2021-11-06		AL-UAB-GX1383	Severe acute respiratory syndrome coronavirus 2 (SARS-CoV-2), taxid:2697049
<input type="checkbox"/>	OL373894	Moates,D., et al.	2021-11-06		AL-UAB-GX1384	Severe acute respiratory syndrome coronavirus 2 (SARS-CoV-2), taxid:2697049
<input type="checkbox"/>	OL373895	Moates,D., et al.	2021-11-06		AL-UAB-GX1385	Severe acute respiratory syndrome coronavirus 2 (SARS-CoV-2), taxid:2697049

Download Results

Step 1 of 3: Select Data Type

Sequence data (FASTA Format)
 Nucleotide

 Coding Region

 Protein

Accession List
 Nucleotide

 Protein

 Assembly

Current table view result
 CSV format

 XML format

Download Results

Step 2 of 3: Select Records

 Download Selected Records

 Download All Records

[Next](#)
[Back](#)
[Next](#)

ENTREZ NUCLEOTIDE

E-UTILS

- Searching a Database
 - Input: Entrez database (&db); Any Entrez text query (&term)
 - Output: List of UIDs matching the Entrez query
 - Example: *Get the PubMed IDs (PMIDs) for articles about SARS-CoV-2 published in Nature in 2021*
 - [https://eutils.ncbi.nlm.nih.gov/entrez/eutils/esearch.fcgi?db=pubmed&term=nature\[journal\]+SARS-CoV-2+2021\[pdat\]](https://eutils.ncbi.nlm.nih.gov/entrez/eutils/esearch.fcgi?db=pubmed&term=nature[journal]+SARS-CoV-2+2021[pdat])

Popular Resources

PubMed

Bookshelf

PubMed Central

BLAST

Nucleotide

Genome

SNP

Gene

Protein

PubChem

STEP 2: PARSING ANNOTATION DATA

0 10 20 30 40 50 60 70 80 90 100 110 120
114 TGACAAAG ↓
115 >gb|CY039912:26-1522| /Human/NP/H1N1/USA/1991/// Influenza A Virus (A/Maryland/12/1991(H1N1)) segment 5, complete sequence ↓
116 ATGGCGTCTCAGGCCACCAAACGATCATATGAACAA TGGAGACTGGTGGGAACGCCAGGATGCCACAG ↓
117 AAATCAGAGCATCTGTGGAAAGAATGATTGGTGAAT GGAAGATTCTACATCCAAATGTGCACTGAACT ↓
118 CAAACTTAGTGAATTAGGGGACGACTAATTCAAATAACATAACAAATAGAGAGAATGGTGCCTCTGCT ↓
119 TTTGATGAGAGAAGCATAAATACCTAGAACAGACATCCCAGTGTGGGAAGGATCTAAGAAAAGTGGAG ↓
120 ACTCATCCTTTATGACAAAGAAGAATAAG ↓
121 GCGGTCTTAECTCATATCATGATTTGGCAC ↓
122 TTGTCGCACTGGAAATGGATCCAGAATGT ↓
123 GC1C1C1AA1GCAAGGT1CAACACT1CCCAGAAAGG1C1GGGGCC1CAGGTGCTGCAGTGAAAGGAGTTGG ↓
124 AACAAATAGCAATGGAGTTAACAGAATGATCAAACGTGGAATCAAAGACC1AAACTCTGGAGGGGTGAA ↓
125 AATGGCGAAGGACA↓
126 CCCAGAGGGCAATGAI↓
127 TTTCCTGGCACGGTCA↓
128 TATGGGCTTGCAAGTAGCAAGTGGGATGACTTGAAAGAGAAGGGATATTCACTGGTGGGATAGACCCCT ↓
129 TCAAATTACTTCAAAACAGTCAGTGTCACTGGCTGATCAGACCAAATGAAAACCCAGCCCAAGAGTC ↓
130 ATTGGTGTGGATGCGCATGGCCTGCTGCATTGAGGATTAAAGAGTATCAAGCTCATAAGAGGGAAAG ↓
131 AAAGTGGTTCCAAGAGGAAAGCTTCCACAAGAGGGGTTCAAGAT ↓
132 TGGACTCTAGTACCCCTAGAAACTAAGAACGAGATACTGGCCATA ↓
133 TCAACAGAAGGCATCCGGGGCAGATCAGTGTGCAACTCAT ↓
134 GAAAGAGCAACCGTTATGGCAGCTTCAGCGGGAAACAATGAGGG ↓
135 TTATAAGGATGATGGAAAGCGCAAAGCCAGAAGATTGCTCTCCAGGGGCGGGGAGTCTCGAGCTCTC ↓
136 GGACGAAAAGGCAACGAAACCCGATGTCGCTCTTGCACATGAGTAATGAAGGGTCTTATTCTCGGA ↓
137 GACATGCAAGAGGAGTATGACAGT ↓
138 >gb|CY039914:7-2157| /Human/PA/H1N1/USA/1991/// Influenza A Virus (A/Maryland/12/1991(H1N1)) segment 3, complete sequence ↓
139 ATGGAAAGACTCTGTACCGACAGTGCCTCAACCCGATGA TGTGAACTTGCAGAAAAGACAAATGAAAGAAT ↓
140 ACGGGGAGAAACCAAAAGATCGAAACAAACAAATTGCA CGATATGCACACATATGGAAGTATGCTTCAT ↓
141 GTATTCAGACTTCACTTCATCAATGAGCGGGGTGAATCAATAATCATAGAGCCTGGTACTCTAATGCA ↓
142 CTGCTGAAACACAGATTGAAATAATTGAAGGAAGGGATCGGAATATGGCATGGACGGTGGTAAACAGTA ↓
143 TTGGGATTCCTTCGTCAGTCGAGAGAGGGCGAAGAGACAATTGAAGA AAGATTGAAATCAGAGGGACG ↓
144 ACTACUTGGGAGAAGGCCAATAAGATAAAGTCT ↓
145 AGAAAATGCCAACAAAAGCCGACTACACACTAG ↓
146 ACTATAAGACAAGAGATGGCTAGCAGAGGGTCT ↓
147 TTGGGATTCCTTCGTCAGTCGAGAGAGGGCGAAGAGACAATTGAAGA AAGATTGAAATCAGAGGGACG ↓
148 ATGCGAAAGCTTGCTGAI ↓
149 ATGGGATTGAGGCCAATG ↓
150 CGAGCCTTCTTGAAGAC ↓
151 AAATTTCTTGTGGATGCTTAAAGCAATTGAGGA1CCAAGCC1GAAGGGAGAAGGAATACCGC ↓
152 TTATGATGCGGTAGTCATGAAGACGTTCTTGGGGTGAAGGAAGGAACTTACCATCATTAAAGCCACATGA ↓
153 AAAAGGGATAATTCAAAATTATCTTGGCATGGAAAGCAAGTGTGGCAGAAAATACAGGACTTTGAGGAT ↓
154 GWGAAAAAAATTCCGAGGGTTAAAATATGAAAAAAACAAAGTC ↓
155 TGCCCCAGAAAAGGTGGACTTTGATGATTGCAAAATGTGAG ↓
156 ACCAGAACTTAGATCGCTGCAAGTGGATACAAAATGAATC ↓

>gb|CY039912:26-1522| /Human/NP/H1N1/USA/1991///
/Human/NP/H1N1/USA/1991///
Influenza A virus (A/Maryland/12/1991(H1N1)) segment 5, complete sequence

>gb|CY039914:7-2157| /Human/PA/H1N1/USA/1991///
/Human/PA/H1N1/USA/1991///
Influenza A virus (A/Maryland/12/1991(H1N1)) segment 3, complete sequence

BEFORE

0 10 20 30 40 50 60 70 80 90 100 110 120

114 TGACAAG ↴
115 >gb|CY039912:26-1522| /Human/NP/H1N1/USA/1991/// Influenza A Virus (A/Maryland/12/1991(H1N1)) segment 5, complete sequence ↴
116 ATGGCGTCICAAGGCACCAACGAACTATGAACAAAAGGAGACTGGTGGGGAACGCCAGGAACGCCACAG ↴
117 AAATCAGAGCATCTGCGGAAGAATGATTGGTGGAACTGGAAAGATTCTACATCCAAATGTGCACTGAAC ↴
118 CAAACTTAGTGAATTGAGGGACGACTAATTCAAATAGCATACAATAGAGAGAAATGGTGTCTCTGCT ↴
119 TTTGATGAGAGAAGGAATAAAATACCTAGAAGAGCATTCCCAGTGCTGGGAAGGATCTAAGAAAATGGAG ↴
120 GACCCATATAGAAGAGTAGACGGAAAGTGGATGAGAGAACCTCATCCTTATGACAAAGAAGAAATAAG ↴
121 GAGAGTTGGCGCAAGCAAACATGGTGAAGATGCAACACGCCGTCTTAACATCATGATTGCGCAC ↴
122 TCCAATCTGAACGATGCCACCTATCAGAGAACAGAGCGCTTGTGCACTGGAATGGATCCCAGAAATGT ↴
123 GCTCTCTAATGCAAGGTTAACACTTCCAGAGGCTGGGGCCGAGGTGCTCAGTGAAGAAGGGAGTTGG ↴
124 AACAAATAGCAATGGAGTTAACAGAATGATCAAACGTTGAATGACCGAAACTCTGGAGGGGTGAA ↴
125 AATGGGCGAAGGACAAGGATTGCAATATGAAAGAATGTCATAATTCTCAAGGAAAGTTCAGACAGCTG ↴
126 CCCAGAGGGCAATGATGGATCAAGTGAGAGAACAGTGGAAACCCAGGAAATGCTGAAATTGAAGATCTCAT ↴
127 TTTCCTGGCACGGTCACTTAACTCAAGGGGTCACTTGCAACATAAGTCTTGCTGCCCTGCTTGTG ↴
128 TATGGGTTGCAGTAGCAAGTGGCATGACTTGTAAAGAGAAGGATATTCACTGGTCGGGATAGACCCCT ↴
129 TCAAATTACTTCAAAACAGTCAGTGTTCAGCTGATCAGACCAAATGAAAACCCAGCCCACAAGAGTC ↴
130 ATTGGTGTGGATGGCATGCCACTCTGCTGCATTGAGGATTAAAGAGTATCAAGCTCATAAGAGGGAAAG ↴
131 AAAGTGGTCCAAGAGGAAAGCTTCCACAAGAGGGGTTCAAGTGTCTCAAATGAGAATGTTGAAGGCTA ↴
132 TGGACTCTAGTACCCAGAACTAAGAACGAGATACTGGGCATAAGGACCAGAACGGGAGAAATACCAA ↴
133 TCAACAGAAGGCATCCGGGGCCAGTCAGTGTGCAACCTACATTCTCAGTGCAACGGAATCTCCCTTT ↴
134 GAAAGAGCAACCGTTATGGCAGCTTCAGCGGAACAATGAGGGACGGACATCAGACATGCGAACGGAAG ↴
135 TTATAAGGATGATGGAAAGCGCAAAGCCAGAAAGGATTTGTCTTCCAGGGCGGGGAATCTCGAGCTCTC ↴
136 GGACGAAAAGGCAACGAACCGATGTCGCTTCCCTTGACATGAGTAATGAAGGGTCTTATTCTCGGA ↴
137 GACAATGGCAAGGGAGTATGACAGT ↴
138 >gb|CY039914:7-2157| /Human/PA/H1N1/USA/1991/// Influenza A Virus (A/Maryland/12/1991(H1N1)) segment 3, complete sequence ↴
139 ATGGAAGACTTCGTACGACAGTGTCAACCCGATGATTGGTAACCTGCAAGAAAAGACAATGAAAAGAAT ↴
140 ACGGGGAGAACCCAAAGATGCAAAACAAACAAATTGCACTGCGATATGCACACATATGAAAGTATGCTTCAT ↴
141 GTATTCAAGACTTCACTCAATGACCGGGGTGAATCAATAATCATAGAGCCTGTTGACTCTAATGCA ↴
142 CTGCTGAAACACAGATTGAAATAATTGAAAGGAAGGGATCGGAATATGGCATGGACGGTGGTAAACAGTA ↴
143 TTTGTAACACTACAGGGGTTGGGAAACCAAGGTATCTCCAGATCTATATGACTACAAAGAAAATAGATT ↴
144 CATTGAGATTGGTGTGACAAGAACAGAACAGTCCATATATACTACCTGGAGAACGCAATAAGATAAAAGTCT ↴
145 GAAGGTACGCACATCCATTTCATTTCACTTACAGGAGAACGAAATGGCAACAAAAGCCACTACACACTAG ↴
146 ATGAAGAAAAGTAGGCCAGAACAAAAACAGGTATTACTATAAGACAAGAGATGGCTAGCAGAGGTCT ↴
147 TTGGGATTCTTCTCGTCAAGTCTCCACCAAACCTTCTCAAGTTGACAACTTTAGAGCCTATGTAG ↴
148 ATGCGAAAGCTTGTGATCAAAGTCTCCACCAAACCTTCTCAAGTTGACAACTTTAGAGCCTATGTAG ↴
149 ATGGATTGAGCCGAAATGGCTACATTGAGGGCAAACCTTCCAAATGTCAGAGAATGCTAGAAT ↴
150 CGAGCCTTCTGAAAGACAACACACCGACCACTTAGACTGCCAAGTGGGCCCTCCCTGTTTCAAAGGTCC ↴
151 AAATTCTTTGATGGATGCTAAATTAAGCATTGAGGATCCAAGCCATGAAGGAGAACGGAATACCGC ↴
152 TTATGATGCGGTCAAGTGCATGAAGACGTTCTGGGTGAAAGAACCTACCATCATTAAGCCACATGA ↴
153 AAAAGGATAAAATCAAATTATCTTGGCATGGAAGCAAGTGTGGCAGAAATACAGGACTTGGAGGAT ↴
154 GWGAAAAAAATCCGAGGGTAAAAATGAAAAACAGTCCACTAAATGGCACCTTGGTGAAAATA ↴
155 TGGCCCCAGAAAAGGTGGACTTGTGATGATTGCAAAAATGTGAGTGTGATCTGAAGCAATATGATGTGATGA ↴
156 ACCAGAACTTAGATCGCTGCAAGTGGATAACAAATGAATTCAACAAAGCATGTGAACTGACTGATTG ↴

AFTER

0 10 20 30 40 50 60
1 >HA/S67220_HA/Swine/H1N1/USA/1992/A↓
2 AUGAAGGCAAUACCAUUAGUCUUGCUAUACAUUUACAGCCGAAAUGCAGACACACUAUGUAUAGGT
3 >PB2/AJ564805_PB2/Human/H1N1/Fiji/1983/A↓
4 AUGGAAAGAAUAAAAGAGCUAAGGAUCUGAUGUCGCAGUCGCCACUCGCGAGAUACUUACAAAAAC
5 >NA/CY039911_NA/Human/H1N1/USA/1991/A↓
6 AUGAAUACAAACCAAAAGAAUAAUACCAUCGGGACAGUCUGUCUGAUAGUUGGAAUAGUUAGUCUAUU
7 >NP/CY039912_NP/Human/H1N1/USA/1991/A↓
8 AUGGCGUCUCAAGGCACCAACGAUCAUAUGAACAAAUGGAGACUGGUGGGAACGCCAGGAUGGCCAC
9 >PB1/CY039915_PB1/Human/H1N1/USA/1991/A↓
10 AUGGAUGUCAAUCGACUUUACUUUCCUGAAAGUGGCCAGCACAAAUGCUALAAGCACAACGUUUCC
11 >M1/GQ404564_M1/Avian/H1N1/Czech Republic/2008/A↓
12 AUGAGUCUUCUAACCGAGGUUCGAAACGUACGUUCUCUCCAUCACCGUCGGGGCCCCUCAAAGCCGA
13 >M2/GQ404564_M2/Avian/H1N1/Czech Republic/2008/A↓
14 AUGAGUCUUCUAACCGAGGUUCGAAACGCCUACCAGAAACGGAUGGGAGUGCAGAACGCAAUCAAGT
15 >M2/GQ404616_M2/Swine/H1N1/Czech Republic/1992/A↓
16 AUGAGUCUUCUAACCGAGGUUCGAAACGCCUACCAGAAACGAAUGGGAGUGCAGAACGCAAUCAAGT
17 >M2/GQ404581_M2/Swine/H1N1/Czech Republic/1957/A↓
18 AUGAGCCUUCUAACCGAGGUUCGAAACGCCUACCAGAAACGAAUGGGGGUGCAGAACGCAAUCAAGT
19 >M2/GQ404585_M2/Swine/H1N1/Belarus/1965/A↓
20 AUGAGUCUUCUAACCGAGGUUCGAAACGCCUACAGAAACGAAUGGGGGUGCAGAACGCAACGGUCAAGT
21 >M2/GQ404587_M2/Swine/H1N1/USA/1968/A↓
22 AUGAGCCUUCUAACCGAGGUUCGAAACGCCUACAAAAGCGAAUGGGGGUGCAGAACGCAACGCAAUCAAGT

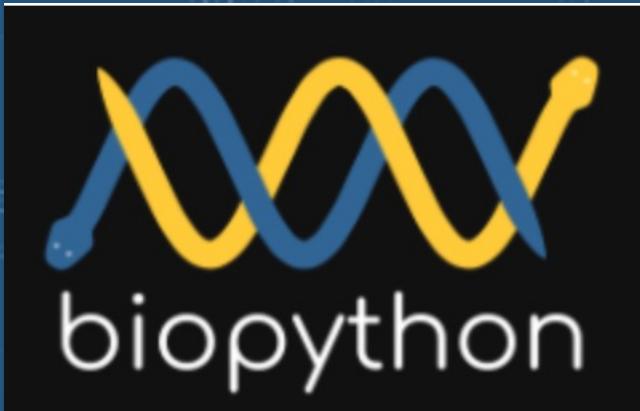
MAKE MYSQL INPUT FILES

```
0   10   20   30   40   50   60   70   80   90   100  110
+-----+-----+-----+-----+-----+-----+-----+-----+-----+
1 CREATE TABLE influ_A_nt_20091207 (←
2 accession varchar (50) NOT NULL default '',←
3 host varchar(20) default NULL,←
4 gene varchar(20) default NULL,←
5 subtype varchar(20) default NULL,←
6 country varchar(20) default NULL,←
7 year int(11) NOT NULL default '0',←
8 len int(11) NOT NULL default '0',←
9 sequence text,←
10 species varchar(5) default NULL,←
11 PRIMARY KEY (accession) ←
12 ) TYPE = MyISAM;←
13 INSERT INTO influ_A_nt_20091207 VALUES ('S67220_HA', 'swine', 'HA', 'H1N1', 'USA', 1992, 1728, 'AUGAAGGCCAUACCAUAGUC
14 INSERT INTO influ_A_nt_20091207 VALUES ('AJ564805_PB2', 'human', 'PB2', 'H1N1', 'Fiji', 1983, 2295, 'AUGGAAAGAAUAAAAAC
15 INSERT INTO influ_A_nt_20091207 VALUES ('CY039911_NA', 'human', 'NA', 'H1N1', 'USA', 1991, 1407, 'AUGAAUACAAACCAAAGAA
16 INSERT INTO influ_A_nt_20091207 VALUES ('CY039912_NP', 'human', 'NP', 'H1N1', 'USA', 1991, 1494, 'AUGGCGUCUCAAGGCACC
17 INSERT INTO influ_A_nt_20091207 VALUES ('CY039915_PB1', 'human', 'PB1', 'H1N1', 'USA', 1991, 2271, 'AUGGAUGUCAUCCGAC
18 INSERT INTO influ_A_nt_20091207 VALUES ('GQ404564_M1', 'avian', 'M1', 'H1N1', 'Czech Republic', 2008, 756, 'AUGAGUCU
19 INSERT INTO influ_A_nt_20091207 VALUES ('GQ404564_M2', 'avian', 'M2', 'H1N1', 'Czech Republic', 2008, 291, 'AUGAGUCU
20 INSERT INTO influ_A_nt_20091207 VALUES ('GQ404616_M2', 'swine', 'M2', 'H1N1', 'Czech Republic', 1992, 291, 'AUGAGUCU
21 INSERT INTO influ_A_nt_20091207 VALUES ('GQ404581_M2', 'swine', 'M2', 'H1N1', 'Czech Republic', 1957, 288, 'AUGAGCCU
22 INSERT INTO influ_A_nt_20091207 VALUES ('GQ404585_M2', 'swine', 'M2', 'H1N1', 'Belarus', 1965, 291, 'AUGAGUCUU
23 INSERT INTO influ_A_nt_20091207 VALUES ('GQ404587_M2', 'swine', 'M2', 'H1N1', 'USA', 1968, 291, 'AUGAGCCUUCUAACCGAGGU
24 INSERT INTO influ_A_nt_20091207 VALUES ('GQ404620_M1', 'swine', 'M1', 'H1N1', 'Germany', 1993, 756, 'AUGAGUCUUCUAACCO
25 INSERT INTO influ_A_nt_20091207 VALUES ('CY021709_HA', 'human', 'HA', 'H1N1', 'USA', 1945, 1701, 'AUGAAGGCCAAGACUACUGO
26 INSERT INTO influ_A_nt_20091207 VALUES ('CY021710_M1', 'human', 'M1', 'H1N1', 'USA', 1945, 759, 'AUGAGUCUU
27 INSERT INTO influ_A_nt_20091207 VALUES ('CY021710_M2', 'human', 'M2', 'H1N1', 'USA', 1945, 294, 'AUGAGUCUU
28 INSERT INTO influ_A_nt_20091207 VALUES ('CY021711_NA', 'human', 'NA', 'H1N1', 'USA', 1945, 1410, 'AUGAAUCCAAAUCAGAAA
29 INSERT INTO influ_A_nt_20091207 VALUES ('CY021712_NP', 'human', 'NP', 'H1N1', 'USA', 1945, 1497, 'AUGGCGUCCCCAAGGCACC
```

Practice #1

Extract sequence
using Entrez from Biopython

What is Biopython..?



- ✓ **Biopython** is a python library to make the **bioinformatics** works easier.
- ✓ Using this, you can import necessary information from bioinformatics file formats (FASTA, FASTAQ, BAM, VCF..) and utilize the tools used in bioinformatics.

Practice #1: Call Entrez in Python

1. Step 1: import Entrez

- from Bio import Entrez

2. Step 2: enter your e-mail

- Entrez.email = "your_email_address"

Practice #1: Call ID list using `esearch`

3. Step 3: Call esearch to find IDs

- `keyWord = "avian influenza"`
- `minDate='2021/11/01'`
- `maxDate='2021/11/01'`
- `search_results = Entrez.esearch(db = "nucleotide", term = keyWord, mindate=minDate , maxdate=maxDate)`

Practice #1: Get ID list using esearch

3. Step 4: get a list of IDs out of esearch

- `records = Entrez.read(search_results)`
- `print(records.keys())`

- `identifiers = records['IdList']`
- `count = int(records["Count"])`
- `print("* Found %i results" % count)`

Count	IdList
RetMax	TranslationSet
RetStart	TranslationStack
QueryKey	QueryTranslation
WebEnv	

Practice #1: Get Sequence using efetch

4. Step 5: use efetch to retrieve entries

- `fetch_handle = Entrez.efetch(db="nucleotide", id = identifiers, rettype="fasta", retmode="text", retmax=str(count))`
- `result_text = fetch_handle.read()`

```
>OK205887.1 Influenza A virus (A/chicken/Veracruz/28159-398/1995(H5N2)) segment 8  
nuclear export protein (NEP) and nonstructural protein 1 (NS1) genes, complete cds  
AGCAAAGCAGGGTGACAAATACATAATGGACTCCAACGATAACCTCGTTCAAGTAGATTGCTATCT  
ATGGCACATACGAAAGCTACTCAGCATGAGAGACATGTGATGTCCTTGTGACAGACTCAGAAGA  
GACCAAAAGGCATTGAAGGAAGAGGCAGCACACTGGACTCGACCTCGAGCCACATAAGAAGGCA  
AAAAGATTGTTGAAGACATCTAAAGACTGAGACGGATGAATTCTCAAATTGCAATTGCACTCCAGCC  
TGCTCCTCGGTATATTACCGATATGAGCATAGAGGAATAAGCAGGGATGGTACATGCTCATGCCAAGG  
CAGAAAATAACAGGAGGCCGTAGATAAATGGATCAGGCCATCATGGACAAGAAATAACTCAAGG  
CAAACCTCTGTCCTATTGATCACTGGAACATTAGTCTACTGAGGGCTTCACAGACGATGGGC  
CATTTGAGCTGAATATCTCCCTTCTCTATCGCAGGACATTCTACAGAGGTGTCAAATTGCAATT  
GGAATCTCATCGGTGGACTGGAATGATAACTCAATTGAGCTCTGAAAATATAAGAGATTG  
CTTGGGGAGTCCGTGATGAGATAAGGGGACCTTCACTCCCTCAAAGCAGAAACGCTACATGGCGAGAAG  
AATTGAGTCAAAGTTGAAGAGATCAGATGGCTAATTGAGCTGAGAGTGTAGAAACATATTAAACCAAACTGA  
GAACAGCTCGAGCAGATAACGTTCTGCAGGCATTGCAACTCTTGAAGTTGAGAGTGTGAGATAAGG  
ACATTTCCTTCAGCTTATTAGTACTAAAAAACCCCTGTTCTACT
```

```
>OK205886.1 Influenza A virus (A/chicken/Veracruz/28159-398/1995(H5N2)) segment 7  
matrix protein 2 (M2) and matrix protein 1 (M1) genes, complete cds  
AGCAAAGCAGGTAGATATTGAAAGATGAGTCTCTAACCGAGGTGAAACGTAGCTCTCTATCGTC  
CCGTAGGGCCCCCTCAAGCCGAGATCGCAGAGACTTAAGAGATGTCTTCAGGGAAACACCGATC  
TTGAGGCACTCATGGAAATGGCTAAAGACAAGACCAATCTGTACCTCTGACTAGGGGGATTAGGATT  
TGTGTTACGCTCACCGTGCCAGTGAGCGAGGACTGCAGCTAGACGCTTGTCAAATGCCCTTAAT  
GGGAATGGGGATCCAACACATGGACAGAGGGTCAAAGTGTACAGGAAGCTAAAAGGAAATAACAT  
TCCATGGGGCAAAAGAATGGCACTTGGCAGTACTCGTGTGACTTGCAGTTGCATGGCCTCATATA  
CAACAGAAGGGAAAGAATGGTACCCGAAGTGGCATTGGCTGGCTGCGCCACATGTGAGCAAATTGCT  
GATTCCCAAGCATGGCTCTCACAGAAATGGTACAAGCCACCAAGCCACTGATTAGACATGAAACAGAA  
TGGTACTGGCCAGTACTACGGCAAAGGCATGGAGCAAATGGCAGGGTCAAGTGAACAGCAGCAGAGGC  
TATGGAGGTTGCTAGTCAGGCTAGACAGATGGTCAGGCAATGAGGACATTGGAAACCCATCTAGCTCC  
AGTGTGGCTAAAGATGTCTCTGAAAATTGCAAGGCCTACCAGAAACGGATGGGAGTGTGCAAATAC  
AGCGATTCAAGTGTACTCTCGTTATTGCGCAAGCATCTGGGATCTGCACTTGTGATATTGTTGATTG  
TTGATGCTCTTCTCAAATGCAATTGCTGCTTAAACGGTTGAAAAGAGGGCCTCTACGGGA  
AGGAGTGCCTGAATCTAGGGGAAGAAATCGGCAGGAACAGCAGAGTGTGCTGGATGTTGACGATGGT  
CATTGGTCAACATAGAGCTGGAGTAAAAACTACCTTGTCTACT
```

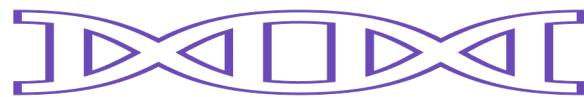
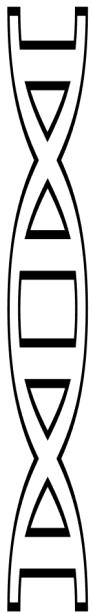
```
>OK205885.1 Influenza A virus (A/chicken/Veracruz/28159-398/1995(H5N2)) segment 6  
neuraminidase (NA) gene, complete cds  
AGCAAAGCAGGGTGAAGAGTGAATCCAACATCGAAAGATAATAACAATTGGCTCGTCTCTAACCATT  
GCAACAGTATGTTCTCATGCAGATTGCCATCTAACACGACTGTGACACTGCATTCAAGCAAACG  
AATGCAGCATCCCCGCTAACAAACCAAGTAGTGCCTGACATGTAACCAATCATATAAGAGAGGAACATAGACTA
```

Practice #2

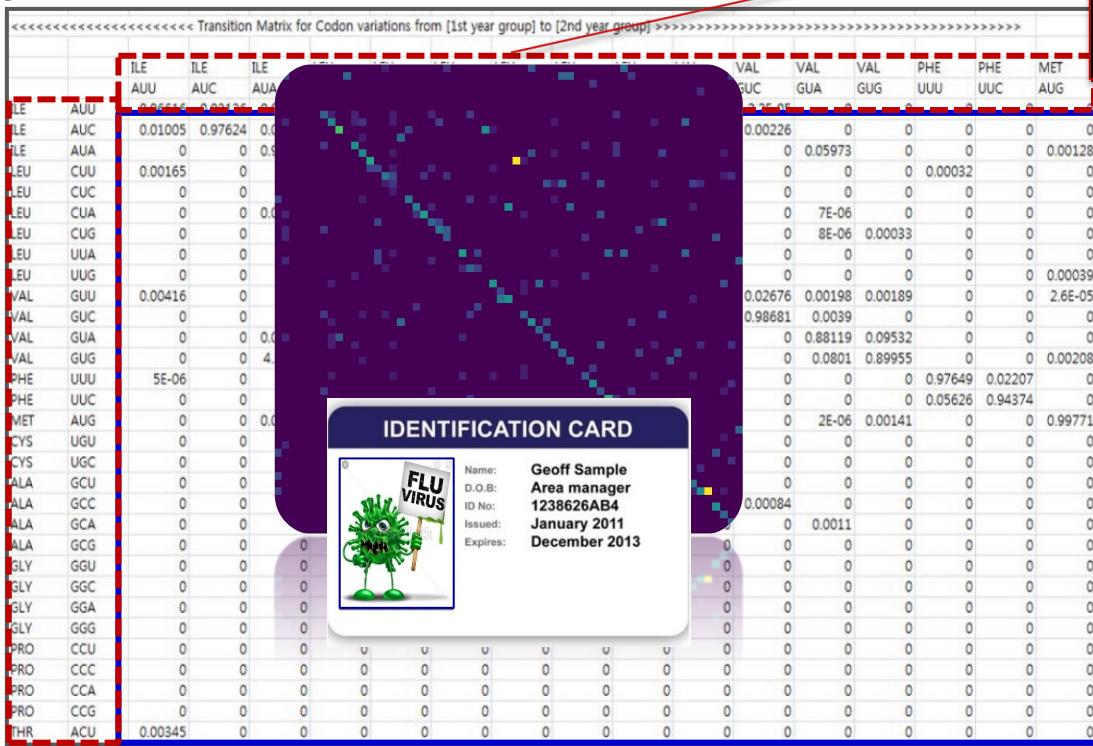
Pilot study using
influenza sequences

Create Variant vector

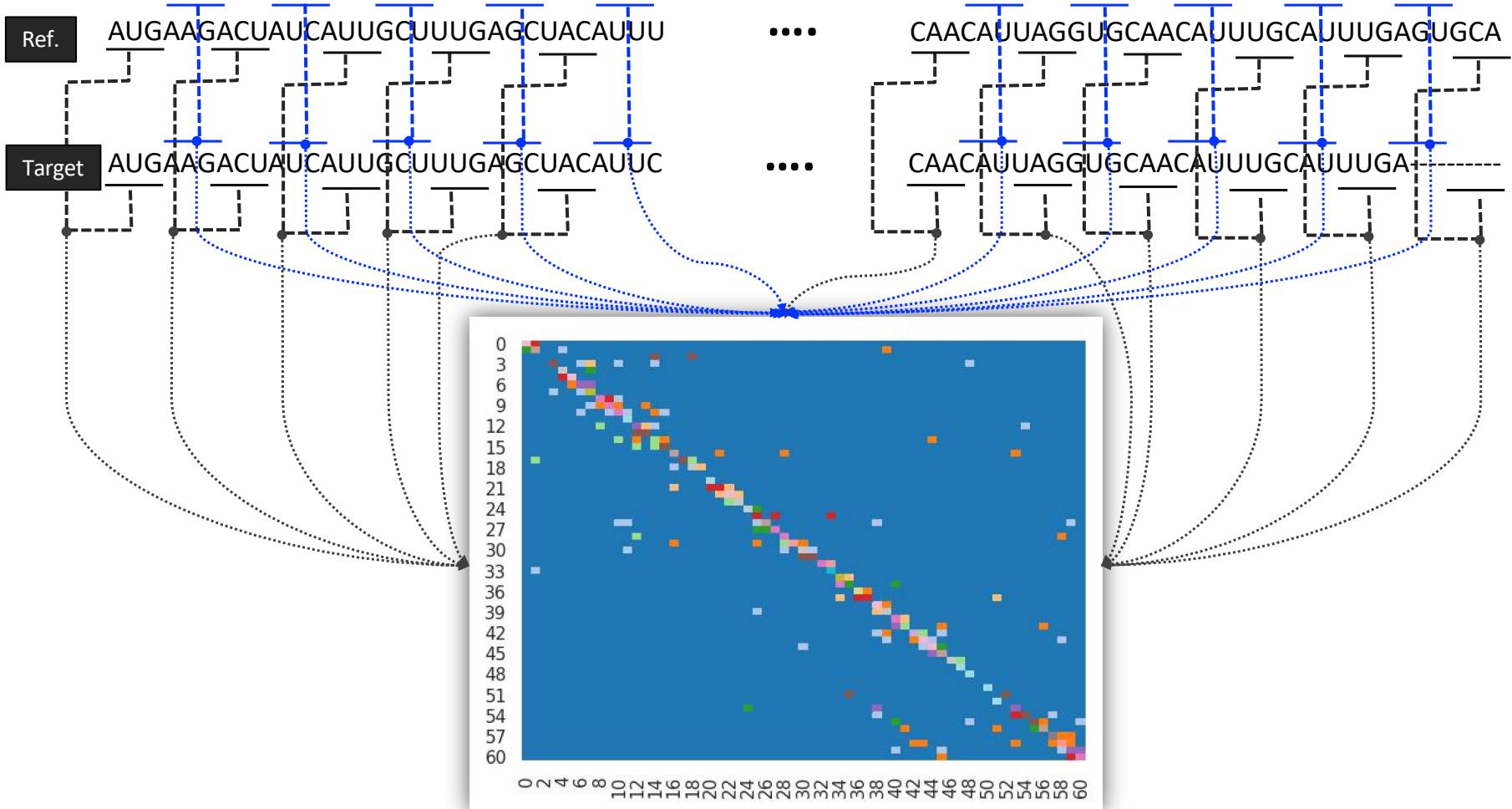
Target sequence



Reference sequence



Create Variant vector



What is synonymous codon..?

		Second base					
		U	C	A	G		
First base	U	UUU Phenylalanine (Phe)	UCU Serine	UAU Tyrosine (Tyr)	UGU Cysteine (Cys)	U Stop Codon	
	C	UUC	UCC	UAC	UGC	UGA Stop Codon	C
	C	UUA Leucine (Leu)	UCA	UAA	UGA	UGG Tryptophan (Trp)	A
	C	UUG	UCG	UAG			G
	C	CUU Leucine (Leu)	CCU Proline	CAU Histidine (His)	CGU Arginine		U
First base	C	CUC	CCC	CAC	CGC		C
	C	CUA	CCA	CAA	CGA		A
	C	CUG	CCG	CAG	CGG		G
	A	AUU Isoleucine (Ile)	ACU Threonine	AAU Asparagine	AGU Serine		U
	A	AUC	ACC	AAC	AGC Serine		C
First base	A	AUA	ACA	AAA	AGA Arginine		A
	A	AUG Methionine (Met) Start codon	ACG	AAG Lysine	AGG Arginine		G
	G	GUU Valine	GCU Alanine	GAU Aspartic acid	GGU Glycine		U
	G	GUC	GCC	GAC	GGC Glycine		C
	G	GUA (Val)	GCA	GAA Glutamic acid	GGA Gly		A
First base	G	GUG	GCG	GAG Glu	GGG G		G
	G						

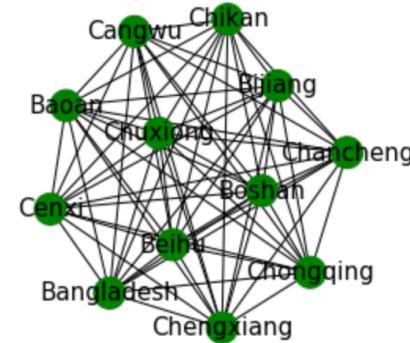
https://www.researchgate.net/figure/Synonymous-Codons-of-20-Amino-Acids_fig2_324469014

Tracking the spread of influenza between countries through cluster analysis

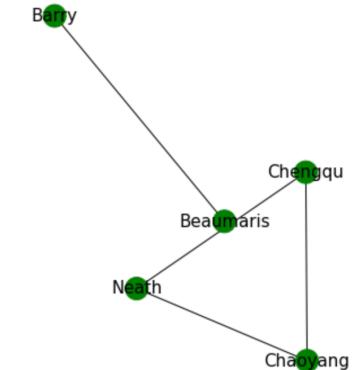
H3N2 2019 (Correlation Coefficient > 0.9)



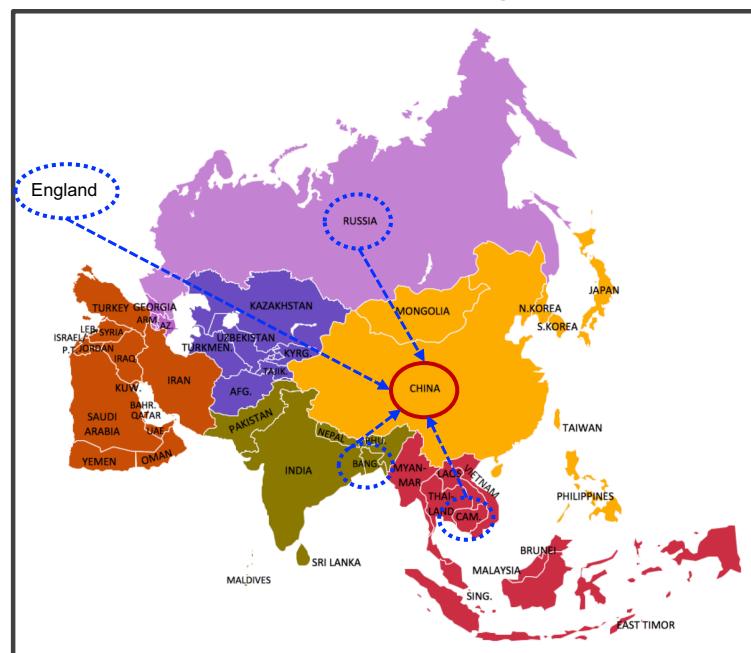
China & Russia & Cambodia



China & Bangladesh



China & England

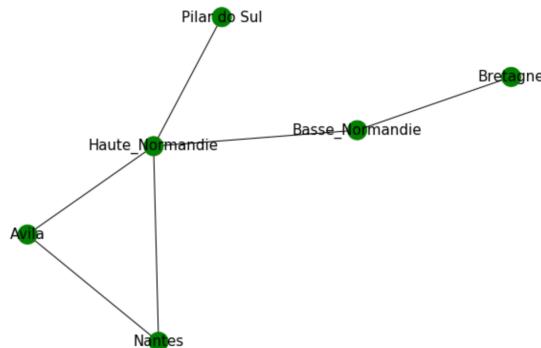


Tracking the spread of influenza between countries through cluster analysis

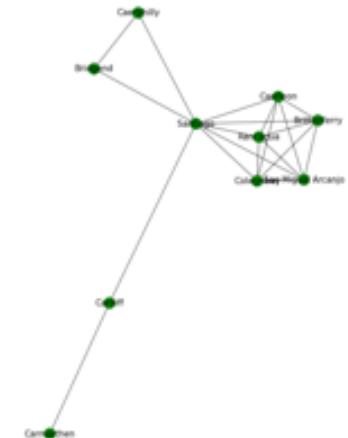
H3N2 2019 (Correlation Coefficient > 0.9)



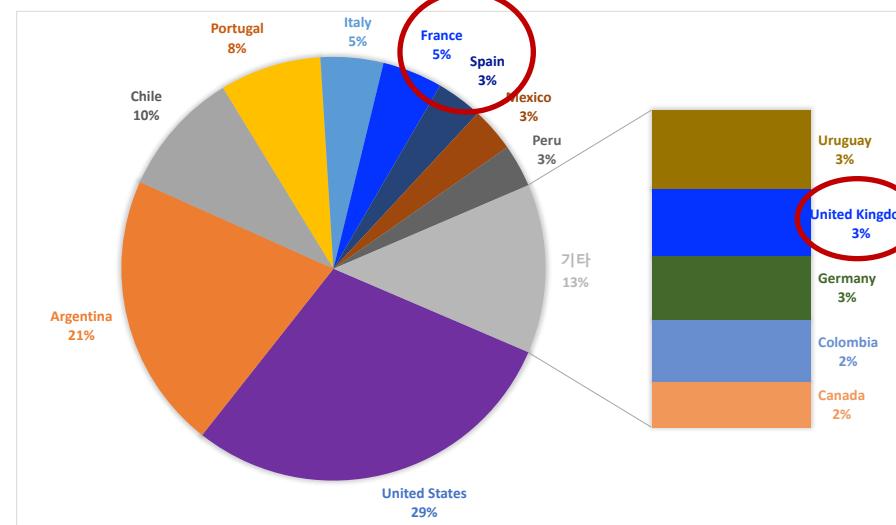
Brazil & France & England



Brazil & Spain & France



Southern America & Western Europe





Assignments





Assignments

1. **Install** biopython library in python and try crawling various data including sequence data from NCBI.
2. **Design** a pilot study related to COVID-19 to be conducted using the big data and technologies learned during the remainder of the semester.
3. **Armanç** from KIST, **Chen Jingyu** from KRICT, please prepare a 5-minute presentation next week.

THANK YOU

