

Ramapo College

Virus Protein Interactions in Plants Database

Final Report

Ritwik Katiyar

Advanced Bioinformatics

Dr. Stuart

10<sup>th</sup> May 2019

## Abstract

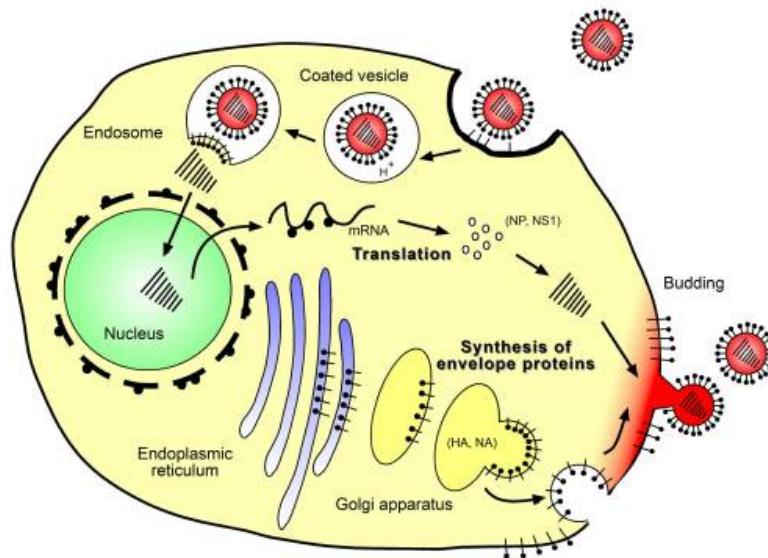
Viral proteins and viral Si-RNA's (Small interfering RNA's) are important to understand in order to come up with ways to counteract viruses. Hence, a database which could record viral proteins along with the type of host the virus infects could prove beneficial in understanding what the viruses are capable of and thereby giving a step ahead against counteracting viruses and protecting agricultural crops for viral diseases.

## Introduction & Background

In 1892 a Russian botanist, Dimitrii Ivanovsky conducts and experiment showing that a pathogen other than bacteria was infecting the tobacco plants and that the organism was able to pass through porcelain filters that would normally retain bacteria. Then in 1898 Marcus Beijerinck made similar observations and concludes that the pathogen must be a distinct organism and officially coins the term “Virus” (Lecoq). Since then, the study of viruses has extended into its own subfield of virology. Many different viruses have been indentified, named and sequenced. Various viral databases have also been created that record the potential viral hosts, record viral sequence, or simply record the existence of a specific virus. Since the discovery of viruses we have been able to create numerous vaccinations and nearly eliminate major viral diseases in humans. We have also created numerous vaccinations for domestic animals as well as for wild animals and even for plants; although, plant vaccinations are simply referred to as pesticides. With the use of pesticides we have insured the safety of our crops. However, there is still a lot that we do not understand about viruses. (Smith)

Viruses are simple parasitic orgasms that require a host in order to survive. Since viruses do not possess a proper cellular structure, they require a host to replicate. Viruses come in various shapes and sizes and comprise mostly of either DNA or RNA. The nucleic acid sequence can also be single or double stranded DNA or RNA. Viruses also possess the capability to encode various types of proteins ranging from as little as 3-4 proteins to 100-200 proteins (Lodish). Viruses depend on a relatively simple mechanism to infect their hosts. First a virus breaks and enters through the membrane (or membranes) of the host cell. Since, viruses are usually extremely small and have a simple structure, a single virus only contains limited amount of RNA or DNA, in other words a very limited number of genes. Viruses use these limited genes

to not only code for proteins that can be used to simply allow the virus to replicate within a host cell, but also influence the proteins produced by the host cell (Smith) (Lodish). This allows the virus to hijack the host cell and replicate. Eventually replicating enough to burst out of the cell and infect other neighboring cells.



**Figure 1: Replication cycle of influenza-A virus.** Binding and entry of the virus, fusion with endosomal membrane and release of viral RNA, replication within the nucleus, synthesis of structural and envelope proteins, budding and release of virions capable of infecting neighboring epithelial cells. (Kamps)

Obviously many plants, animals and even some single celled organisms possess a natural anti viral defense systems. However, many viruses also come equipped with a counter host anti viral defense mechanisms.

One of the most important strategies of plants against viral infections is known as si-RNA- mediated (small interfering RNA mediated) gene silencing. Si-RNA-mediated gene silencing allows for defensive signal to spread to other neighboring cells otherwise the mechanism itself is similar to the mi-RNA (micro-RNA) mediated gene silencing. Plant miRNAs play two main functions in the antiviral defense; mi-RNAs target the viral RNA/DNA and

prevent the virus from reproducing, and miRNAs trigger the synthesis/biogenesis of si-RNA (which is responsible for the plant cells primary antiviral response). However, viruses have a counter defense system against the plant cells known as viral suppressors of RNA silencing (VSRs), which interfere with host RNA silencing in multiple ways. (However, VSRs additional functionality includes assistance in viral replication, encapsidation and movement.) VSRs primarily use two methods in dealing with plants defense system; first being the suppression and the assembly of AGOs (Argonautes; protein family that plays vital role in RNA silencing processes) into RISCs (RNA-induced silencing complex; meant to silence the viral DNA/RNA) and the second being interacting with AGOs (Argonautes) to degrade the protein. The figure below shows the two viral counters to the plant anti viral mechanisms in action. (Rui) (Kamthan)

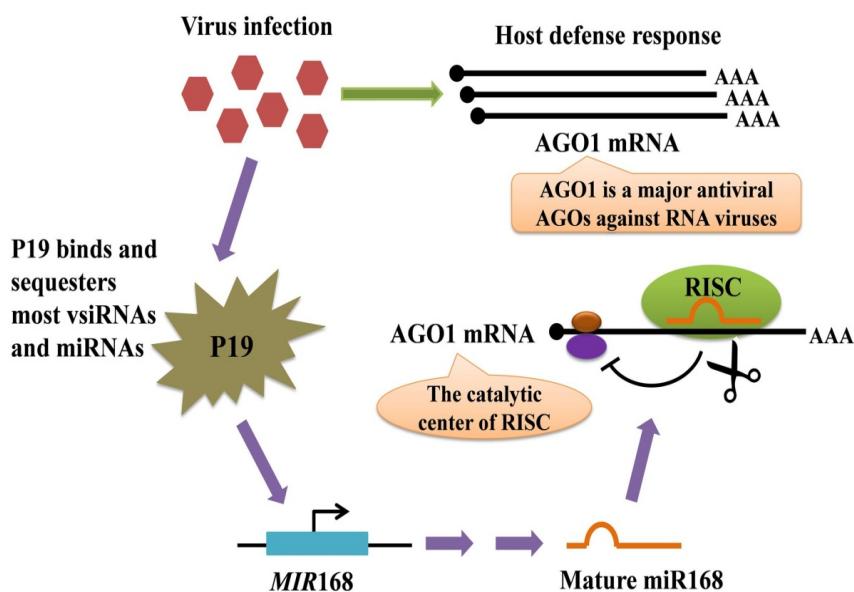


Figure 2: represents a model for the regulation of AGO mRNA level mediated by protein P19-induced miR168. The AGO mRNAs are made when a virus infection is detected within the cell. While the plant cell is creating mRNAs for AGO the virus produces P19 VSR which binds to the virus- encoded siRNAs and with host cell miRNAs which prevents the miRNA loading into AGO. The P19 protein however, does not bind to MIR168 micro RNA (MIR168, represses the AGO1 mRNA); resulting in low production of AGO1 mRNA. (Rui) (Pratt)

Based on the information, it is clear that viral proteins and viral si-RNA's are important to understand in order to come up with ways to counteract viruses, and there by protect

agricultural crops form viral diseases. Hence, a database which could record viral proteins along with the type of host the virus infects could prove beneficial in understanding what the viruses are capable of and thereby giving a step ahead against countering viruses and protecting agricultural crops; and thereby preventing massive crop loss due to viral outbreaks.

## Requirements & Materials Used

Requirements & Materials	Summary
Visual Studio Code	Used to write all the code
Python and Bio-python	Python was used as a base to bio-python and postgresql. Bio-python was used to grab viral sequence information and viral IDs
Flask, HTML and CSS	To visualize all the data within the data and create an app to showcase all the data
Postgresql (Elephant SQL)	Database where the data was stored
R, R-studio and Taxize	Used to grab scientific names and common names of plants
Beautiful Soup	Used to grab information from Plant viruses online website
NCBI and Conserved Domains Database	NCBI was used from bio-python and CDD was used to search for various motifs and domains within a sequence of translated CDS's
Microsoft Excel	To better organize the raw data gathered from various sources
Matplotlib and Pandas	To create the graph on the home page of the website.

Table 1: Shows a comprehensive list of requirements for the following methods of setting up the database

## Method

1. The first step to creating any such database is to develop a schema for the database.

---

```

Database : Vpipdb
Plants(Scientific_Name, Common_Name, PRIMARY KEY:(Scientific_Name))
Infect(Infect_Unique_ID, Scientific_Name, Virus_name, PRIMARY KEY:(Infect_Unique_ID))
Viruses(Virus_Name, Virus_ID, Sequence, PRIMARY KEY:(Virus_ID))
CDS(Virus_ID, CDS_UNIQUE_ID ,CDS_Location, CDS_translatedseq, PRIMARY KEY:(CDS_UNIQUE_ID))
CDS_Analyze(UNIQUE_ID ,CDS_UNIQUE_ID, Start, End, E-Value, Bit-score, Accession#, Motifs_And_Domain_Name, PRIMARY KEY:(UNIQUE_ID))

```

Figure 3: displays the schema of the proposed database.

2. Next to select the plants for the database. Since there are numerous agricultural crops grown around the world and a lack of data linking a specific virus to a particular crop; it is important to focus on the top grown crops around the world. To do that, the data about the amount of agricultural crops grown around the world was taken from Food and agricultural organization (FAO). This data was then sorted and only the top 10 plants were kept.
3. FAO presents the data with the use of common names hence to get the scientific names of the data a simple R script was used that would grab the genus of the plants based on the common name.

```

#Finds taxonomic information on the following common names within the NCBI database
library("taxize")
set_entrez_key("8f31f1938108f0e82af1c7406f19d44fac08")
Sys.getenv("ENTREZ_KEY")

name <- c('Sugarcane', 'Corn', 'Rice', 'Wheat', 'Cucumbers', 'Tomatoes', 'Carrots', 'Spinach', 'Papaya')
uids <- get_uid(commonname)
uids.found <- as.uid(uids[!is.na(uids)])
commonname.found <- commonname[!is.na(uids)]
sciname <- sci2comm(name, db = 'ncbi')
names(sciname) <- commonname.found
sciname

```

Figure 4: Shows the R-script that finds taxonomic information from common names within NCBI

Saccharum	Sugarcane
Zea	Corn
Oryza	Rice
Glycine	Soybean
Triticum	Wheat
Cucumis	Cucumbers,Melons
Solanum	Tomatoes,Potatos,Eggplants
Brassica	Cabbages,Cauliflowers
Daucus	Carrots
Spinacia	Spinach

Figure 5: The data was then organized using excel

4. With this information in hand a script was written using beautiful soup to gather information from the Plant and viruses online website. The data that was extracted from Plant and viruses online website was the data on plant susceptibility to a specific virus , all the plant species that the data was available for and all the viruses that were recorded by the website.

```

import lxml
from urllib.request import urlopen
from bs4 import BeautifulSoup

url = "http://sdb.im.ac.cn/vide/hostlist.htm"
html = urlopen(url)
soup = BeautifulSoup(html, 'lxml')
type(soup)

all_links = soup.find_all("a")
for link in all_links:
    if link.get("href") != None:
        temp = link.get("href")
        str(temp)
        f = open("Host_names.txt", "a")
        f.write(temp + "\n")
    else:
        f = open("Host_names.txt", "a")
        f.write("None" + "\n")

```

Figure 6: The first script used to get the list of hosts used by the website

```

family035.htm#Abelia grandiflora
family082.htm#Abelmoschus esculentus
family082.htm#Abelmoschus manihot
family078.htm#Abrus precatorius
family082.htm#Abutilon
family082.htm#Abutilon hirtum
family082.htm#Abutilon indicum
family082.htm#Abutilon theophrasti
family044.htm#Acanthospermum hispidum
family002.htm#Acer palmatum
family044.htm#Achillea filipendulina
family044.htm#Achillea ptarmica
family025.htm#Adansonia digitata

```

Figure 7: Shows the output from the script

```

Saccharum officinarum
Zea diploperennis
Zea mays
Zea mays ssp. mays
Zea mays ssp. mexicana
Zea perennis
Oryza australiensis
Oryza barthii
Oryza cubensis
Oryza glaberrima

```

Figure 8: Using Excel the output data was organized

```

import lxml
from urllib.request import urlopen
from bs4 import BeautifulSoup

url = "http://sdb.im.ac.cn/vide/sppindex.htm"
html = urlopen(url)
soup = BeautifulSoup(html, 'lxml')
type(soup)

all_links = soup.find_all("a")
for link in all_links:
    f = open("Virus_names.txt", "a")
    f.write(str(link) + "\n")

```

Figure 9: the second script used to extract all viruses that infect the said host from the website

Once the list of all plant and the viruses used by the Plant Viruses Online website was gathered, the list was then organized and arranged using excel to contain appropriate data. The viruses were also then linked with their host plants; figure 10 shows the final version of what the data looked like.

Saccharum officinarum	Maize dwarf mosaic potyvirus
Saccharum officinarum	Sugarcane bacilliform badnavirus
Saccharum officinarum	Sugarcane Fiji disease fijivirus
Saccharum officinarum	Sugarcane mosaic potyvirus
Saccharum officinarum	Sugarcane streak monogeminivirus
Zea mays	Barley stripe mosaic hordeivirus
Zea mays	Barley yellow dwarf luteovirus
Zea mays	Barley yellow striate mosaic cytorhabdovirus
Zea mays	Bermuda grass etched-line marafivirus
Zea mays	Brome mosaic bromovirus

Figure 10: shows how the data looks like once it's been organized in the tab delimited format. The first column shows the species of plant. The second column shows the virus that infects the plant.

5. With the use of a similar script to that shown in figure 4, the common names were determined for each and every plant.

Saccharum officinarum	Sugar Cane
Zea mays	Corn
Zea mays ssp. Mays	Common Corn-Variations
Zea mays ssp. mexicana	Corn- Subspecies
Oryza australiensis	Australian Rice
Oryza barthii	Wild Rice
Oryza rufipogon	Wild Red Rice
Oryza glaberrima	African Rice
Oryza latifolia	Broadleaf Rice
Oryza longistaminata	Longstamin wild rice
Oryza nivara	Wild Asian Rice
Oryza perennis	Brown Bread Rice
Oryza punctata	Red Rice

Figure 11: The data for the plants table mentioned in the schema in figure 3. Scientific Names on the left and the common names on the right

6. Once the viruses were connected to their respective host plants the NCBI viruses ID's was determined using a script. Then another script was used to grab the sequences for the viruses based on the determined ID.

```
import csv
import Bio
from Bio import Entrez

Entrez.email = "rkatiyar@ramapo.edu"
textfile = open("10.txt", "r")
reader = csv.reader(textfile)
allRows = [row for row in reader]

for i in range(len(allRows)):
    handle = Entrez.esearch(db = "nucleotide", term = allRows[i], retmax = 1)
    record = Entrez.read(handle)
    handle.close()
    idList = record["IdList"]
    f = open("idrecords10.txt", "a")
    f.write(str(idList)+ "\n")
```

Figure 12: Shows the script that extracts a single NCBI ID from the search of the virus

```
#Gets sequences from the NCBI- ID
import csv
import Bio
from Bio import SeqIO
from Bio import Entrez

Entrez.email = "rkatiyar@ramapo.edu"
Entrez.api_key = "8f31f1938108f0e82af1c7406f19d44fac08"

textfile = open("ID_only_1.txt", "r")
reader = csv.reader(textfile)
all_ID = [row for row in reader]

for i in range(len(all_ID)):
    handle = Entrez.efetch(db = "nucleotide", id = all_ID[i], rettype = "gb", retmode = "text")
    record = SeqIO.read(handle, "genbank")
    handle.close()
    Sequence = record.seq
    f = open("seq_list.txt", "a")
    f.write(str(Sequence) + "\n")
```

Figure 13: shows the script used to extract sequence information based on the ID obtained from previous script

7. The resulting data was then arranged using excel

Maize dwarf mosaic potyvirus	['18490052']	AAAAACAACAAGACTCAACACAACAAACACGA
Maize rayado fino marafivirus	['14141972']	GTCGACGTCGCATTCTGCACCAGCTTCGCTCGCCAGAA
Rice tungro bacilliform badnavirus	['18026839']	TGGTATCAGAGGGATGTTGAACTTTAAGGGAAAATAGA
Rice dwarf phytoreovirus	['20428565']	GGCAAAACCTCGCCATGGCTTATCTAACGACGTCAAGAA
Rice ragged stunt oryzavirus	['20428612']	GATAAAATCTCATGACTCTATTAGTGATCACCGAGCAGAC
Rice grassy stunt tenuivirus	['9635243']	ACACAAAGTCTGGCAATTACAAACAAAGAAAAACTTA
Barley stripe mosaic hordeivirus	['19744922']	GTAAAAGAAAAGGAACAAACCTGTTGTTGACGCTA
Glycine mosaic comovirus	['944542958']	TATTTAAATCTTTATAAGATTGATAACCGCAATCATAA
Glycine mottle carmovirus	['216905810']	GGGTAACCCAGCCAGTTATCCACCATTTAACCTTCAGGA
Soybean mosaic potyvirus	['1591440922']	AAATTTAAACTAGTTATAAAAGACAACAAACAAATTAAAC
Barley yellow striate mosaic cytorhabdovirus	['1586082777']	CACGACCAGTGTGCTATAATTGATTATTGGTGTACTCTC
Abutilon mosaic bigeminivirus	['1464307913']	ACCGGATGGCCCGAAATTGGTGTCCAGAACTTTAAT
Melon Ournia ourmaviruses	['194351521']	CCAGATTACGGTATCTTCGACACCGCAAGAGCGAACTT
Heracleum latent trichovirus	['1464311295']	TCCCTCCGATTATGAGGCCTTCGATCGTAGGTCAAGATGA
Tomato spotted wilt tospovirus	['1389531354']	AGAGCAATTGTCAAATTATTCACACCTTAACACTCAG
Peach rosette mosaic nepovirus	['1159189057']	GGAAAAAAACCAAAGTTGTTTCTTGACTGCCCTTGT
Alfalfa mosaic alfamovirus	['1561349898']	ACTATGCTGCCTTGCACAGCTCAACTGCCAAACCTCC
Datura Colombian potyvirus	['1572772023']	CAAACCTGGGTGTGGAATGGCTCACTAAAGCTGAATTG
Beet pseudo-yellows closterovirus	['268529022']	ATAAAATATCCTTAGGGTAAAGAAAGTTTCCCCCCC

Figure 14: shows how the data looked in excel. The first column shows the name of the virus. The second column shows the ID found. The third column shows the complete sequence of the virus.

8. After the sequences were obtained the next step was to obtain information on potential CDS regions for each virus. Not all viruses contain a CDS region and for viruses without a CDS region no further data was necessary to obtain.

```
import csv
import Bio
from Bio import Entrez
from Bio import SeqIO

Entrez.email = "rkatiyar@ramapo.edu"
Entrez.api_key = "8f31f1938108f0e82af1c7406f19d44fac08"

textfile = open("Virus_ID.txt", "r")
reader = csv.reader(textfile)
all_ID = [row for row in reader]

for i in range(len(all_ID)):
    with Entrez.efetch(db="nucleotide", rettype="gb", retmode="text", id= all_ID[i]) as handle:
        for seq_record in SeqIO.parse(handle, "gb"):
            f = open("CDS_loc.txt", "a")
            f.write(str(all_ID[i]) + "\n")
            for feature in seq_record.features:
                if feature.type == 'CDS':
                    f.write(str(feature.location) + "\n")
```

Figure 15: Shows the script for getting the CDS region for each virus

18490052 CDS_01	139	9265 MAGTWTHVTHWQPNLDNPRDVRIMELFAAKGQVYDEKRALEHNSKLRRQAQVVDEPMITVQPKCAQIV
18490052 CDS_02	2682	2922 NLCRSVESSVDRIIIVWKILRNMACVQGQEILQAFNPAKKRFRRLCQYISASNIKFSAEKQSQQLYFNQTPPRI
14141972 CDS_03	96	6180 MSSFLRGHLLSGVESLPTTHTRDITAPIESLATPLRSLERYPSIPKEFHSFLHTCGVDISGGFHAAHPHPVHK
14141972 CDS_04	301	1561 MPLTPTPSIRPSRPTSFSMSGPTTLGVRQTCSRSSPSSPDSPSTSSTGSCPCTPPGTPLRTSRTARPSS
18026839 CDS_05	67	667 VLKRNLTSQNIESRYEKLFLDLAVWGKEKKQYLLSTDNI5FYCYDTSKTSESERKHTFHSDNKQLNSIVDLIJKHSI
18026839 CDS_06	663	996 MSADYPTEKALEKFKNLESDTAAKDKFNWWFTLENKTTADVNLASKGLVQLYALQEIDKKINNLTAQVKLQPSI
18026839 CDS_07	992	6026 MSLRPFTGTSRTITQDSTSSENKKGKNSTKRELIEEVNDQDVNFWDWKLSGKPKNLKQSI
18026839 CDS_08	6046	7216 MNIEPYSHIIDKKNKVPIYDQGNLFHTEKSARLHESRGLLDHLTFSSNTTERVRLHLADYLLESERESYKNEW
20428565 CDS_09	14	3365 MAYPNDRVNVWWDVYNFRDPVNREHLIRDNGLTVRNLTNMLTMERDDQQLIAQLSNMMKSLSIGIEKAQ
20428612 CDS_10	11	3779 MTLLVITEQTIHSLCLDHGETNQIAIEKQLEKPELFSYITDAEPLATEGVFVGPDICGCNCITHTRFVPDYVAKPPPYC
20428612 CDS_11	490	1471 MPSVRASDPKQQLYPIAEQLSQEDVRNIAKTHRSPTLRRTQTKFCVGRGELSVDTLAPLQYSEPGDSLRLGSIDI
9635243 CDS_12	74	614 MSKSHSDVVGTVGSLNRYLFDYMDIPDRISQKLRREITDPKTCNASKIPVLVLAEEEVSRMIDHDKDGYTKVQVK
9635243 CDS_13	1527	2505 MALLQKLGSKVSNSKRMSMSPAMIPLDSINQDLDVDPQEQDAKNNKEGKKKLDVSMPLTGKLPGLKKQVDTG
19744922 CDS_14	89	686 MPNVNLSATAKGGGHYIEDQWDTQVVEAGVFFDDWWVHVHEAWNKFQDNLRGINSVASSRSRQVAEYLALDRLI
19744922 CDS_15	803	2390 MDMTKTVEEKKTNGTDISKVGFENSTIPKVTGEMGGDSSTSLLKETLKVAQDQPLSVNDGAKSKLDSSDRQ
19744922 CDS_16	2358	2754 MKTTVGSRPNKYWPIVAGIGVGLFAYLIFSNQKHSTESGDNIHFKANGGSYRDGSKSISYNRNHPFAYGNASSP
19744922 CDS_17	2563	3031 MAMPHPLECCCPQCLPSSFPIYGEQEIPCSSETQAEATTPTVEKTRVANVLTDILDDHHYAILASLFIIALWLIIYLSII
944542958 CDS_18	205	3502 MSIPEKTYTWFSDNVVKVPTAWQTDGHTWARICELRVERFPNADFSRFDFGQTEWRTRRDISDIFTLHTTC
216905810 CDS_19	39	2280 MMIFQFTMPVPRIDGPGRPSMGMALRAVGKFLVNCCSAFGDNWADSTRNRTQHVCALQYFGDLPICEIVKS
216905810 CDS_20	2279	2486 MDKSPQRGRSRSRQTOGPKGPKPENKQIQVAAHAAVDKARGKPPGGDHGDFIVAHVTVTNINFINI

Figure 16: Shows how the data looks in excel. The first column shows the Virus\_ID. The second column shows the Unique\_ID given to every CDS region found. The third column shows the length/ region of CDS within the viral genome. Finally, the last column is the translated protein sequence.

## 9. Using NCBI: CDD data on possible Motifs and Domains within the Translated Protein

Sequence was obtained. NCBI:CDD allows for batch sequence searches which looked like:

The screenshot shows the NCBI Batch Web CD-Search Tool. The main interface has a search bar at the top where users can enter protein sequences. Below the search bar are several search options: 'Search mode' (set to 'Automatic'), 'Expect Value threshold' (set to 0.01), 'Composition-corrected scoring' (checkbox checked), 'Apply low-complexity filter' (checkbox unchecked), and 'Maximum number of hits' (set to 500). There is also a checkbox for 'Include related sequences'. At the bottom of the search form, there is a 'Browse' button for file uploads and a 'Submit' button.

Figure 17: Shows the CDD batch search.

Vpipmtdom_01	VpipCDS_01	260	693	9.84E-104	340.818 cl20022	Peptidase C6 superfamily
Vpipmtdom_02	VpipCDS_01	2254	2715	2.34E-98	326.635 cl02808	RT like superfamily
Vpipmtdom_03	VpipCDS_01	2806	3038	1.58E-91	297.593 cl02961	Poty coat superfamily
Vpipmtdom_04	VpipCDS_01	1495	1768	5.16E-68	231.607 cl07169	Poty PP superfamily
Vpipmtdom_05	VpipCDS_01	1986	2219	3.71E-57	199.161 cl24133	Peptidase C4 superfamily
Vpipmtdom_06	VpipCDS_01	2	232	2.52E-54	191.006 pfam01577	Peptidase S30
Vpipmtdom_07	VpipCDS_01	708	1160	1.11E-53	196.399 cl16319	Potyvirid-P3 superfamily
Vpipmtdom_08	VpipCDS_01	1198	1329	4.91E-12	66.0336 cl28899	DEAD-like helicase N superfamily
Vpipmtdom_09	VpipCDS_01	1186	1334	2.03E-24	103.341 smart00487	DEXOc
Vpipmtdom_10	VpipCDS_01	1372	1468	3.04E-09	55.6805 smart00490	HELIc
Vpipmtdom_11	VpipCDS_03	45	325	2.26E-78	261.839 pfam01660	Vmethyltransf
Vpipmtdom_12	VpipCDS_03	722	817	2.55E-16	76.3005 cl05113	Peptidase C21 superfamily
Vpipmtdom_13	VpipCDS_03	1450	1639	3.45E-10	64.2033 cl03049	RdRp 2 superfamily
Vpipmtdom_14	VpipCDS_03	1840	1995	4.22E-10	60.6161 cl03052	Tymo coat superfamily
Vpipmtdom_15	VpipCDS_03	1090	1140	0.00610147	37.156 cl38915	DEAD-like helicase C superfamily
Vpipmtdom_16	VpipCDS_03	908	1138	6.33E-42	154.07 pfam01443	Viral helicase1
Vpipmtdom_17	VpipCDS_04	2	417	0.00019676	43.7737 PHA03247	PHA03247
Vpipmtdom_18	VpipCDS_05	85	194	0.00705485	36.7863 cl26680	SIdE superfamily
Vpipmtdom_19	VpipCDS_06	1	110	2.97E-71	207.393 cl05711	RTBV P12 superfamily
Vpipmtdom_20	VpipCDS_07	1240	1389	4.84E-44	157.757 cd01647	RT LTR
Vpipmtdom_21	VpipCDS_07	1488	1614	2.58E-16	76.3773 cl14782	RNase H like superfamily
Vpipmtdom_22	VpipCDS_07	976	1073	2.38E-05	44.2496 cl03033	retropepsin like
Vpipmtdom_23	VpipCDS_07	1248	1392	1.97E-35	133.572 pfam00078	RVT 1
Vpipmtdom_24	VpipCDS_08	1	389	0	823.578 cl05630	RTBV P46 superfamily
Vpipmtdom_25	VpipCDS_12	9	168	2.78E-21	85.1108 cl20276	Tenui NCP superfamily
Vpipmtdom_26	VpipCDS_13	67	271	3.12E-24	99.4412 cl03993	Tenui NS4 superfamily
Vpipmtdom_27	VpipCDS_14	18	183	1.52E-64	196.31 pfam00721	TMV coat
Vpipmtdom_28	VpipCDS_15	266	359	7.12E-06	46.0084 cl28899	DEAD-like helicase N superfamily

Figure 18: shows the final results that were obtained from the NCBI: CDD database. These results are edited and arranged based on the schema. Look at image 1 for reference.

10. Now that all the data was gathered, the next step was to place the data in the database so that further queries can be written to allow for greater analysis of the data.

```

1  #Create tables and populates them
2  import csv
3  import psycopg2
4  import psycopg2.extras
5  from urlparse import urlparse, uses_netloc
6  import configparser
7  #initialize connection to the database
8  def connect_to_db(conn_str):
9      uses_netloc.append('postgres')
10     url = urlparse(conn_str)
11
12     conn = psycopg2.connect(database=url.path[1:],
13                            user=url.username,
14                            password=url.password,
15                            host=url.hostname,
16                            port=url.port)
17
18     return conn
19 #Load connection infomration from Config.ini
20 config = configparser.ConfigParser()
21 config.read('config.ini')
22 connection_string = config['database']['postgres_connection']
23 conn = connect_to_db(connection_string)
24 cursor = conn.cursor()
25 #Create the tables and columns
26
27 cursor.execute("""
28     CREATE TABLE IF NOT EXISTS Plants (ScientificName VARCHAR(255) NOT NULL, CommonName VARCHAR(255) NOT NULL,
29     PRIMARY KEY(ScientificName))
30     """)
31 cursor.execute("""
32     CREATE TABLE IF NOT EXISTS Infect (Infection_ID VARCHAR(20) NOT NULL, ScientificName VARCHAR(255) NOT NULL,
33     VirusName VARCHAR(255) NOT NULL , PRIMARY KEY(Infection_ID),
34     FOREIGN KEY (ScientificName) REFERENCES Plants(ScientificName)
35     ON DELETE CASCADE ON UPDATE CASCADE
36     DEFERRABLE INITIALLY DEFERRED)
37     """)
38 cursor.execute("""
39     CREATE TABLE IF NOT EXISTS Viruses (Name VARCHAR(255) NOT NULL, ID INTEGER NOT NULL,
40     Sequence TEXT, PRIMARY KEY (ID))
41     """)
42 cursor.execute("""
43     CREATE TABLE IF NOT EXISTS CDS (VirusID INTEGER NOT NULL, CDSID VARCHAR(20) NOT NULL, CDSStart INTEGER,
44     CDSStop INTEGER, CDSTranslatedSeq TEXT, PRIMARY KEY (CDSID),
45     FOREIGN KEY (VirusID) REFERENCES Viruses(ID)
46     ON DELETE CASCADE ON UPDATE CASCADE
47     DEFERRABLE INITIALLY DEFERRED)
48     """)
49 cursor.execute("""
50     CREATE TABLE IF NOT EXISTS CDSResult (CDSResultID VARCHAR(20) NOT NULL, CDSUniqueID VARCHAR(20) NOT NULL,
51     MotifStart INTEGER, MotifStop INTEGER, EValue DOUBLE PRECISION, BitScore DOUBLE PRECISION,
52     AccessionNo VARCHAR(50), MotifName VARCHAR(225), PRIMARY KEY (CDSResultID),

```

```

MotifStart INTEGER, MotifStop INTEGER, EValue DOUBLE PRECISION, BitScore DOUBLE PRECISION,
AccessionNo VARCHAR(50), MotifName VARCHAR(225), PRIMARY KEY (CDSResultID),
FOREIGN KEY (CDSUniqueID) REFERENCES CDS(CDSID)
ON DELETE CASCADE ON UPDATE CASCADE
DEFERRABLE INITIALLY DEFERRED)
""")
#Insert the data into the tables
with open('Plants.csv', 'r') as Plants:
    reader = csv.reader(Plants)
    for row in reader:
        cursor.execute('INSERT INTO Plants VALUES (%s,%s)',row)

with open('Infect.csv', 'r') as Infects:
    reader = csv.reader(Infects)
    for row in reader:
        cursor.execute('INSERT INTO Infect VALUES (%s,%s,%s)',row)

with open('Viruses.csv', 'r') as Viruses:
    reader = csv.reader(Viruses)
    for row in reader:
        cursor.execute('INSERT INTO Viruses VALUES (%s,%s,%s)',row)

with open('CDS.csv', 'r') as CDS:
    reader = csv.reader(CDS)
    for row in reader:
        cursor.execute('INSERT INTO CDS VALUES (%s,%s,%s,%s,%s)',row)

with open('CDS_Motifs.csv', 'r') as CDS_Motifs:
    reader = csv.reader(CDS_Motifs)
    for row in reader:
        cursor.execute('INSERT INTO CDSResult VALUES (%s,%s,%s,%s,%s,%s,%s,%s)',row)
# Commit and close the connection after the queries
conn.commit()
conn.close()

```

Figure 19 &amp; 20: Shows the scrip used to seed the database.

11. Next to display the data in a table properly from the database a Flask application was set up and quires were written to allow user to further use the gathered data and to create a web interface.

```

from flask import Flask, render_template, request, redirect
from .Queries import *
app = Flask(__name__)
#Calls function that connects to the database and Loads all data onto the database
#At initialization of the app the user is first directed to the home page
##### Main Page #####
@app.route('/')
def Home():
    return render_template('Home.html')
##### Home Pages #####
@app.route('/Home/Common')
def Home_Common():
    return render_template('Search/Home-Common.html')

@app.route('/Home/Virus')
def Home_Virus():
    return render_template('Search/Home-Virus.html')

@app.route('/Home/Motif')
def Home_Motif():
    return render_template('Search/Home-Motif.html')
##### Search Results #####
@app.route('/Scientific Name/Result/', methods=['GET'])
def search_results_Scientific():
    searchs = request.args.get('search')
    Results = search_scientific(searchs)
    return render_template("Result/Scientific-Result.html",Results = Results)

@app.route('/Common Name/Result/', methods=['GET'])
def search_results_Common():
    searchs = request.args.get('search')
    Results = search_full(searchs)
    return render_template("Result/Common-Result.html",Results = Results)

@app.route('/Virus/Result/', methods=['GET'])
def search_results_Virus():
    searchs = request.args.get('search')
    Results = search_virus(searchs)
    return render_template("Result/Virus-Result.html",Results = Results)

@app.route('/Motif/Result/', methods=['GET'])
def search_results_Motif():
    searchs = request.args.get('search')
    Results = search_motif(searchs)
    return render_template("Result/Motif-Result.html",Results = Results)
##### About Page #####
@app.route('/About/')
def About():
    return render_template('About/About.html')
##### View Plants #####
@app.route('/Plants/')
def Plants():

```

Figure 21: Shows the main flask app code. For full code simply open the “Main.py” file in the application folder.

```

#Sends queries to the database to receive information
import psycopg2
import pandas as pd
import psycopg2.extras
from urllib.parse import urlparse, uses_netloc
import configparser
#initialize connection to the database
def connect_to_db(conn_str):
    uses_netloc.append('postgres')
    url = urlparse(conn_str)
    conn = psycopg2.connect(database=url.path[1:],
                           user=url.username,
                           password=url.password,
                           host=url.hostname,
                           port=url.port)
    return conn
#Load connection information from Config.ini
config = configparser.ConfigParser()
config.read('config.ini')
connection_string = config['database']['postgres_connection']
conn = connect_to_db(connection_string)
#Get all data from plants
def get_plants():
    with conn.cursor(cursor_factory=psycopg2.extras.RealDictCursor) as cursor:
        cursor.execute('SELECT * FROM Plants')
        for d in cursor:
            yield d
#Get all data from viruses
def get_viruses():
    with conn.cursor(cursor_factory=psycopg2.extras.RealDictCursor) as cursor:
        cursor.execute('SELECT * FROM Viruses')
        for d in cursor:
            yield d
#Get Just the scientific name and virus name from infects
def get_infects():
    with conn.cursor(cursor_factory=psycopg2.extras.RealDictCursor) as cursor:
        cursor.execute('SELECT scientificname,virusname FROM Infect')
        for d in cursor:
            yield d
#Get all data from the CDS as well as the virus table.
#To display the name of the virusID
def get_CDS():
    with conn.cursor(cursor_factory=psycopg2.extras.RealDictCursor) as cursor:
        cursor.execute('SELECT * FROM cds JOIN viruses ON cds.virusid = viruses.id')
        for d in cursor:
            yield d
#Get all data from the motifs_domain
def get_motifs_domain():
    with conn.cursor(cursor_factory=psycopg2.extras.RealDictCursor) as cursor:
        cursor.execute('SELECT * FROM CDSresult\\
                      JOIN cds ON cdsresult.cdsuniqueid = cds.cdsid\\
                      JOIN viruses ON cds.virusid = viruses.id')

```

Figure 22: The figure shows the queries that were then called to the database and displayed on the flask application using HTML (to display on browser) and CSS (to style the display).

## Results

### Over All Statistics for the application

10 Plant genuses and 73 Plant species were recorded
193 Viruses were recorded that infect at-least one plant species
Over all 764 Recorded infections
486 Coding regions identified for various viruses
717 Various motifs identified from the coding regions

### Folder structure of the application

a.

Name	Date modified	Type	Size
__pycache__	5/5/2019 5:28 PM	File folder	
Data_Backup	5/3/2019 2:28 PM	File folder	
Other_Scripts	5/3/2019 2:28 PM	File folder	
static	5/5/2019 7:48 PM	File folder	
templates	5/4/2019 10:36 PM	File folder	
config	4/20/2019 7:35 PM	Configuration sett...	1 KB
_init_	4/16/2019 12:07 PM	PY File	0 KB
Graphs	5/3/2019 12:02 PM	PY File	2 KB
Main	5/5/2019 2:47 PM	PY File	3 KB
Queries	5/5/2019 5:28 PM	PY File	5 KB

b.

CDS	4/20/2019 7:44 PM	Microsoft Office E...	271 KB
CDS_Motifs	4/14/2019 8:25 AM	Microsoft Office E...	51 KB
Infect	4/20/2019 8:21 PM	Microsoft Office E...	44 KB
Plants	4/20/2019 8:23 PM	Microsoft Office E...	3 KB
Viruses	4/13/2019 8:29 PM	Microsoft Office E...	791 KB

c.

Script-1	5/3/2019 12:02 PM	R Source File	1 KB
Script-2	5/3/2019 12:02 PM	R Source File	1 KB
Script-3	5/3/2019 12:02 PM	PY File	1 KB
Script-4	5/3/2019 12:03 PM	PY File	1 KB
Script-5	3/23/2019 2:02 PM	PY File	1 KB
Script-6	3/24/2019 1:59 PM	PY File	1 KB
Script-7	3/24/2019 11:58 PM	PY File	1 KB
Seed_database	5/3/2019 12:04 PM	PY File	4 KB

Figure 23 – a: Shows the folder structure of the application and how the scripts were arranged. The \_\_pycache\_\_ folder is created at use of the application and holds information on the usage of the app. The Data\_Backup folder (figure 23 - b) contains all the data

that was gathered in csv format and tab-delimited format. The Other\_Scripts folder (figure 23- c) contains all the scripts that were used for gathering the data including script used to seed the database. Static and Templates contain CSS and HTML files for the application. `__init__.py` file simply initializes the application and should be empty. The `Graphs.py` script creates the graph based on the data in the database. `Main.py` is the main application and is called to run the application. `Queries.py` connects the database and the application together and calls queries to the database.

```
Microsoft Windows [Version 10.0.17763.475]
(c) 2018 Microsoft Corporation. All rights reserved.

C:\Users\ritwi>cd Documents\College\Bioinformatics\Project\App

C:\Users\ritwi\Documents\College\Bioinformatics\Project\App>set FLASK_APP=Main.py

C:\Users\ritwi\Documents\College\Bioinformatics\Project\App>flask run
 * Serving Flask app "Main.py"
 * Environment: production
   WARNING: Do not use the development server in a production environment.
   Use a production WSGI server instead.
 * Debug mode: off
 * Running on http://127.0.0.1:5000/ (Press CTRL+C to quit)
```

Figure 24: Shows how to run the application once everything was set up.

## Application

### Home and About pages

The application contains a home page, about page, a browse section (allows the user to browse the all the data that was gathered for the project) and a search bar (which can be set according to what the user wants to search for).

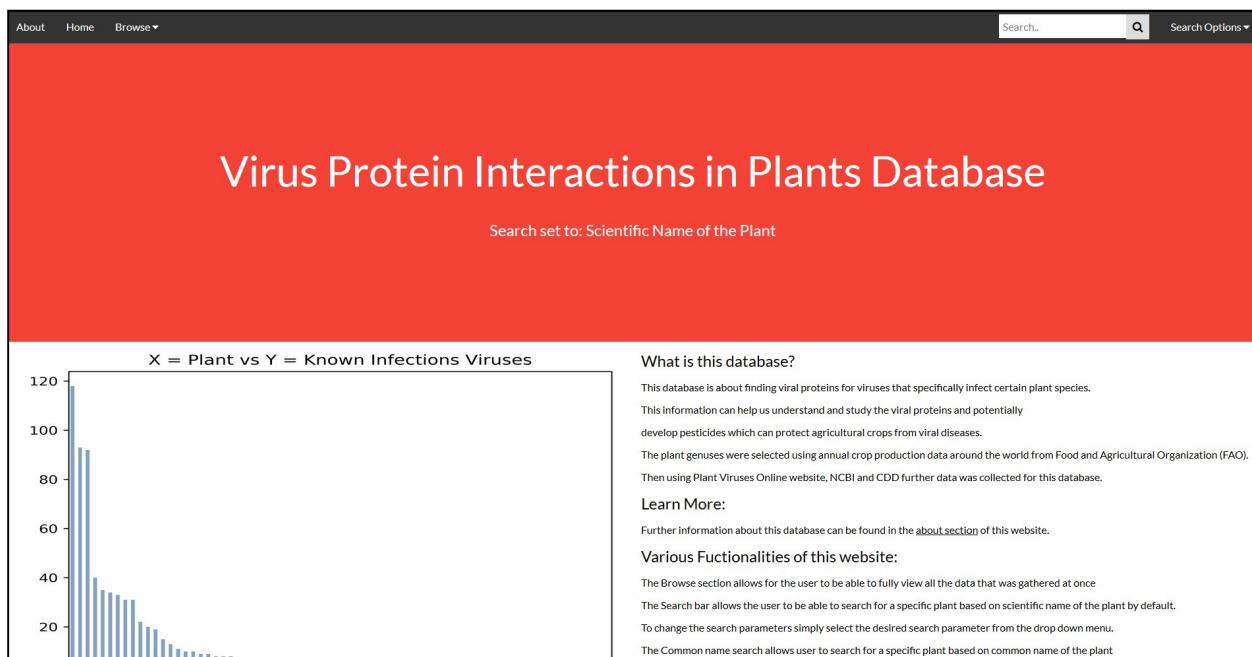


Figure 25: Shows the home page of the application. The graph showcases all the recorded infections for a particular plant species. The home page of the application goes in more detail about the database and about the graph itself.

About    Home    Browse▼    Search..  Search Options▼

## Virus Protein Interactions in Plants Database - About

### What is this database?

In short, this database is about finding viral proteins for viruses that specifically infect certain plant species.

### Why is this database useful?

This information can help us understand and study the viral proteins and potentially develop pesticides which can protect agricultural crops from viral diseases.

### How was the data collected?

First 'Popular' plants that are agriculturally grown around the world was taken from Food and agricultural organization (FAO) Then the data was organized based on the amount of yield. Only the top 10 highly produced plants were selected. Then the Genus name was determined from the common name using the R-package, Taxize

Once the Genuses were determined then using beautiful soup the data for species and their susceptibility to various viruses was scraped from plant viruses online website.

With the use of Biopython the viruses were then searched and a corresponding ID was obtained for the viruses. Using the obtained ID the NCBI (National Center For Biotechnology information) database was searched for viral sequences and any possible CDS's (Coding regions) the virus may have.

After obtaining the CDS's (Coding regions) for the viruses the data was then sent to CDD's (Conserved Domains Database) batch search tool where the expect value threshold was set to 0.01. The CDD's batch search tool then gave the data on domains and motifs for the CDS's.(Coding regions)

This, and all other data was then placed in ElephantSQL (Postgresql) and the database was created

### Future Updates and changes

First: Since the database currently relies on the data obtained by plants and viruses online website, to increase accuracy and efficiency of this database the viral susceptibility data will be obtained by this database itself.

Second: A better more interactive GUI. So that the data isn't presented only in tables

Third: Addition of more plant genera

### Resources Used

[Food and Agricultural Organization](#)  
[Plant Virus Online](#)  
[NCBI \(National Center For Biotechnology Information\)](#)

Figure 26: Shows the about page of the application. The about page summarizes how the data was obtained and what resources were used for the creation of the application.

### The Browse section

Contains 5 different tables that showcase all the data from the database.

About    Home    Browse▼    Search..  Search Options▼

### Plants Table

The table displays the list of all the plants.  
Which includes their common and scientific names.

Scientific Name	Common Name
Saccharum officinarum	Sugar Cane
Zea mays	Corn
Zea mays ssp. Mays	Common Corn-Variations
Zea mays ssp. mexicana	Corn- Subspecies
Oryza australiensis	Australian Rice
Oryza barthii	Wild Rice
Oryza rufipogon	Wild Red Rice
Oryza glaberrima	African Rice
Oryza latifolia	Broadleaf Rice
Oryza longistaminata	Longstamin wild rice
Oryza nivara	Wild Asian Rice
Oryza perennis	Brown Bread Rice
Oryza punctata	Red Rice
Oryza sativa	Asian Rice

Figure 27: Shows the first table within the browse section.

About	Home	Browse ▾	Search..	Search Options ▾
<b>Virus Table</b>				
Name	ID	Sequence		
Brome mosaic bromovirus	1491772236	ATGTCGACTTCAGGAACCTGGTAAGATGACTCGCGC< >		
Cacao yellow mosaic tymovirus	1464310171	ACCACGTTCAATAGCCAACATGTCTCCGATCTCATTCT< >		
Zucchini yellow mosaic potyvirus	1591442859	TATATAGAGATGAGAAATGCAGAGGCACCATACATGCC< >		
Clover yellow mosaic potexvirus	1569273958	ACCCCATAGCCACGGTTAAGTTGCCAGATTGCGAAC< >		
Potato T trichovirus	1444394066	ACAATTAAAATGTCTTCTCCTCAGAACGCCAGCAG< >		
Clitoria yellow vein tymovirus	1464307030	ATGTCGACCGACGTATCGTTCTGCCAACCTCTCATC< >		
Lettuce infectious yellows closterovirus	258676964	GTAATAACTTACCAAATTTCGTCCAAGGGAGACAG< >		
Beet mild yellowing luteovirus	19881389	ACAAAAGAAACCCAGCGAGGATCTAGCAGTCTATGCAAT< >		
<i>Cereal northern mosaic cytorhabdovirus</i>	9625417	CACGACCACTGATCGAACAAACCTGAATCATGGTGACC		

Figure 28: Shows the second table within the browse section. The scroll under the sequences allows the user to view the whole sequence without the need to scroll horizontally on the page itself.

About	Home	Browse ▾	Search..	Search Options ▾
<b>Infection Table</b>				
Plant Name	Virus Name			
Saccharum officinarum	Maize dwarf mosaic potyvirus			
Saccharum officinarum	Sugarcane bacilliform badnavirus			
Saccharum officinarum	Sugarcane Fiji disease fijivirus			
Saccharum officinarum	Sugarcane mosaic potyvirus			
Saccharum officinarum	Sugarcane streak monogeminivirus			
Zea mays	Barley stripe mosaic hordeivirus			
Zea mays	Barley yellow dwarf luteovirus			
Zea mays	Barley yellow striate mosaic cytorhabdovirus			
Zea mays	Bermuda grass etched-line marafivirus			
Zea mays	Brome mosaic bromovirus			
Zea mays	Cassia yellow blotch bromovirus			
Zea mays	Cereal northern mosaic cytorhabdovirus			

Figure 29: Shows the third table within the browse section.

About Home Browse ▾ Search..  Search Options ▾

### CDS Table

Shows the CDS regions that were found for each virus

Virus ID	Virus Name	CDS start	CDS Stop	CDS Translated Sequence
1586082777	Barley yellow striate mosaic cytorhabdovirus	3911	4412	MSRVTRFLFLKLDVEMEVDFGD< >
18490052	Maize dwarf mosaic potyvirus	139	9265	MAGTWTHVTHKWQPNLDNPR< >
18490052	Maize dwarf mosaic potyvirus	2682	2922	NLCRSVESSVDRIIIVWKILRNMA< >
14141972	Maize rayado fino marafivirus	96	6180	MSSFLRGGHLLSGVESLPTTHRI< >
14141972	Maize rayado fino marafivirus	301	1561	MPLTPTPSIRPSRPTFSMSGPTT< >
18026839	Rice tungro bacilliform badnavirus	67	667	VLKRNLTSQNIESRYEKFELDLA< >
18026839	Rice tungro bacilliform badnavirus	663	996	MSADYPTFKEALEKFKNLESDTA< >
18026839	Rice tungro bacilliform badnavirus	992	6026	MSLRPFTGTSRTITQDSTSESNIKI< >

Figure 30: Shows the forth table within the browse section. The scroll under the sequences allows the user to view the whole sequence without the need to scroll horizontally on the page itself.

About Home Browse ▾ Search..  Search Options ▾

### Motifs & Domain Table

Table displays the list of all the Motifs and Domains found within a CDS region.

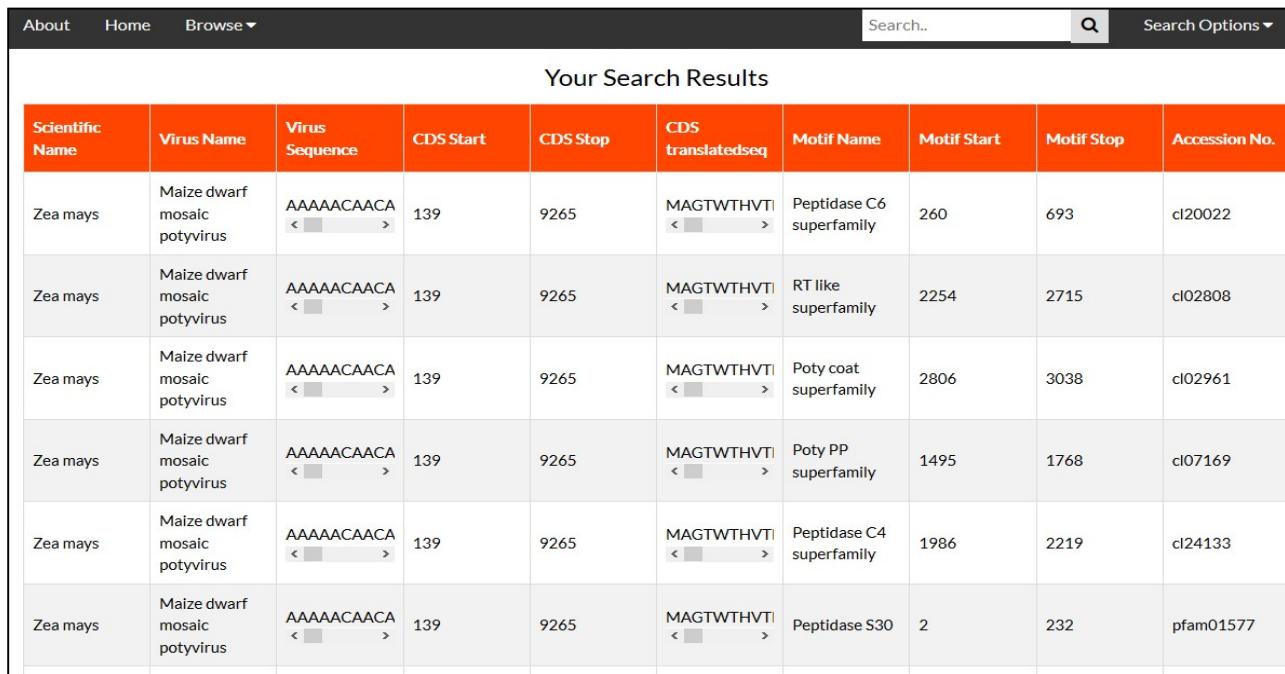
The table includes the Unique CDS ID, start and stop values within the translated region, E-value, Bit-score and the accession number of the motifs in other databases

Virus Name	Motif Name	Motif Start Point	Motif End Point	E-Value	Bit-Score	Accession Number
Maize dwarf mosaic potyvirus	Peptidase C6 superfamily	260	693	9.84e-104	340.818	cl20022
Maize dwarf mosaic potyvirus	RT like superfamily	2254	2715	2.34e-98	326.635	cl02808
Maize dwarf mosaic potyvirus	Poty coat superfamily	2806	3038	1.58e-91	297.593	cl02961
Maize dwarf mosaic potyvirus	Poty PP superfamily	1495	1768	5.16e-68	231.607	cl07169
Maize dwarf mosaic potyvirus	Peptidase C4 superfamily	1986	2219	3.71e-57	199.161	cl24133
Maize dwarf mosaic potyvirus	Peptidase S30	2	232	2.52e-54	191.006	pfam01577
Maize dwarf mosaic potyvirus	Potyvirid-P3 superfamily	708	1160	1.11e-53	196.399	cl16319
Maize dwarf mosaic potyvirus	DEAD-like helicase N superfamily	1198	1329	4.91e-12	66.0336	cl28899
Maize dwarf mosaic potyvirus	DEXDc	1186	1334	2.03e-24	103.341	smart00487
Maize dwarf mosaic potyvirus	HELICc	1372	1468	3.04e-09	55.6809	smart00490
Maize rayado fino marafivirus	Vmethyltransf	45	325	2.26e-78	261.839	pfam01660

Figure 31: Shows the fifth table within the browse section. The accession numbers for the motifs can be used in the CDD website to search for more information on the motifs that were found.

## Search Section

The search section allows the user to search by the scientific name by default. However, it also allows the user to search by common name, virus and motif as well. The search option for scientific and common name of the plant is for the user to see the search results for a specified plant rather than needing to browse through the browse section of the application. The Virus search allows the user to search for a specific virus and view the virus sequence, cds and motifs found for the cds that a virus might posses. The motif search results allows user to search for various viruses that may possess a specific motif; this search also provides the accession number for the motifs.



The screenshot shows a web-based search interface. At the top, there is a navigation bar with links for 'About', 'Home', and 'Browse'. To the right of the navigation bar are search input fields labeled 'Search..', a magnifying glass icon, and a 'Search Options' dropdown menu. Below the navigation bar, the main content area has a title 'Your Search Results'. A table is displayed with the following columns: Scientific Name, Virus Name, Virus Sequence, CDS Start, CDS Stop, CDS translatedseq, Motif Name, Motif Start, Motif Stop, and Accession No. There are six rows of data in the table, all corresponding to 'Zea mays' and 'Maize dwarf mosaic potyvirus'. Each row shows the same sequence: 'AAAAACAAACA' with a greyed-out segment between positions 139 and 9265. The 'CDS translatedseq' column shows 'MAGTWTHVTI' with a cursor at position 139. The 'Motif Name' column lists various superfamily names like 'Peptidase C6 superfamily', 'RT like superfamily', 'Poty coat superfamily', 'Poty PP superfamily', 'Peptidase C4 superfamily', and 'Peptidase S30'. The 'Accession No.' column contains accession numbers such as 'cl20022', 'cl02808', 'cl02961', 'cl07169', 'cl24133', and 'pfam01577'.

Scientific Name	Virus Name	Virus Sequence	CDS Start	CDS Stop	CDS translatedseq	Motif Name	Motif Start	Motif Stop	Accession No.
Zea mays	Maize dwarf mosaic potyvirus	AAAAACAAACA < [ ] >	139	9265	MAGTWTHVTI < [ ] >	Peptidase C6 superfamily	260	693	cl20022
Zea mays	Maize dwarf mosaic potyvirus	AAAAACAAACA < [ ] >	139	9265	MAGTWTHVTI < [ ] >	RT like superfamily	2254	2715	cl02808
Zea mays	Maize dwarf mosaic potyvirus	AAAAACAAACA < [ ] >	139	9265	MAGTWTHVTI < [ ] >	Poty coat superfamily	2806	3038	cl02961
Zea mays	Maize dwarf mosaic potyvirus	AAAAACAAACA < [ ] >	139	9265	MAGTWTHVTI < [ ] >	Poty PP superfamily	1495	1768	cl07169
Zea mays	Maize dwarf mosaic potyvirus	AAAAACAAACA < [ ] >	139	9265	MAGTWTHVTI < [ ] >	Peptidase C4 superfamily	1986	2219	cl24133
Zea mays	Maize dwarf mosaic potyvirus	AAAAACAAACA < [ ] >	139	9265	MAGTWTHVTI < [ ] >	Peptidase S30	2	232	pfam01577

Figure 32: Shows an example of the default search results where “Zea Mays” plant was searched.

Common Name	Scientific Name	Virus Name	Virus Sequence	CDS Start	CDS Stop	CDS translatedseq	Motif Name	Motif Start	Motif Stop	Accession No.
Sugar Cane	Saccharum officinarum	Maize dwarf mosaic potyvirus	AAAAACAAAC <  >	139	9265	MAGTWTHV <  >	Peptidase C6 superfamily	260	693	cl20022
Sugar Cane	Saccharum officinarum	Maize dwarf mosaic potyvirus	AAAAACAAAC <  >	139	9265	MAGTWTHV <  >	RT like superfamily	2254	2715	cl02808
Sugar Cane	Saccharum officinarum	Maize dwarf mosaic potyvirus	AAAAACAAAC <  >	139	9265	MAGTWTHV <  >	Poty coat superfamily	2806	3038	cl02961
Sugar Cane	Saccharum officinarum	Maize dwarf mosaic potyvirus	AAAAACAAAC <  >	139	9265	MAGTWTHV <  >	Poty PP superfamily	1495	1768	cl07169
Sugar Cane	Saccharum officinarum	Maize dwarf mosaic potyvirus	AAAAACAAAC <  >	139	9265	MAGTWTHV <  >	Peptidase C4 superfamily	1986	2219	cl24133

Figure 33: Shows an example where “Sugar Cane” was searched by setting the search settings to “common name” by using the dropdown menu next to the search bar.

Virus Name	Virus Sequence	CDS Start	CDS Stop	CDS translatedseq	Motif Name	Motif Start	Motif Stop
Heracleum latent trichovirus	TCCCTCCGATTATG <  >	2008	2602	MDGISRSARIRNAV <  >	Tricho coat superfamily	10	197

Figure 34: Shows an example where “Heracleum latent trichovirus” was searched by setting the search settings to “viruses” by using the dropdown menu next to the search bar.

Motif Name	Accession No.	Virus Name
Viral Hsp90	pfam03225	Beet pseudo-yellows closterovirus
Viral Hsp90	pfam03225	Beet yellows closterovirus

Figure 35: Shows an example where “Viral Hsp90” Motif was searched by setting the search settings to “Motif” by using the dropdown menu next to the search bar.

## Discussion & Conclusion

Overall the application was connected and the data was retrieved from the database properly. The overall purpose of the application was fulfilled; however there is a lot of room for improvements. Future updates for the application would involve creating a proper GUI for viewing the sequences and results. The results would not contain redundancy. The reason for redundancy in the current version of the application is simply due to the join statements on the various tables, which are displayed without any edits made to the joined tables. The search bar needs to be made more user friendly. Currently to change the search parameter one must use the drop down menu. The drop down menu simply directs the user to a version of the home page where the search bar itself has been changed based on what the user wishes to search for.

Another update that needs to be made for the search bar is to make the searches not case sensitive. Currently the searches are very case sensitive and require the user to enter the exact search term. For example, in Figure 33: if the user had searched for “sugar cane” instead of “Sugar Cane” no results would be shown to the user. Additionally, more plant genera need to be added to the database. The reason why only 10 genera were used in the current version of the database is because Elephant SQL has a limit of 128 MB of free usage. Once the data exceeds 128 MB the user is required to purchase more space for the database. Finally the application currently requires data from Plant and viruses online website; the website itself isn’t updated all too much. Hence, to remove reliance on Plant and viruses online website, an algorithm that generates viral infection data based on plant species is required, within the application. Creating this algorithm would remove reliance on Plant and viruses online website. Despite these future improvements the database is able to record viral proteins along with the type of host the virus usually infects. This could prove beneficial in understanding what the viruses are capable of and thereby giving us a step ahead against countering viruses and protecting agricultural crops.

## References

1. (Lecoq) Lecoq, H. "Discovery of the First Virus, the Tobacco Mosaic Virus: 1892 or 1898?" *Comptes Rendus De L'Academie Des Sciences. Serie III, Sciences De La Vie*, U.S. National Library of Medicine, Oct. 2001, [www.ncbi.nlm.nih.gov/pubmed/11570281](http://www.ncbi.nlm.nih.gov/pubmed/11570281).
2. (Smith) Smith ,H. "Mechanisms of Viral Pathogenicity" *Department of microbiology, University of Birmingham, Birmingham, England*, Vol.36,NO.3 Bacteriological Reviews. American Society of Microbiology, Sept. 1972,  
<https://mmbre.asm.org/content/mmbre/36/3/291.full.pdf>
3. (Kamps) Kamps, Bernd Sebastian, et al. *Influenza Book | Pathogenesis and Immunology*, [www.influenzareport.com/ir/pathogen.htm](http://www.influenzareport.com/ir/pathogen.htm).
4. (Lodish) Lodish H, Berk A, Zipursky SL, et al. Molecular Cell Biology. 4th edition. New York: W. H. Freeman; 2000. Section 6.3, Viruses: Structure, Function, and Uses.  
Available from:
5. (Rui) Rui, Liu Sheng, et al. "MicroRNA-Mediated Gene Silencing in Plant Defense and Viral Counter-Defense." *Frontiers in Microbiology* , vol. 08, no. 1664-302x, 2017, p. 1801., doi:10.3389/fmicb.2017.01801.
6. (Pratt) Pratt, Ashley J and Ian J MacRae. "The RNA-induced silencing complex: a versatile gene-silencing machine" *Journal of biological chemistry* vol. 284,27 (2009): 17897-901.
7. (Kamthan) Kamthan, Ayushi et al. "Small RNAs in plants: recent development and application for crop improvement" *Frontiers in plant science* vol. 6 208. 2 Apr. 2015, doi:10.3389/fpls.2015.00208