

# Orthopoxvirus Genome Analysis for Species Classification and Virulence Prediction

G Gnanashree  
BL.SC.U4AIE24014  
CSE-AIE, Amrita School of Engineering  
Amrita Vishwa Vidyapeetham, Bengaluru  
Email: [gnanashreegekkala@gmail.com](mailto:gnanashreegekkala@gmail.com)

PVS Lalitha  
BL.SC.U4AIE24031  
CSE-AIE, Amrita School of Engineering,  
Bengaluru  
Email: [wby1606@gmail.com](mailto:wby1606@gmail.com)

Y Likhitha Sree  
BL.SC.U4AIE24056  
CSE-AIE, Amrita School of Engineering,  
Bengaluru  
Email: [likhithayarraguntla@gmail.com](mailto:likhithayarraguntla@gmail.com)

## Abstract

*Orthopoxviruses are closely related double-stranded DNA viruses that include Variola, Monkeypox, Cowpox, and Vaccinia viruses. Although these viruses share high genomic similarity, they differ in host range and virulence. Understanding these differences at the genomic level is essential for evolutionary analysis and public health monitoring.*

*In this study, a systematic genome-based computational analysis was performed on four Orthopoxvirus species. Complete genome sequences were retrieved and curated to construct a high-quality, non-redundant dataset using a k-mer-based cosine similarity approach. Sequence-derived features including nucleotide composition, k-mer frequencies, and dinucleotide bias were extracted to quantitatively represent genomic variation.*

*Two machine learning models were developed: a multi-class classification model for species identification, and a virulence prediction model to categorize genomes based on pathogenic potential. Feature importance analysis was conducted to determine the genomic attributes contributing most significantly to species differentiation and virulence variation. The results demonstrate that alignment-free genomic features combined with supervised learning provide reliable classification and meaningful biological insights.*

## Keywords

*Orthopoxvirus, Genome Analysis, k-mer Analysis, Machine Learning, Species Classification, Virulence Prediction*

## I. INTRODUCTION

Orthopoxviruses are large double-stranded DNA viruses belonging to the Poxviridae family and include medically significant species such as Variola virus, Monkeypox virus, Cowpox virus, and Vaccinia virus. Despite sharing substantial genomic similarity, these viruses differ in host specificity, transmission patterns, and virulence. Studying their genomic variations is important for understanding viral evolution, monitoring emerging outbreaks, and supporting vaccine and therapeutic research.

Genome-based analysis provides a robust approach for distinguishing closely related viral species. Traditional alignment-based sequence comparison methods can become computationally intensive when handling complete genomes and large datasets. In contrast, alignment-free approaches such as k-mer analysis enable efficient and scalable genome comparison while preserving informative sequence patterns.

Recent advances in bioinformatics and machine learning have enabled automated feature extraction and predictive modelling from genomic data. However, accurate modelling depends on well-curated, non-redundant datasets and biologically meaningful feature representation.

In this study, a curated dataset is constructed of complete genomes from four Orthopoxvirus species and extract sequence-based features

including nucleotide composition, k-mer frequencies, and dinucleotide bias. Using these features, two supervised machine learning models are developed: one for multi-class species classification and another for virulence prediction. Additionally, feature importance analysis is performed to identify genomic determinants that contribute to species differentiation and variations in pathogenic potential.

## II. METHODOLOGY

### A. Data Collection

For this project, genome sequences of four Orthopoxvirus species — Variola, Monkeypox, Cowpox, and Vaccinia — were retrieved from the NCBI GenBank database in FASTA format.

The sequences were retrieved using a Python script with the Biopython library (Entrez module), which automated querying the NCBI database and downloading sequences in FASTA format. A total of 250 sequences per virus were collected, resulting in 1,000 sequences overall. Both complete and partial nucleotide sequences were included to capture the diversity of strains, covering multiple years and geographic regions.

As a result of this process, the sequences were successfully collected in a structured and organized dataset, ready for further analysis

### B. Data Curation

To ensure a non-redundant dataset, we used a k-mer based cosine similarity method, which is a commonly used approach in bioinformatics for identifying highly similar nucleotide sequences. This method converts each sequence into a vector of k-mer counts — short sub sequences of length k — allowing sequences to be compared numerically rather than by direct alignment. In this study, k = 6 was chosen. The cosine similarity between each new sequence and previously kept sequences is calculated, and sequences with similarity above a predefined threshold of 0.992 (for Variola, Cowpox, and Vaccinia) are considered redundant and excluded. Monkeypox sequences were treated differently: all sequences were initially kept due to the low diversity observed in preliminary analyses. The program reads sequences from FASTA files using Biopython's SeqIO module, converts

sequences to vectors, performs similarity comparisons, and writes the nonredundant sequences back to new FASTA files for each virus.

### C. Dataset Specificity Validation

To confirm that the curated sequences were specific to each virus, a few representative sequences from each dataset were checked using NCBI BLASTn. BLAST compares a query sequence against a database to find similar sequences, allowing us to verify that the sequences matched primarily with the intended virus.

For this validation, five sequences per virus were selected from the non-redundant dataset to cover different strains and geographic regions. These sequences were submitted as queries in FASTA format, and the BLAST results were analyzed to ensure that the top hits corresponded to the target virus.

### D. Feature Extraction

After completing dataset curation and validation, the next step was to extract informative features from each complete Orthopoxvirus genome so that meaningful comparisons could be made across species. Because the study aims to classify species and predict virulence categories, the feature extraction process was designed to capture characteristics of the genome that reflect evolutionary differences and potential functional variation.

- *Mononucleotide Composition*

Mononucleotide composition refers to the relative proportion of A,T,G,C within each complete genome. For every genome sequence, the total occurrences of each nucleotide were counted and normalized by the genome length to obtain percentage values.

These features provide a basic structural profile of each genome. Variations in nucleotide composition may reflect evolutionary divergence and host adaptation patterns, which can assist the ML models in distinguishing between different viral species.

- *GC Content and Strand Asymmetry*

GC content was calculated as the proportion of guanine and cytosine bases relative to total genome length. GC skew

and AT skew were computed using strand asymmetry formulas.

GC content relates to genome stability and mutation tendencies, while skew values reflect strand-level bias. Together, these metrics add additional discriminatory information and help capture evolutionary differences across species.

- *Dinucleotide Frequencies*

All dinucleotide combinations were quantified using overlapping two-base windows across each genome and normalized by sequence length.

Dinucleotide frequencies capture short-range nucleotide dependencies beyond single-base composition. These patterns often vary across viral species and may reflect adaptation to host immune pressure, making them useful for classification and virulence analysis.

- *Observed-to-Expected Ratios*

For each dinucleotide, an observed-to-expected ratio was calculated using mononucleotide-based expectations. This measure highlights enrichment or suppression of specific dinucleotides. Such suppression patterns are linked to immune recognition and provide biologically relevant signals for modeling virulence.

- *Trinucleotide (3-mer) Frequencies*

All trinucleotide combinations were enumerated using a sliding window approach. Since trinucleotides correspond to codons, their distribution reflects codon usage patterns.

Codon bias can indicate host adaptation and translational efficiency. Differences in these patterns assist in distinguishing species and may relate to differences in viral fitness.

- *Tetranucleotide (4-mer) Frequencies*

All four-base combinations were calculated using overlapping windows across the genome.

4-mers capture more specific local motifs beyond codon-level information. These features improve resolution when

distinguishing closely related viral genomes.

- *k-mer Entropy*

Entropy was also computed at the k-mer level to assess motif diversity.

This metric captures local sequence complexity and complements basic compositional features.

- *Genome Length*

The total nucleotide length of each complete genome was recorded.

Genome length differences may arise from insertions or deletions and contribute to structural variation between species.

- *Repeat Density*

Repetitive regions were identified by detecting recurring sequence patterns and estimating their proportional representation.

Repeat density reflects aspects of genome organization and may vary across Orthopoxvirus species.

- **BERT Model**

To complement compositional features, a transformer-based model, DNABERT, was used to generate contextual genome representations. DNABERT is an adaptation of the BERT architecture for DNA sequences, where overlapping k-mers are treated as tokens. In this study, complete Orthopoxvirus genomes were segmented into overlapping k-mers and provided as input to the pretrained DNABERT model. Through self-attention layers, the model captures contextual relationships and long-range dependencies within the genome. The final hidden layer embeddings were aggregated using mean pooling to produce a fixed-length vector (e.g., 768 dimensions) representing each genome. Unlike frequency-based features, these embeddings encode deeper sequence patterns that may help distinguish closely related Orthopoxvirus species and capture subtle genomic signals associated with virulence

variation. The resulting embedding vectors were used as input to both the species classification and virulence prediction models.

### III. IMPLEMENTATION DETAILS

#### A. Dataset Retrieval

```

1  from Bio import Entrez, SeqIO
2
3  Entrez.email = "wby1606@gmail.com"
4
5  viruses = [
6      "Variola virus",
7      "Monkeypox virus",
8      "Cowpox virus",
9      "Vaccinia virus"
10 ]
11
12 sequences_per_virus = 250
13
14 for virus in viruses:
15     print(f"\nSearching for {virus}...")
16
17     search_handle = Entrez.esearch(
18         db="nucleotide",
19         term=f"(virus)[Organism]",
20         retmax=sequences_per_virus
21     )
22     search_results = Entrez.read(search_handle)
23     ids = search_results["IdList"]
24
25     print(f"Found {len(ids)} sequences")
26
27     if len(ids) == 0:
28         print("No sequences found, skipping.")
29         continue
30
31     fetch_handle = Entrez.efetch(
32         db="nucleotide",
33         id=ids,
34         rettype="fasta",
35         retmode="text"
36     )
37
38     sequences = list(SeqIO.parse(fetch_handle, "fasta"))
39
40     filename = virus.replace(" ", "_") + ".fasta"
41     SeqIO.write(sequences, filename, "fasta")
42
43     print(f"Saved to {filename}")

```

---

```

from Bio import SeqIO
import numpy as np
from itertools import product

input_files = {
    "variola": "Orthopox_dataset/Variola_virus.fasta",
    "monkeypox": "Orthopox_dataset/Monkeypox_virus.fasta",
    "cowpox": "Orthopox_dataset/Cowpox_virus.fasta",
    "vaccinia": "Orthopox_dataset/Vaccinia_virus.fasta"
}

k = 6
thresholds = {
    "variola": 0.992,
    "monkeypox": 0.992,
    "cowpox": 0.992,
    "vaccinia": 0.992
}

def generate_kmers(k):
    return ''.join(p) for p in product('ATGC', repeat=k)

all_kmers = generate_kmers(k)
kmer_index = {kmer: idx for idx, kmer in enumerate(all_kmers)}

def seq_to_vector(seq):
    vector = np.zeros(len(all_kmers), dtype=np.float32)
    for i in range(len(seq) - k + 1):
        kmer = seq[i:i+k]
        if kmer in kmer_index:
            vector[kmer_index[kmer]] += 1
    norm = np.linalg.norm(vector)
    if norm > 0:
        vector = vector / norm
    return vector

for virus_name, file_path in input_files.items():
    records = list(SeqIO.parse(file_path, "fasta"))

    if virus_name.lower() == "monkeypox":
        kept_records = records
        removed_count = 0
        non_redundant_count = len(records)
    else:
        kept_records = []
        genome_vectors = []
        threshold = thresholds[virus_name]

        for record in records:
            seq = str(record.seq).upper()
            vec = seq_to_vector(seq)

            is_redundant = False
            for existing_vec in genome_vectors:
                sim = np.dot(vec, existing_vec)
                if sim >= threshold:
                    is_redundant = True
                    break

            if not is_redundant:
                kept_records.append(record)
                genome_vectors.append(vec)

        non_redundant_count = len(kept_records)
        removed_count = len(records) - non_redundant_count

    print("----", virus_name.upper(), "----")
    print("Original sequences : ", len(records))
    print("Removed as redundant : ", removed_count)
    print("Non-redundant obtained : ", non_redundant_count)
    print()

    output_file = virus_name + "_final.fasta"
    SeqIO.write(kept_records, output_file, "fasta")

```

#### B. Data Curation

#### C. Dataset Specificity Validation

```

1  from Bio import SeqIO
2  import numpy as np
3  from sklearn.metrics.pairwise import cosine_similarity
4  from itertools import product
5  import os
6
7  input_files = {
8      "variola": "variola_virus.fasta",
9      "monkeypox": "monkeypox_virus.fasta",
10     "cowpox": "cowpox_virus.fasta",
11     "vaccinia": "vaccinia_virus.fasta"
12 }
13
14 k = 3
15 similarity_threshold = 0.98
16
17 def generate_kmers(k):
18     return [''.join(p) for p in product('ATGC', repeat=k)]
19
20 all_kmers = generate_kmers(k)
21 kmer_index = {kmer: idx for idx, kmer in enumerate(all_kmers)}
22
23 def seq_to_vector(seq):
24     vector = np.zeros(len(all_kmers))
25     for i in range(len(seq) - k + 1):
26         kmer = seq[i:i+k]
27         if kmer in kmer_index:
28             vector[kmer_index[kmer]] += 1
29     return vector
30
31 for virus_name, file_path in input_files.items():
32
33     records = list(SeqIO.parse(file_path, "fasta"))
34     vectors = []
35     kept_records = []
36
37     for record in records:
38         vec = seq_to_vector(str(record.seq).upper())
39
40         if len(vectors) == 0:
41             vectors.append(vec)
42             kept_records.append(record)
43         else:
44             sims = cosine_similarity([vec], vectors)[0]
45             if max(sims) < similarity_threshold:
46                 vectors.append(vec)
47                 kept_records.append(record)
48
49     print("-----", virus_name.upper(), "-----")
50     print("Original:", len(records))
51     print("Non-redundant:", len(kept_records))
52     print()
53
54     output_file = virus_name + ".nonredundant.fasta"
55     SeqIO.write(kept_records, output_file, "fasta")

```

## IV. RESULTS AND ANALYSIS

### A. Data Retrieval Results

Searching for Variola virus...  
 Found 250 sequences  
 Saved to Variola\_virus.fasta

Searching for Monkeypox virus...  
 Found 250 sequences  
 Saved to Monkeypox\_virus.fasta

Searching for Cowpox virus...  
 Found 250 sequences  
 Saved to Cowpox\_virus.fasta

Searching for Vaccinia virus...  
 Found 250 sequences  
 Saved to Vaccinia\_virus.fasta

### B. Dataset Curation Results

```

----- VARIOLA -----
Original sequences      : 250
Removed as redundant   : 184
Non-redundant obtained : 66

----- MONKEYPOX -----
Original sequences      : 250
Removed as redundant   : 0
Non-redundant obtained : 250

----- COWPOX -----
Original sequences      : 250
Removed as redundant   : 167
Non-redundant obtained : 83

----- VACCINIA -----
Original sequences      : 250
Removed as redundant   : 130
Non-redundant obtained : 120

```

The redundancy removal showed differences between the viruses. Variola had many similar sequences, leaving only a small number of non-redundant sequences. In contrast, all Monkeypox sequences were retained, as preliminary analysis showed low redundancy and high diversity among strains.

Cowpox and Vaccinia displayed moderate redundancy, with a substantial number of unique sequences remaining.

This process effectively reduced duplication while maintaining sequence diversity, resulting in a curated dataset suitable for further analyses.

### C. Dataset Specificity Validation Results

Description	Scientific Name	Max Score	Total Score	Query Cover	E value	Per. Ident	Acc. Len	Accession
✓ Variola virus isolate VARV_AFRICA_1964_LUK_complete genome	Variola.virus	200	200	100%	3e-47	100.00%	186197	PR405584.1
✓ Tateropox virus_complete genome	Tateropox....	200	200	100%	3e-47	100.00%	198000	NC_002911
✓ Variola virus strain Somalia 1977_complete genome	Variola.virus	200	200	100%	3e-47	100.00%	186231	DQ417590.1
✓ Variola virus strain China Horn 1948_complete genome	Variola.virus	200	200	100%	3e-47	100.00%	186668	DQ417592.1
✓ Camelpox virus_CMS_complete genome	Camelpox....	200	200	100%	3e-47	100.00%	202205	AY009098.1
✓ Orthopoxvirus camelpox strain Camelpox virusK2/Beineau2023_compo_Orthopoxv...	Orthopoxv...	200	200	100%	3e-47	100.00%	202273	PY020573.1
✓ Variola virus strain Pakistan 1969 (Rally Lekene)_complete sequence	Variola.virus	200	200	100%	3e-47	100.00%	186595	PR405590.1
✓ Variola virus strain Ethiopia 1972 (Eth17 R14-X72 Addis)_complete ...	Variola.virus	200	200	100%	3e-47	100.00%	186640	DQ414251
✓ Variola virus strain India 1953 (New Delhi)_complete genome	Variola.virus	200	200	100%	3e-47	100.00%	186662	DQ414281
✓ Camelpox virus strain D1795/20_complete genome	Camelpox....	200	200	100%	3e-47	100.00%	201938	MZ308583.1
✓ Variola virus isolate VARV_64-7_complete genome	Variola.virus	200	200	100%	3e-47	100.00%	186216	PR405594.1
✓ Cowpox virus isolate Brighton_Red-1956/227_complete genome	Cowpox.vi...	920	5849	100%	0.0	100.00%	221599	PX585982.1
✓ Cowpox virus_complete genome	Cowpox.vi...	835	12975	94%	0.0	98.33%	224499	NC_003653.2
✓ Cowpox virus CPXV/Bos taurus/Basrah/ASRAH/2024 DNA_compl..._Cowpox.vi...	Cowpox.vi...	835	12948	94%	0.0	98.33%	LC889885.1	
✓ Mutant Cowpox virus strain recombinant Brighton Red clone deltaCPXV..._Cowpox.vi...	Cowpox.vi...	835	12975	94%	0.0	98.33%	220600	QH599271.1
✓ Cowpox virus BASRH/2 DNA_complete genome	Cowpox.vi...	835	12742	94%	0.0	98.33%	224499	LC890709.1
✓ Cowpox virus BASRH/BIF/2 DNA_complete genome	Cowpox.vi...	835	12874	94%	0.0	98.33%	224499	LC890707.1
✓ Cowpox virus BASRH/BIF/2 DNA_complete genome	Cowpox.vi...	835	12874	94%	0.0	98.33%	224499	LC890708.1
✓ Cowpox virus isolate CPXV_K4207_complete genome	Cowpox.vi...	830	12864	94%	0.0	98.12%	221439	KY549150.2
✓ Cowpox virus isolate CPXV_Catbox3/97_complete genome	Cowpox.vi...	819	10743	100%	0.0	97.70%	216799	KY549143.2
✓ Cowpox virus isolate CPXV_Catbox5v1_complete genome	Cowpox.vi...	817	13084	100%	0.0	97.51%	224952	KY549144.1
✓ Cowpox virus isolate CPXV_1639_complete genome	Cowpox.vi...	815	12925	100%	0.0	97.50%	224936	KY549148.1

<input checked="" type="checkbox"/> Monkeypox virus isolate MPXV/Germany/2022/ON/RK062, partial genome	Monkeypox...	924	924	100%	0.0	100.00%	190602	QG274881.1
<input checked="" type="checkbox"/> Monkeypox virus isolate MoxV/Human/USA/WA-JW-089807/2022, partial genome	Monkeypox...	924	1848	100%	0.0	100.00%	197192	OP528486.1
<input checked="" type="checkbox"/> Monkeypox virus isolate NY-NCPH-001311, partial genome	Monkeypox...	924	1848	100%	0.0	100.00%	195665	PP751827.1
<input checked="" type="checkbox"/> Monkeypox virus isolate MPXV/Human/USA/CA-LACPHL-MA00079/2022, partial genome	Monkeypox...	924	1848	100%	0.0	100.00%	195594	OP442540.1
<input checked="" type="checkbox"/> Monkeypox virus isolate MPXV/UZ_REGA_162/Belgium/2022, partial genome	Monkeypox...	924	1848	100%	0.0	100.00%	197425	QR586971.1
<input checked="" type="checkbox"/> Monkeypox virus isolate MPXV_USA_2022_LA0036, partial genome	Monkeypox...	924	1848	100%	0.0	100.00%	193707	PV584927.1
<input checked="" type="checkbox"/> Monkeypox virus isolate MoxV/US/ACA-COPH-1M/100041/2022, partial genome	Monkeypox...	924	1848	100%	0.0	100.00%	196057	PQ279975.1
<input checked="" type="checkbox"/> Monkeypox virus isolate MoxV/US/ACA-COPH-1M/100041/2022, partial genome	Monkeypox...	924	1848	100%	0.0	100.00%	195603	PQ483292.1
<input checked="" type="checkbox"/> Monkeypox virus isolate MPXV/US/AL-RIPHL-050-0202/2024, partial genome	Monkeypox...	924	1848	100%	0.0	100.00%	195595	OP515214.1
<input checked="" type="checkbox"/> Monkeypox virus isolate MPXV/US/AL-RIPHL-050-0202/2024, complete genome	Monkeypox...	924	1848	100%	0.0	100.00%	197205	OP515214.1
<input checked="" type="checkbox"/> Monkeypox virus isolate MPXV/UZ_REGA_143/Belgium/2022, partial genome	Monkeypox...	924	924	100%	0.0	100.00%	194440	OR886952.1
<input checked="" type="checkbox"/> Monkeypox virus isolate MoxV/US/ACA-COPH-1M/100041/2022, partial genome	Monkeypox...	924	1848	100%	0.0	100.00%	195609	PQ279972.1
<input checked="" type="checkbox"/> Monkeypox virus isolate MoxV/US/ACA-COPH-1M/100041/2022, partial genome	Monkeypox...	924	1848	100%	0.0	100.00%	195595	OP524545.1
<input checked="" type="checkbox"/> Vaccinia virus strain Dryvax clone DPP13, complete genome	Orthopoxv...	918	918	100%	0.0	99.80%	194800	JN654980.1
<input checked="" type="checkbox"/> Vaccinia virus strain 3737, complete genome	Orthopoxv...	918	918	100%	0.0	99.80%	199377	QG277945.1
<input checked="" type="checkbox"/> Vaccinia virus strain Dryvax clone DPP19, complete genome	Orthopoxv...	918	918	100%	0.0	99.80%	198609	JN654981.1
<input checked="" type="checkbox"/> Vaccinia virus strain Dryvax clone DPP12, complete genome	Orthopoxv...	913	913	100%	0.0	99.60%	198741	JN654979.1
<input checked="" type="checkbox"/> Vaccinia virus strain Acambia clone 3, complete genome	Orthopoxv...	913	913	100%	0.0	99.60%	197418	AY313848.1
<input checked="" type="checkbox"/> Vaccinia virus isolate VACV-Wyeth_A211, complete genome	Orthopoxv...	911	911	100%	0.0	99.60%	199801	OP751802.1
<input checked="" type="checkbox"/> Vector syn/VACV-SFV, complete sequence	Vector_syn...	907	907	100%	0.0	99.40%	199901	MW990419.1
<input checked="" type="checkbox"/> Vector syn/VACV-Delta1-3, complete sequence	Vector_syn...	907	907	100%	0.0	99.40%	200013	MW990422.1
<input checked="" type="checkbox"/> Orthopoxvirus vaccinia strain Acambia/2013, complete genome	Orthopoxv...	907	907	100%	0.0	99.40%	199150	MT227314.1
<input checked="" type="checkbox"/> Vaccinia virus strain Acambia clone 2000, complete genome	Orthopoxv...	907	907	100%	0.0	99.40%	199234	AY313847.1
<input checked="" type="checkbox"/> Vaccinia virus strain Dryvax clone DPP18, complete genome	Orthopoxv...	907	907	100%	0.0	99.40%	198820	JN654982.1

The BLAST validation confirmed that the curated sequences are specific to their respective viruses.

For each representative sequence:

- Variola sequences matched primarily with other Variola sequences, showing 100% identity and full query coverage. Only a few hits aligned with closely related Orthopoxviruses like Taterapox or Camelpox, reflecting natural sequence similarity within the genus.
- Cowpox sequences showed high identity (around 98–100%) with other Cowpox sequences, confirming their specificity.
- Monkeypox sequences matched exclusively with other Monkeypox sequences, showing low redundancy and high diversity within the species.
- Vaccinia sequences aligned strongly with other Vaccinia sequences, showing high identity and full coverage, with minimal alignment to unrelated sequences.

Across all viruses, no significant hits were observed with unrelated organism. These results demonstrate that the nonredundant sequences are highly specific to each virus.