# RBP-ID

## Development of a Web-Based Bioinformatics Tool for the Identification of Phage Receptor-Binding Proteins

Presented by Cátia Rosário
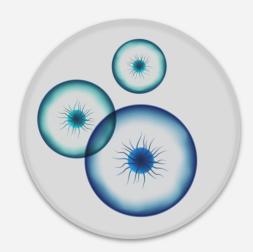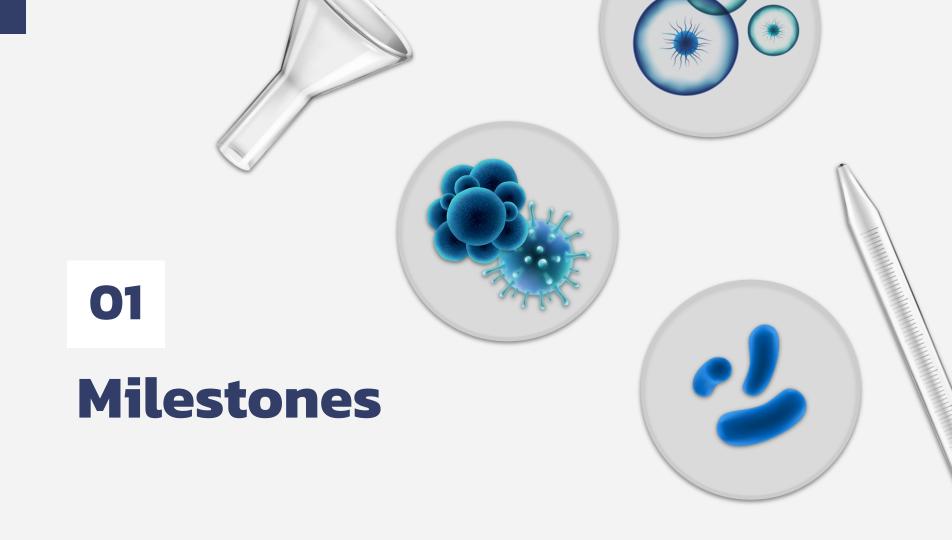Oriented by Silvio Santos and Óscar Dias

# Table of contents

# 01

# Milestones

# Milestones reached
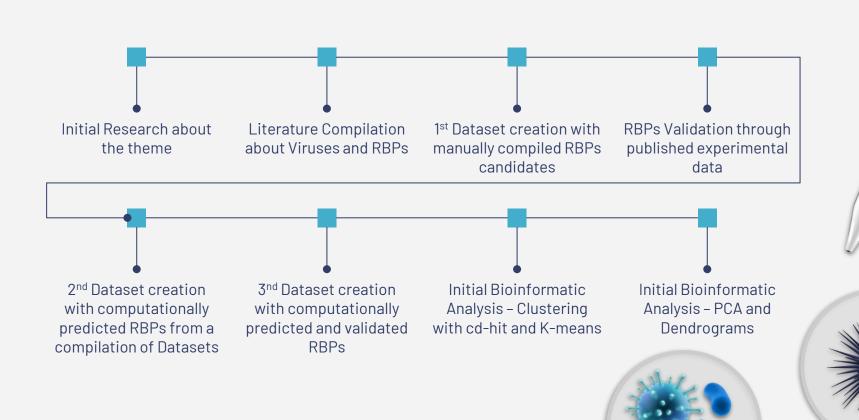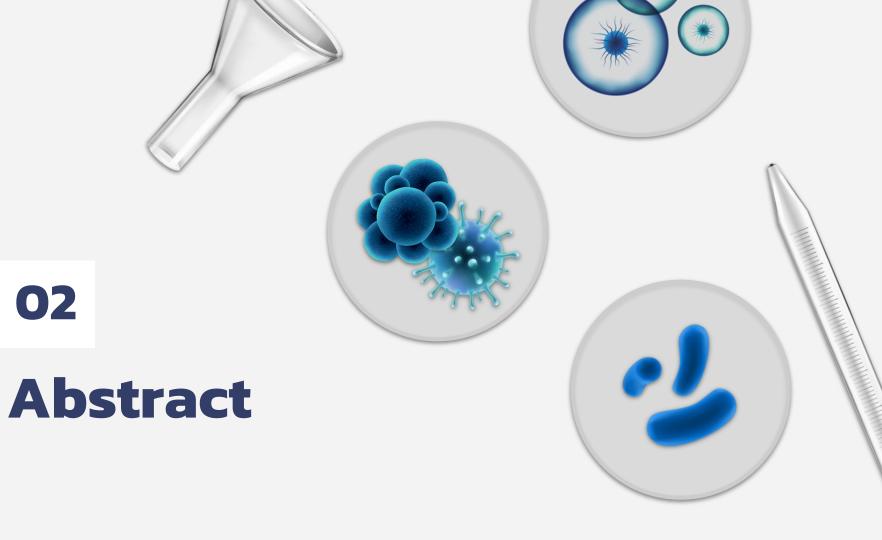
Initial Research about the theme

Literature Compilation about Viruses and RBPs

1st Dataset creation with manually compiled RBPs candidates

RBPs Validation through published experimental data

2nd Dataset creation with computationally predicted RBPs from a compilation of Datasets

3nd Dataset creation with computationally predicted and validated RBPs

Initial Bioinformatic Analysis – Clustering with cd-hit and K-means

Initial Bioinformatic Analysis – PCA and Dendrograms

**02**

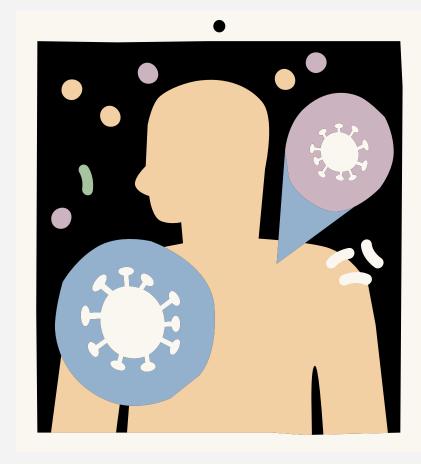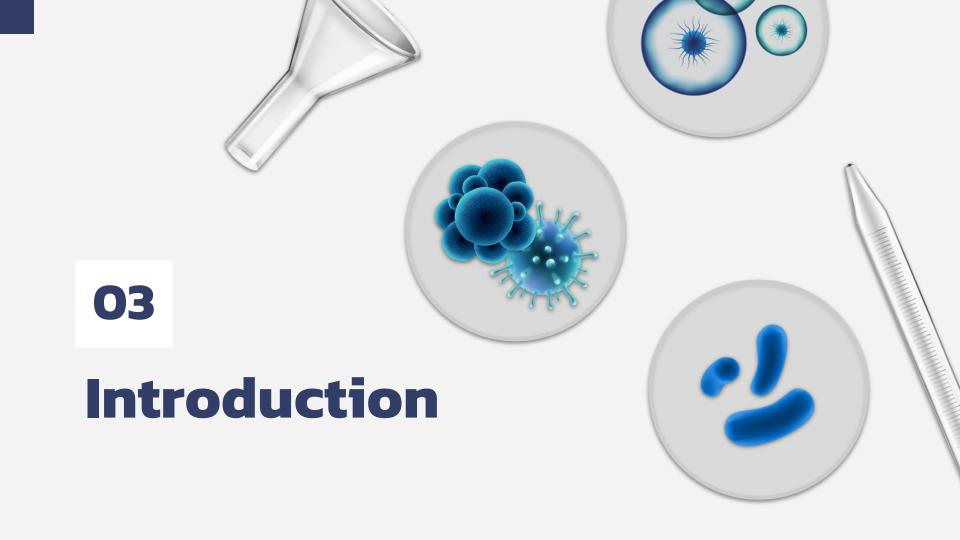# Abstract

# Abstract

- **Bacteriophages:** Viruses that infect bacteria with key roles in biotechnology and medicine.

- **Receptor-Binding Proteins (RBPs):** Determine host specificity by mediating bacterial recognition and attachment.

- **Challenge:** RBPs exhibit high sequence diversity and low conservation, making identification difficult.

- **Solution – RBP-ID:** A web-based bioinformatics tool for systematic RBP prediction in phage genomes.

03

# Introduction

# Introduction

## Background

**Phage display & delivery systems:** Technologies for gene editing and therapeutic biomolecule delivery.

**Adaptation mechanisms:** Co-evolution with bacteria through receptor mutations and CRISPR-Cas defence systems.

**Computational tools for phage genome annotation:**
- Hybrid approaches: HMM searches and ML classifiers
- Host prediction via RBP analysis: Random Forest and SVMs, XGBoost
- PHYPred: Functional annotation tools for phage-encoded enzymes
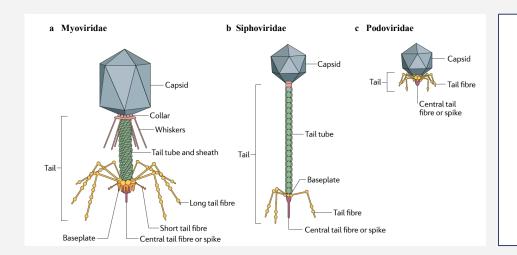- PhANNs

## Importance

**Biotechnological & therapeutic:** Antibiotic resistance solutions, microbiome engineering, biosensing and precision medicine

**Phages vs. antibiotics:** Target specific pathogens without harming beneficial bacteria, a key advantage in antimicrobial resistance (AMR)

**Synthetic biology advances:** Engineered phages and phage-derived enzymes broaden applications in medicine, agriculture, and industry

**Need for advanced bioinformatics tools**

a  **Myoviridae**

Capsid
Collar
Whiskers
Tail tube and sheath
Tail
Long tail fibre
Short tail fibre
Baseplate
Central tail fibre or spike

b  **Siphoviridae**

Capsid
Tail tube
Tail
Baseplate
Tail fibre
Central tail fibre or spike

c  **Podoviridae**

Capsid
Tail
Tail fibre
Central tail fibre or spike

# Morphotypes

# Tail Fibre World

# Short, Long, Spikey, Contractile, Non-Contractile

# HOW DO VIRUSES REPLICATE

The virus's genetic material reprograms the host's DNA and starts to replicate itself

When the host cell is filled with virus copies, it bursts, releasing them into the bloodstream to repeat the process

The virus invades the host cell

The virus injects its genetic material into the host cell

**04**

# Results

**4.1**
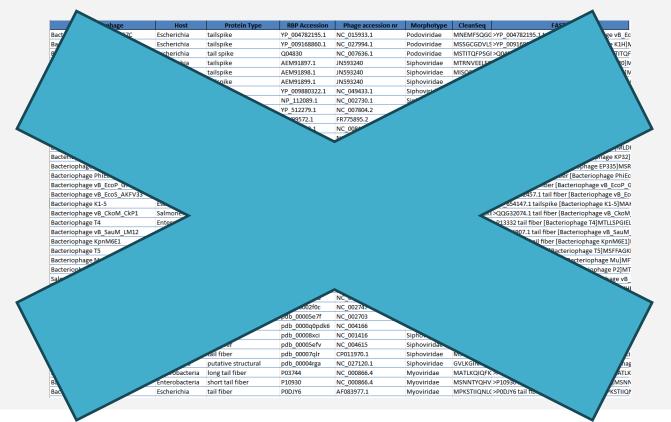
# 1st Dataset

# 1st Dataset – Validated RBPs

```python
import pandas as pd
from Bio import Entrez, SeqIO
import time

Entrez.email = "catiarosario10@gmail.com"
excel_path = "RBP_validated_final.xlsx"
sheet_name = "Sheet1"
phage_id_col = "Bacteriophage"
protein_id_col = "RBP Accession"
desc_col = "Protein Type"
source_col = "Bacteriophage"
seq_col = "RBP sequence (aa)"
output_fasta_file = "validated_RBP.fasta"
output_excel_file = "validated_processed.xlsx"

def clean_sequence(seq):
    if not isinstance(seq, str):
        return ""
    lines = seq.strip().splitlines()
    lines = [line for line in lines if not line.startswith(">")]
    return "".join(lines).replace(" ", "").upper()

def make_fasta_header(row):
    desc = row.get(desc_col) if pd.notna(row.get(desc_col)) else ""
    source = row.get(source_col) if pd.notna(row.get(source_col)) else ""
    return f">{row[protein_id_col]} {desc} [{source}]".strip()

def make_full_fasta(row):
    header = make_fasta_header(row)
    seq = row["CleanSeq"]
    return f"{header}\n{seq}"

def get_taxid_and_accession(phage_name):
    try:
        search = Entrez.esearch(db="nucleotide", term=phage_name, retmode="xml")
        record = Entrez.read(search)
        search.close()
        if not record["IdList"]:
            return None, None
        accession_id = record["IdList"][0]
        handle = Entrez.efetch(db="nucleotide", id=accession_id, rettype="gb", retmode="text")
        seq_record = SeqIO.read(handle, "genbank")
        handle.close()
        taxid = None
        for feature in seq_record.features:
            if feature.type == "source":
                for ref in feature.qualifiers.get("db_xref", []):
                    if ref.startswith("taxon:"):
                        taxid = ref.split(":")[1]
                        break
        return taxid, accession_id
    except Exception as e:
        print(f"[{phage_name}] TaxID/Accession fetch error: {e}")
        return None, None
```

```python
df = pd.read_excel(excel_path, sheet_name=sheet_name)
df = df.dropna(subset=[phage_id_col, protein_id_col, seq_col])

df["CleanSeq"] = df[seq_col].apply(clean_sequence)
df["FASTA"] = df.apply(make_full_fasta, axis=1)

seq_to_phageids = {}
for _, row in df.iterrows():
    seq = row["CleanSeq"]
    phageid = row[phage_id_col]
    seq_to_phageids.setdefault(seq, []).append(phageid)

def get_duplicate_phageids(row):
    seq = row["CleanSeq"]
    phageid = row[phage_id_col]
    duplicates = [pid for pid in seq_to_phageids.get(seq, []) if pid != phageid]
    return ", ".join(duplicates) if duplicates else ""

df["DuplicatePhageIDs"] = df.apply(get_duplicate_phageids, axis=1)

df_unique = df.drop_duplicates(subset=["CleanSeq"]).copy()

required_cols = ["Bacteriophage", "Phage accession nr", "Protein Type", "Host", "RBP Accession", "Morphotype"]
missing_cols = [col for col in required_cols if col not in df_unique.columns]
if missing_cols:
    raise ValueError(f"Missing columns: {missing_cols}")

df_unique["AccessionFetched"] = df_unique["Phage accession nr"]
df_unique["Protein Type"] = df_unique["Protein Type"]
df_unique["Bacteriophage"] = df_unique["Bacteriophage"]

with open(output_fasta_file, "w", encoding="utf-8") as fasta_out:
    for fasta_entry in df_unique["FASTA"]:
        fasta_out.write(fasta_entry + "\n")

cols_to_save = ["Bacteriophage", "Host", "Protein Type", "RBP Accession", "Phage accession nr", "Morphotype", "CleanSeq", "FASTA"]

df_unique.to_excel(output_excel_file, columns=cols_to_save, index=False)
print(f"FASTA saved to '{output_fasta_file}'")
print(f"Excel saved to '{output_excel_file}'")
```

**Workflow data processing - VALIDATED:**
- Add FASTA format when missing
- Create a column with only aa to remove duplicates
- Create column with FASTA format to create FASTA file
- Create column with duplicated identifiers when found
- Remove duplicate lines based on sequence
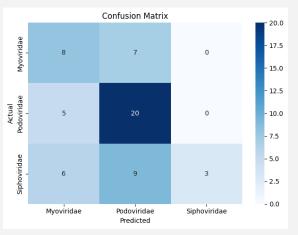- Create FASTA file and Excel file

# Validated RBPs – Processement with Python

**Validated RBPs after Processement with Python**

# Validated RBPs Analysis K-means and Confusion Matrix

```python
import numpy as np
from sklearn.cluster import KMeans
from sklearn.feature_extraction.text import TfidfVectorizer
from sklearn.metrics import silhouette_score, classification_report
import matplotlib.pyplot as plt
import seaborn as sns

# === Load Data ===
df = pd.read_excel("validated_processed.xlsx")
sequences = df["CleanSeq"].astype(str).tolist()

# === Vectorize Sequences with TF-IDF 3-4-mers ===
vectorizer = TfidfVectorizer(analyzer='char', ngram_range=(3, 4))
X = vectorizer.fit_transform(sequences)

# === Find Optimal K Using Silhouette Score ===
best_k = None
best_score = -1
score_dict = {}

print("\nFinding optimal number of clusters...")
for k in range(2, 10):
    kmeans = KMeans(n_clusters=k, random_state=42, n_init='auto')
    labels = kmeans.fit_predict(X)
    score = silhouette_score(X, labels)
    score_dict[k] = score
    print(f"Silhouette score for k={k}: {score:.4f}")
    if score > best_score:
        best_score = score
        best_k = k

print(f"\nBest number of clusters: {best_k} (score={best_score:.4f})")

# === Final KMeans with Best K ===
kmeans = KMeans(n_clusters=best_k, random_state=42, n_init='auto')
df["BestKMeans_Cluster"] = kmeans.fit_predict(X)

# === Majority Vote Mapping from Cluster to Morphotype ===
def assign_morphotypes_by_majority(df, cluster_col, morph_col):
    mapping = {}
    for cluster in df[cluster_col].unique():
        subset = df[df[cluster_col] == cluster]
        most_common = subset[morph_col].mode()
        if not most_common.empty:
            mapping[cluster] = most_common.iloc[0]
        else:
            mapping[cluster] = "Unknown"
    return mapping

def apply_cluster_to_morphotype_mapping(df, cluster_col, mapping):
    return df[cluster_col].map(mapping).fillna("Unknown")

cluster_to_morph = assign_morphotypes_by_majority(df, "BestKMeans_Cluster", "Morphotype")
df["Predicted_morphotype"] = apply_cluster_to_morphotype_mapping(df, "BestKMeans_Cluster", cluster_to_morph)

print(f"\nCluster to Morphotype Mapping:\n{cluster_to_morph}")

# === Evaluate ===
valid_rows = df["Predicted_morphotype"].notna()
print("\nClassification Report:")
print(classification_report(df.loc[valid_rows, "Morphotype"], df.loc[valid_rows, "Predicted_morphotype"], zero_division=0))

# === Confusion Matrix ===
conf_mat = pd.crosstab(df["Morphotype"], df["Predicted_morphotype"], rownames=["Actual"], colnames=["Predicted"])
sns.heatmap(conf_mat, annot=True, cmap="Blues", fmt='d')
plt.title("Confusion Matrix")
plt.tight_layout()
plt.show()

# === Save Results ===
df.to_excel("phage_clustered_with_morphotypes.xlsx", index=False)
print("\nSaved results to 'phage_clustered_with_morphotypes.xlsx'")
```

**Performance Analysis:**

Podoviridae:

Highest recall and f1-score
↓
Reliably identified

Myoviridae and Siphoviridae:

Weaker performance;
Low recall for Siphoviridae
↓
Potential issues in classification

Classification Report:

|  | precision | recall | f1-score | support |
|---|---|---|---|---|
| Myoviridae | 0.42 | 0.53 | 0.47 | 15 |
| Podoviridae | 0.56 | 0.80 | 0.66 | 25 |
| Siphoviridae | 1.00 | 0.17 | 0.29 | 18 |
| | | | | |
| accuracy | | | 0.53 | 58 |
| macro avg | 0.66 | 0.50 | 0.47 | 58 |
| weighted avg | 0.66 | 0.53 | 0.49 | 58 |

```
Morphotype
Podoviridae      25
Siphoviridae     18
Myoviridae       15
Name: count, dtype: int64
Predicted_morphotype
Podoviridae      36
Myoviridae       19
Siphoviridae      3
Name: count, dtype: int64
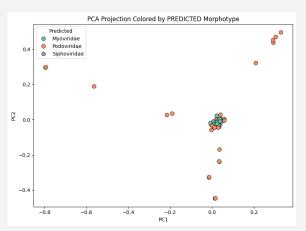```

# Validated RBPs – PCA and Clusters

```python
# === PCA Visualization ===
pca = PCA(n_components=2)
X_pca = pca.fit_transform(X.toarray())

# Add coordinates to dataframe
df["PC1"] = X_pca[:, 0]
df["PC2"] = X_pca[:, 1]

# === Plot PCA Results ===
fig, axes = plt.subplots(1, 2, figsize=(16, 6))  # 1 row, 2 columns

# Plot by True Morphotype
sns.scatterplot(data=df, x="PC1", y="PC2", hue="Morphotype", palette="Set2",
                s=60, edgecolor='k', ax=axes[0])
axes[0].set_title("PCA Projection Colored by TRUE Morphotype")
axes[0].legend(title="Morphotype", loc='best')

# Plot by Predicted Morphotype
sns.scatterplot(data=df, x="PC1", y="PC2", hue="Predicted_morphotype", palette="Set2",
                s=60, edgecolor='k', ax=axes[1])
axes[1].set_title("PCA Projection Colored by PREDICTED Morphotype")
axes[1].legend(title="Predicted", loc='best')

plt.tight_layout()
plt.show()
```

**PCA - TRUE Morphotype Clusters**:
**Podoviridae**: Well-separated
**Siphoviridae**: Cluster closer to Myoviridae
**Myoviridae**: Not clearly distinct; Slightly overlapping with Siphoviridae

- Distinct clustering for Podoviridae
- Confusion/blending between Siphoviridae and Myoviridae

**PCA - PREDICTED Morphotype Clusters**:
**Podoviridae**: New clustering observed, likely misclassified from the TRUE representation.
**Siphoviridae:** Predicted model seems less accurate
**Myoviridae**: Some appear clustered together with Podoviridae

- Predictive model not capture morphotype distinctions effectively



PCA Projection Colored by PREDICTED Morphotype



PCA Projection Colored by TRUE Morphotype

# 1ˢᵗ Dataset – Dendrogram K-means

```
# === Dendrogram Visualization ===
# Compute the distance matrix
distance_matrix = pairwise_distances(X.toarray(), metric='cosine')

# Create the linkage matrix
Z = linkage(distance_matrix, method='ward')

# Plotting the dendrogram
plt.figure(figsize=(10, 6))
labels = df.apply(lambda row: f"{row['Morphotype']}_{row['RBP Accession']}", axis=1).values
dendrogram(Z, labels=labels)
plt.title('Dendrogram of Clusters Validated RBPs')
plt.xlabel('Cluster ID')
plt.ylabel('Distance')
plt.show()
```

For global compositional similarity - identifying potential functional groups:

- Cluster shape: Compact, evenly sized clusters
- Separation between clusters: Distinct groups shared sequence features
- Reflecting functional or structural similarity, even among non-homologous proteins - spacing on the left side indicate relationships among RBPs
- Clustering suggests distinct groups among the validated RBPs – some by morphotype


Dendrogram of Clusters Validated RBPs

# 1ˢᵗ Dataset – Validated RBPs CD-HIT

```python
from Bio import SeqIO
import pandas as pd
def parse_cdhit_output(file_path):
    clusters = {}
    current_cluster = None

    with open(file_path, 'r') as file:
        for line in file:
            line = line.strip()
            if line.startswith(">Cluster"):
                current_cluster = line.split(" ")[1]
                clusters[current_cluster] = []
            elif current_cluster is not None and line:
                if '>' in line:
                    try:
                        seq_info = line.split(",")[1] if ',' in line else line
                        seq_id_part = seq_info.split(">")[1]
                        seq_id = seq_id_part.split("...")[0].strip()
                        if seq_id:
                            clusters[current_cluster].append(seq_id)
                    except Exception as e:
                        print(f"Error parsing line: {line}\n{e}")
    return clusters

def load_fasta_sequences(fasta_path):
    seq_dict = {}
    for record in SeqIO.parse(fasta_path, "fasta"):
        core_id = record.id
        seq_dict[core_id] = str(record.seq)
    return seq_dict

def build_sequences_list(clusters, seq_dict):
    sequences = []
    data = []
    for cluster_id, seq_ids in clusters.items():
        for seq_id in seq_ids:
            seq = seq_dict.get(seq_id, "")
            if seq == "":
                print(f"Warning: Sequence ID '{seq_id}' not found in FASTA file.")
            data.append({'Cluster': cluster_id, 'Sequence ID': seq_id})
            sequences.append(seq)
    df = pd.DataFrame(data)
    return df, sequences

if __name__ == "__main__":
    cluster_file = "cluster_validadas.txt"
    fasta_file = "validated_RBP.fasta"

    clusters = parse_cdhit_output(cluster_file)
    seq_dict = load_fasta_sequences(fasta_file)
    df, sequences = build_sequences_list(clusters, seq_dict)

    print(f"Total sequences parsed: {len(sequences)}")

    df.to_csv("clustered_sequences0.csv", index=False)
```
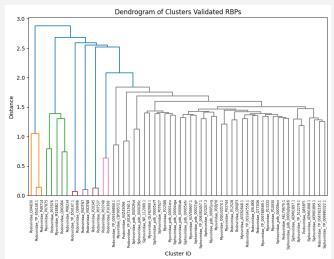
```python
from scipy.cluster.hierarchy import dendrogram, linkage
import numpy as np
from sklearn.feature_extraction.text import TfidfVectorizer

# Vectorize the sequences using TF-IDF
vectorizer = TfidfVectorizer(analyzer='char', ngram_range=(3, 4))
X = vectorizer.fit_transform(sequences)
# Create a linkage matrix using the actual feature matrix
Z = linkage(X.toarray(), method='ward')
# Plotting the dendrogram
plt.figure(figsize=(10, 6))

labels = df.apply(lambda row: f"Cl{row['Cluster']}_{row['Sequence ID']}", axis=1).values
dendrogram(Z, labels=labels)
plt.title('Dendrogram of Clusters Validated RBPs')
plt.xlabel('Cluster ID')
plt.ylabel('Distance')
plt.show()
```
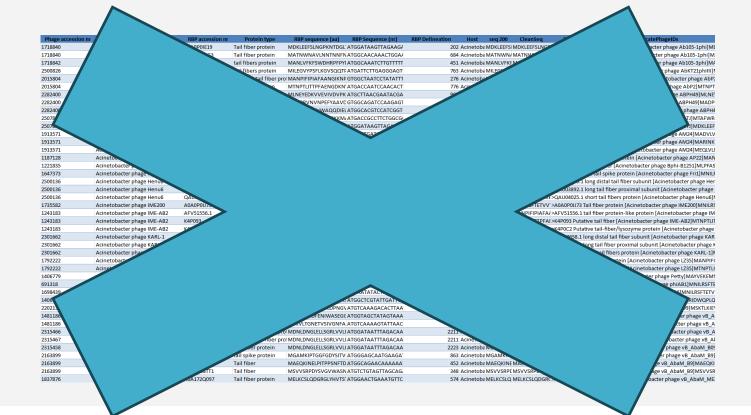
For redundancy removal and homology grouping:
- Pairwise sequence identity
- Cluster shape: Broader clusters
- Separation: Some branches merging at higher distances = low sequence similarity
- Bifurcations at grey represent a bigger diversification in less time comparing with the pink cluster with same ancestry
- Clustering suggests distinct groups among the validated RBPs – some by morphotype



Dendrogram of Clusters Validated RBPs

# 4.2

# 2nd Dataset

**2ⁿᵈ Dataset – Computationally Predicted RBPs**

```python
import pandas as pd
from Bio import Entrez, SeqIO
import time

# === Configuration ===
Entrez.email = "catiarosario10@gmail.com"
excel_path = "RBP_database_NaoValidadas.csv"
sheet_name = "RBP_database_NaoValidadas"
phage_id_col = "Bacteriophage"
protein_id_col = "RBP accession nr"
desc_col = "Protein type"
source_col = "Bacteriophage"
seq_col = "RBP sequence (aa)"
output_fasta_file = "unique_phage_proteins2.fasta"
output_excel_file = "phage_processed2.xlsx"

# === Helper Functions ===

def clean_sequence(seq):
    if not isinstance(seq, str):
        return ""
    lines = seq.strip().splitlines()
    lines = [line for line in lines if not line.startswith(">")]
    return "".join(lines).replace(" ", "").upper()

def make_fasta_header(row):
    desc = row.get(desc_col, "")
    source = row.get(source_col, "")
    desc = desc if pd.notna(desc) else ""
    source = source if pd.notna(source) else ""
    return f">{row[protein_id_col]} {desc} [{source}]".strip()

def make_full_fasta(row):
    header = make_fasta_header(row)
    seq = row["CleanSeq"]
    return f"{header}\n{seq}"
```

```python
def get_taxid_and_accession(phage_name):
    try:
        search = Entrez.esearch(db="nucleotide", term=phage_name, retmode="xml")
        record = Entrez.read(search)
        search.close()
        if not record["IdList"]:
            return None, None
        accession_id = record["IdList"][0]
        handle = Entrez.efetch(db="nucleotide", id=accession_id, rettype="gb", retmode="text")
        seq_record = SeqIO.read(handle, "genbank")
        handle.close()
        taxid = None
        for feature in seq_record.features:
            if feature.type == "source":
                for ref in feature.qualifiers.get("db_xref", []):
                    if ref.startswith("taxon:"):
                        taxid = ref.split(":")[1]
                        break
        return taxid, accession_id
    except Exception as e:
        print(f"[{phage_name}] TaxID/Accession fetch error: {e}")
        return None, None

def get_morphotype_from_taxonomy(taxid):
    try:
        handle = Entrez.efetch(db="taxonomy", id=taxid, retmode="xml")
        records = Entrez.read(handle)
        handle.close()
        lineage_ex = records[0].get("LineageEx", [])
        for entry in lineage_ex:
            if entry["Rank"] == "family":
                return entry["ScientificName"]
        return "Unknown morphotype"
    except Exception as e:
        print(f"[{taxid}] Taxonomy error: {e}")
        return "Error"
```

```python
# === Load data ===
df = pd.read_csv(excel_path)
df = df.dropna(subset=[phage_id_col, protein_id_col, seq_col])

# Step 1: Preserve original sequence
df["OriginalSeq"] = df[seq_col]

# Step 2: Clean sequences
df["CleanSeq"] = df[seq_col].apply(clean_sequence)

# Step 3: Create FASTA format
df["FASTA"] = df.apply(make_full_fasta, axis=1)

# Step 4: Add column with just the sequence part used in FASTA
df["SeqInFASTA"] = df["CleanSeq"]

# Step 5: Map duplicates
seq_to_phageids = {}
for _, row in df.iterrows():
    seq = row["CleanSeq"]
    phageid = row[phage_id_col]
    seq_to_phageids.setdefault(seq, []).append(phageid)

def get_duplicate_phageids(row):
    seq = row["CleanSeq"]
    phageid = row[phage_id_col]
    duplicates = [pid for pid in seq_to_phageids.get(seq, []) if pid != phageid]
    return ", ".join(duplicates) if duplicates else ""

df["DuplicatePhageIDs"] = df.apply(get_duplicate_phageids, axis=1)

# Step 6: Remove duplicate sequences, keep first occurrence
df_unique = df.drop_duplicates(subset=["CleanSeq"]).copy()

# Step 7: Retrieve morphotype and accession
taxid_map = {}
accession_map = {}
morphotype_map = {}

for name in df_unique[phage_id_col].astype(str).unique():
    print(f"Fetching info for: {name}")
    taxid, accession = get_taxid_and_accession(name)
    morphotype = get_morphotype_from_taxonomy(taxid) if taxid else "TaxID not found"
    taxid_map[name] = taxid
    accession_map[name] = accession
    morphotype_map[name] = morphotype
    time.sleep(0.4)

df_unique["TaxID"] = df_unique[phage_id_col].map(taxid_map)
df_unique["AccessionFetched"] = df_unique[phage_id_col].map(accession_map)
df_unique["Morphotype"] = df_unique[phage_id_col].map(morphotype_map)

# Step 8: Write FASTA file
with open(output_fasta_file, "w") as fasta_out:
    for fasta_entry in df_unique["FASTA"]:
        fasta_out.write(fasta_entry + "\n")

# Step 9: Save to Excel
df_unique.to_excel(output_excel_file, index=False)

print(f"FASTA saved to '{output_fasta_file}'")
print(f"Excel saved to '{output_excel_file}'")
```

**Workflow data processing – NON-VALIDATED:**

- Fetch the id of bacteriophages from name
- Add column with ids
- Add FASTA format if missing
- Create column with only aa to remove duplicates
- Create column with FASTA format
- Create new column with duplicated identifiers
- Remove duplicate lines based on sequence
- Retrieve morphotypes
- Create FASTA file and Excel file

# 2nd Dataset Processement with Python

**Computationally Predicted RBPs after Processement with Python**

# 2ⁿᵈ Dataset Analysis – K-means and PCA

```python
import pandas as pd
from sklearn.cluster import KMeans
from sklearn.feature_extraction.text import TfidfVectorizer
from sklearn.metrics import silhouette_score
import matplotlib.pyplot as plt
import seaborn as sns
from sklearn.decomposition import PCA
from sklearn.metrics import pairwise_distances
from scipy.cluster.hierarchy import dendrogram, linkage

# === Load data ===
df = pd.read_excel("phage_processed2.xlsx")
sequences = df["RBP sequence (aa)"].astype(str).tolist()

# === Vectorize sequences using TF-IDF of 3-4-mers ===
vectorizer = TfidfVectorizer(analyzer='char', ngram_range=(3, 4))
X = vectorizer.fit_transform(sequences)

# === Best number of clusters by Silhouette score ===
print("Finding optimal number of clusters by silhouette score...")

best_k = None
best_score = -1
scores = []

for k in range(2, 6):
    kmeans = KMeans(n_clusters=k, random_state=42, n_init='auto')
    labels = kmeans.fit_predict(X)
    score = silhouette_score(X, labels)
    scores.append(score)
    print(f"k={k}, silhouette score={score:.4f}")
    if score > best_score:
        best_score = score
        best_k = k

print(f"\nBest k={best_k} with silhouette score={best_score:.4f}")

kmeans = KMeans(n_clusters=best_k, random_state=42, n_init='auto')
df['Cluster'] = kmeans.fit_predict(X)

print(f"Best number of clusters: {best_k} with silhouette score: {best_score:.4f}")

best_model = KMeans(n_clusters=best_k, random_state=42, n_init='auto')
df['Cluster'] = best_model.fit_predict(X)
```

```
Finding optimal number of clusters by silhouette score...
k=2, silhouette score=0.0163
k=3, silhouette score=0.0252
k=4, silhouette score=0.0324
k=5, silhouette score=0.0393

Best k=5 with silhouette score=0.0393
Best number of clusters: 5 with silhouette score: 0.0393
```
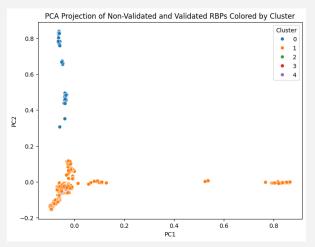
```python
# === PCA ===
pca = PCA(n_components=2)
X_pca = pca.fit_transform(X.toarray())
df['PC1'] = X_pca[:, 0]
df['PC2'] = X_pca[:, 1]

plt.figure(figsize=(8,6))
sns.scatterplot(data=df, x='PC1', y='PC2', hue='Cluster', palette='tab10', legend='full')
plt.title('PCA Projection of Non-Validated RBPs Colored by Cluster')
plt.show()

# === Dendrogram Visualization ===
distance_matrix = pairwise_distances(X.toarray(), metric='cosine')

Z = linkage(distance_matrix, method='ward')

plt.figure(figsize=(30, 10))
labels = df.apply(lambda row: f"Cl{row['Cluster']}_{row['RBP accession nr']}", axis=1).values
dendrogram(Z, labels=labels, truncate_mode='level', p=10, leaf_rotation=90)
plt.xticks(rotation=90, size=8)
plt.title('Dendrogram of Non-Validated RBPs Clustering')
plt.xlabel('Sample Index or Cluster ID')
plt.ylabel('Distance')
plt.show()

# === Save results ===
df.to_excel("rbp_clustered_results_onlyNV.xlsx", index=False)
print("Saved clustering results to rbp_clustered_results_onlyNV.xlsx")
```

```python
best_score = -1
best_params = {}

for k in range(2, 10):
    for init_method in ['k-means++', 'random']:
        model = KMeans(n_clusters=k, init=init_method, random_state=42, n_init='auto')
        labels = model.fit_predict(X)
        score = silhouette_score(X, labels)
        if score > best_score:
            best_score = score
            best_params = {'n_clusters': k, 'init': init_method}

print(f"Best params: {best_params} with silhouette score: {best_score:.4f}")
```

```
Best params: {'n_clusters': 9, 'init': 'random'} with silhouette score: 0.0810
```



PCA Projection of Non-Validated and Validated RBPs Colored by Cluster

**Cluster 0:**
Separated from others → distinct characteristics
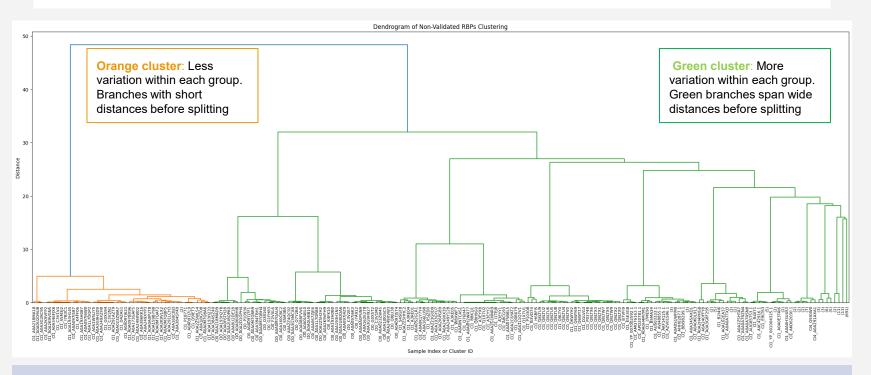**Cluster 1:**
Larger group with some spread along both axes
**Cluster 2, Cluster 3, Cluster 4:**
Less populated and overlap → similarity in characteristics

**Interpretation:**
- Distinct clusters → Non-Validated RBPs may have different characteristics
- **Cluster 0** is notably distinct, potentially representing a unique group of RBPs

# 2nd Dataset Analysis – Dendrogram



Dendrogram of Non-Validated RBPs Clustering

**Orange cluster**: Less variation within each group. Branches with short distances before splitting

**Green cluster**: More variation within each group. Green branches span wide distances before splitting

- Cluster shape: Broad
- Separation between clusters: High
- Cohesion between sub-groups
- Nº Clusters: Low

- Dissimilar sequences that share functional elements

**These clusters likely represent:**
- Possibly novel or divergent RBPs
- Groupings of functionally analogous

# 2nd Dataset – Non-Validated RBPs CD-HIT

```python
from Bio import SeqIO
import pandas as pd
def parse_cdhit_output(file_path):
    clusters = {}
    current_cluster = None

    with open(file_path, 'r') as file:
        for line in file:
            line = line.strip()
            if line.startswith(">Cluster"):
                current_cluster = line.split(" ")[1]
                clusters[current_cluster] = []
            elif current_cluster is not None and line:
                if '>' in line:
                    try:
                        seq_info = line.split(",")[1] if ',' in line else line
                        seq_id_part = seq_info.split(">")[1]
                        seq_id = seq_id_part.split("...")[0].strip()
                        if seq_id:
                            clusters[current_cluster].append(seq_id)
                    except Exception as e:
                        print(f"Error parsing line: {line}\n{e}")
    return clusters


def load_fasta_sequences(fasta_path):
    seq_dict = {}
    for record in SeqIO.parse(fasta_path, "fasta"):
        core_id = record.id
        seq_dict[core_id] = str(record.seq)
    return seq_dict
```

```python
def build_sequences_list(clusters, seq_dict):
    sequences = []
    data = []
    for cluster_id, seq_ids in clusters.items():
        for seq_id in seq_ids:
            seq = seq_dict.get(seq_id, "")
            if seq == "":
                print(f"Warning: Sequence ID '{seq_id}' not found in FASTA file.")
            data.append({'Cluster': cluster_id, 'Sequence ID': seq_id})
            sequences.append(seq)
    df = pd.DataFrame(data)
    return df, sequences

def build_sequences_list(clusters, seq_dict):
    sequences = []
    data = []
    for cluster_id, seq_ids in clusters.items():
        for seq_id in seq_ids:
            seq = seq_dict.get(seq_id, "")
            if seq == "":
                print(f"Warning: Sequence ID '{seq_id}' not found in FASTA file.")
            data.append({'Cluster': cluster_id, 'Sequence ID': seq_id})
            sequences.append(seq)
    df = pd.DataFrame(data)
    return df, sequences

if __name__ == "__main__":
    cluster_file = "cluster_naovalidadas.txt"
    fasta_file = "unique_phage_proteins2.fasta"

    clusters = parse_cdhit_output(cluster_file)
    seq_dict = load_fasta_sequences(fasta_file)
    df, sequences = build_sequences_list(clusters, seq_dict)

    print(f"Total sequences parsed: {len(sequences)}")
    df.to_csv("clustered_sequences2.csv", index=False)
```

```python
from scipy.cluster.hierarchy import dendrogram, linkage
import numpy as np
from sklearn.feature_extraction.text import TfidfVectorizer

# Vectorize the sequences using TF-IDF
vectorizer = TfidfVectorizer(analyzer='char', ngram_range=(3, 4))
X = vectorizer.fit_transform(sequences)  # This is your feature matrix
# Create a linkage matrix using the actual feature matrix
Z = linkage(X.toarray(), method='ward')  # Convert sparse matrix to dense

plt.figure(figsize=(30, 10))
labels = df.apply(lambda row: f"Cl{row['Cluster']}_{row['Sequence ID']}", axis=1).values
dendrogram(Z, labels=labels, truncate_mode='level', p=10, leaf_rotation=90)
plt.xticks(rotation=90, size=5)
plt.title('Dendrogram of Clusters Non Validated')
plt.xlabel('Cluster ID')
plt.ylabel('Distance')
plt.show()
```
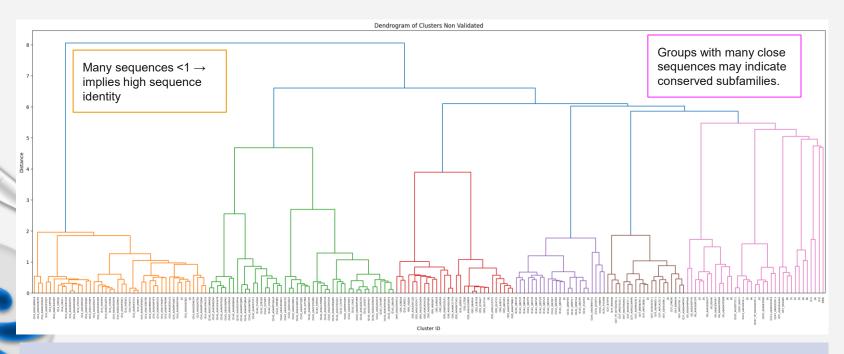
# 2<sup>nd</sup> Dataset – Non-Validated RBPs CD-HIT



Dendrogram of Clusters Non Validated

Many sequences <1 → implies high sequence identity

Groups with many close sequences may indicate conserved subfamilies.
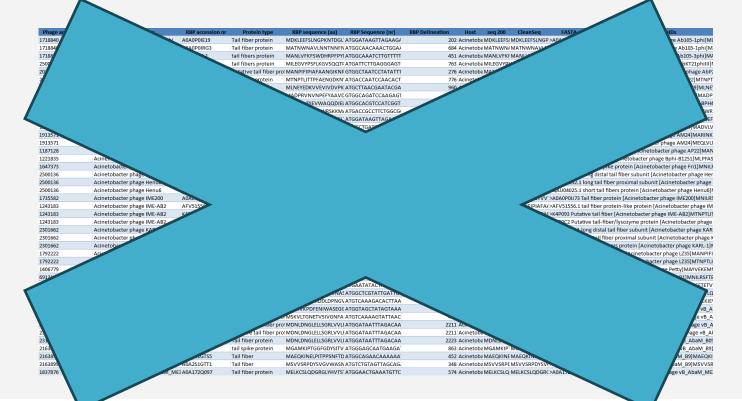
- Cluster shape: Compact, narrow (short branches within)
- Separation between clusters: High (clear distinction)
- Intra-cluster similarity: High (identity ≥ 40%)
- Inter-cluster difference: High (hard cutoff)

Each cluster could correspond to:
- One phage RBP subtype,
- A known structural domain or tail fiber class,
- Evolutionary lineages with strong sequence conservation.

**4.3**

# 3rd Dataset

3ⁿᵈ Dataset – Computationally Predicted and Validated RBPs

```python
import pandas as pd
from Bio import Entrez, SeqIO
import time

# === Configuration ===
Entrez.email = "catiarosario10@gmail.com"
excel_path = "data_test_code.xlsx"
sheet_name = "data_test"
phage_id_col = "Bacteriophage"
protein_id_col = "RBP_Accession"
desc_col = "Protein_Type"
source_col = "Bacteriophage"
seq_col = "RBP_Sequence"
output_fasta_file = "unique_phage_proteins1.fasta"
output_excel_file = "phage_processed1.xlsx"

# === Helper Functions ===

def clean_sequence(seq):
    if not isinstance(seq, str):
        return ""
    lines = seq.strip().splitlines()
    lines = [line for line in lines if not line.startswith(">")]
    return "".join(lines).replace(" ", "").upper()

def make_fasta_header(row):
    desc = row.get(desc_col, "")
    source = row.get(source_col, "")
    desc = desc if pd.notna(desc) else ""
    source = source if pd.notna(source) else ""
    return f">{row[protein_id_col]} {desc} [{source}]".strip()

def make_full_fasta(row):
    header = make_fasta_header(row)
    seq = row["CleanSeq"]
    return f"{header}\n{seq}"
```

```python
def get_taxid_and_accession(phage_name):
    try:
        search = Entrez.esearch(db="nucleotide", term=phage_name, retmode="xml")
        record = Entrez.read(search)
        search.close()
        if not record["IdList"]:
            return None, None
        accession_id = record["IdList"][0]
        handle = Entrez.efetch(db="nucleotide", id=accession_id, rettype="gb", retmode="text")
        seq_record = SeqIO.read(handle, "genbank")
        handle.close()
        taxid = None
        for feature in seq_record.features:
            if feature.type == "source":
                for ref in feature.qualifiers.get("db_xref", []):
                    if ref.startswith("taxon:"):
                        taxid = ref.split(":")[1]
                        break
        return taxid, accession_id
    except Exception as e:
        print(f"[{phage_name}] TaxID/Accession fetch error: {e}")
        return None, None

def get_morphotype_from_taxonomy(taxid):
    try:
        handle = Entrez.efetch(db="taxonomy", id=taxid, retmode="xml")
        records = Entrez.read(handle)
        handle.close()
        lineage_ex = records[0].get("LineageEx", [])
        for entry in lineage_ex:
            if entry["Rank"] == "family":
                return entry["ScientificName"]
        return "Unknown morphotype"
    except Exception as e:
        print(f"[{taxid}] Taxonomy error: {e}")
        return "Error"
```

```python
# === Load data ===
df = pd.read_excel(excel_path, sheet_name=sheet_name)
df = df.dropna(subset=[phage_id_col, protein_id_col, seq_col])

# Step 1: Preserve original sequence
df["OriginalSeq"] = df[seq_col]

# Step 2: Clean sequences
df["CleanSeq"] = df[seq_col].apply(clean_sequence)

# Step 3: Create FASTA format
df["FASTA"] = df.apply(make_full_fasta, axis=1)

# Step 4: Add column with just the sequence part used in FASTA
df["SeqInFASTA"] = df["CleanSeq"]

# Step 5: Map duplicates
seq_to_phageids = {}
for _, row in df.iterrows():
    seq = row["CleanSeq"]
    phageid = row[phage_id_col]
    seq_to_phageids.setdefault(seq, []).append(phageid)

def get_duplicate_phageids(row):
    seq = row["CleanSeq"]
    phageid = row[phage_id_col]
    duplicates = [pid for pid in seq_to_phageids.get(seq, []) if pid != phageid]
    return ", ".join(duplicates) if duplicates else ""

df["DuplicatePhageIDs"] = df.apply(get_duplicate_phageids, axis=1)

# Step 6: Remove duplicate sequences, keep first occurrence
df_unique = df.drop_duplicates(subset=["CleanSeq"]).copy()

# Step 7: Retrieve morphotype and accession
taxid_map = {}
accession_map = {}
morphotype_map = {}

for name in df_unique[phage_id_col].astype(str).unique():
    print(f"Fetching info for: {name}")
    taxid, accession = get_taxid_and_accession(name)
    morphotype = get_morphotype_from_taxonomy(taxid) if taxid else "TaxID not found"
    taxid_map[name] = taxid
    accession_map[name] = accession
    morphotype_map[name] = morphotype
    time.sleep(0.4)

df_unique["TaxID"] = df_unique[phage_id_col].map(taxid_map)
df_unique["AccessionFetched"] = df_unique[phage_id_col].map(accession_map)
df_unique["Morphotype"] = df_unique[phage_id_col].map(morphotype_map)

# Step 8: Write FASTA file
with open(output_fasta_file, "w") as fasta_out:
    for fasta_entry in df_unique["FASTA"]:
        fasta_out.write(fasta_entry + "\n")

# Step 9: Save to Excel
df_unique.to_excel(output_excel_file, index=False)

print(f"FASTA saved to '{output_fasta_file}'")
print(f"Excel saved to '{output_excel_file}'")
```

**Workflow data processing – NON-VALIDATED:**

- Fetch the id of bacteriophages from name
- Add column with ids
- Add FASTA format if missing
- Create column with only aa to remove duplicates
- Create column with FASTA format
- Create new column with duplicated identifiers
- Remove duplicate lines based on sequence
- Retrieve morphotypes
- Create FASTA file and Excel file
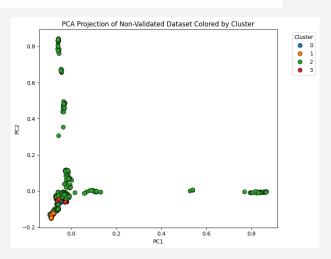
# 3rd Dataset Processement with Python

**Computationally Predicted and Validated RBPs after Processement with Python**

# 3rd Dataset Analysis – K-means and PCA

```python
import pandas as pd
from sklearn.cluster import KMeans
from sklearn.feature_extraction.text import TfidfVectorizer
from sklearn.metrics import silhouette_score
from sklearn.decomposition import PCA
from sklearn.metrics import pairwise_distances
from scipy.cluster.hierarchy import dendrogram, linkage
import matplotlib.pyplot as plt
import seaborn as sns

# === Load dataset ===
df = pd.read_excel("phage_processed1.xlsx")

# === Vectorization ===
sequences = df["CleanSeq"].astype(str).tolist()
vectorizer = TfidfVectorizer(analyzer='char', ngram_range=(3, 4))
X = vectorizer.fit_transform(sequences)

# === Find optimal k using silhouette score ===
print("Finding best k based on silhouette score...")
best_k = None
best_score = -1
for k in range(2, 5):
    kmeans = KMeans(n_clusters=k, random_state=42, n_init='auto')
    labels = kmeans.fit_predict(X)
    score = silhouette_score(X, labels)
    print(f"k={k}, silhouette score={score:.4f}")
    if score > best_score:
        best_score = score
        best_k = k
print(f"Best k = {best_k} with silhouette score = {best_score:.4f}")

# === Final clustering with best k ===
kmeans = KMeans(n_clusters=best_k, random_state=42, n_init='auto')
df['Cluster'] = kmeans.fit_predict(X)
```

```
Finding best k based on silhouette score...
k=2, silhouette score=0.0140
k=3, silhouette score=0.0233
k=4, silhouette score=0.0301
Best k = 4 with silhouette score = 0.0301
```

```
Best params: {'n_clusters': 9, 'init': 'random'} with silhouette score: 0.0774
```

```python
# === PCA Visualization ===
pca = PCA(n_components=2)
X_pca = pca.fit_transform(X.toarray())
df["PC1"] = X_pca[:, 0]
df["PC2"] = X_pca[:, 1]

plt.figure(figsize=(8, 6))
sns.scatterplot(data=df, x="PC1", y="PC2", hue="Cluster", palette="tab10", s=60, edgecolor='k')
plt.title("PCA Projection of Non-Validated Dataset Colored by Cluster")
plt.legend(title="Cluster", bbox_to_anchor=(1.05, 1), loc='upper left')
plt.tight_layout()
plt.show()

# === Dendrogram Visualization ===
distance_matrix = pairwise_distances(X.toarray(), metric='cosine')

Z = linkage(distance_matrix, method='ward')

plt.figure(figsize=(30, 10))
labels = df.apply(lambda row: f"C{row['Cluster']}_{row['RBP_Accession']}", axis=1).values
dendrogram(Z, labels=labels, truncate_mode='level', p=10, leaf_rotation=90)
plt.xticks(rotation=90, size=8)
plt.title('Dendrogram of Both RBPs Clustering')
plt.xlabel('Sample Index or Cluster ID')
plt.ylabel('Distance')
plt.show()

# === Save results ===
df.to_excel("rbp_clustered_results_unsupervised.xlsx", index=False)
print("\nSaved results to 'rbp_clustered_results_unsupervised.xlsx'")
```

```python
best_score = -1
best_params = {}

for k in range(2, 10):
    for init_method in ['k-means++', 'random']:
        model = KMeans(n_clusters=k, init=init_method, random_state=42, n_init='auto')
        labels = model.fit_predict(X)
        score = silhouette_score(X, labels)
        if score > best_score:
            best_score = score
            best_params = {'n_clusters': k, 'init': init_method}

print(f"Best params: {best_params} with silhouette score: {best_score:.4f}")
```



PCA Projection of Non-Validated Dataset Colored by Cluster

**Cluster 0:**
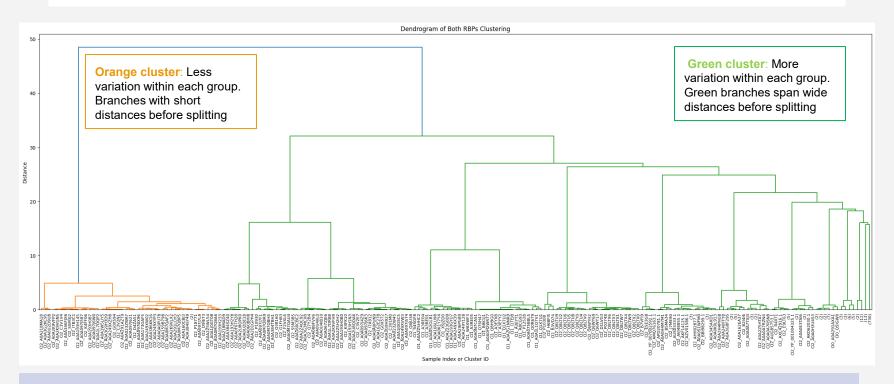Numerous points but overlaps with **Cluster 2**
**Cluster 2:**
Has most points → dominant cluster in this dataset
**Cluster 1, Cluster 3:**
Less populated and overlap → similarity in characteristics

**Interpretation:**
- Distinct clusters → Validated and Non-Validated RBPs may have different characteristics

# 3<sup>rd</sup> Dataset Analysis – Dendrogram



Dendrogram of Both RBPs Clustering

**Orange cluster**: Less variation within each group. Branches with short distances before splitting

**Green cluster**: More variation within each group. Green branches span wide distances before splitting

- Cluster shape: Broad, with higher internal variation
- Separation between clusters: High
- Cohesion between sub-groups
- Nº Clusters: Low

- Dissimilar sequences that share functional elements
**These clusters likely represent:**
- Possibly novel or divergent RBPs
- Groupings of functionally analogous

# 3rd Dataset – Both RBPs CD-HIT

```python
from Bio import SeqIO
import pandas as pd
def parse_cdhit_output(file_path):
    clusters = {}
    current_cluster = None

    with open(file_path, 'r') as file:
        for line in file:
            line = line.strip()
            if line.startswith(">Cluster"):
                current_cluster = line.split(" ")[1]
                clusters[current_cluster] = []
            elif current_cluster is not None and line:
                if '>' in line:
                    try:
                        seq_info = line.split(",")[1] if ',' in line else line
                        seq_id_part = seq_info.split(">")[1]
                        seq_id = seq_id_part.split("...")[0].strip()
                        if seq_id:
                            clusters[current_cluster].append(seq_id)
                    except Exception as e:
                        print(f"Error parsing line: {line}\n{e}")
    return clusters


def load_fasta_sequences(fasta_path):
    seq_dict = {}
    for record in SeqIO.parse(fasta_path, "fasta"):
        core_id = record.id
        seq_dict[core_id] = str(record.seq)
    return seq_dict
```
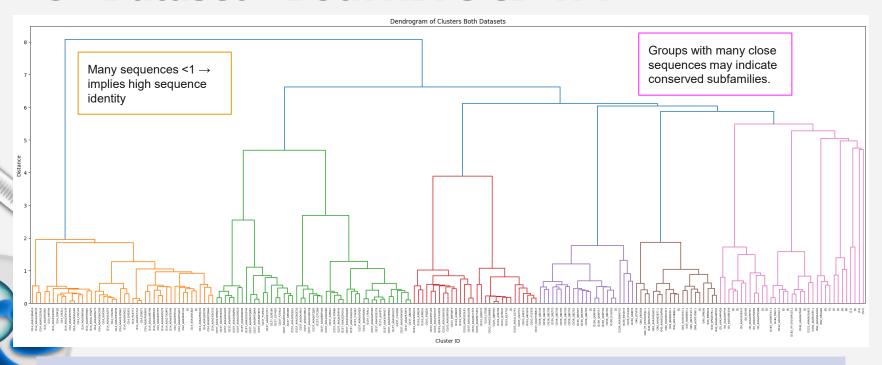
```python
def build_sequences_list(clusters, seq_dict):
    sequences = []
    data = []
    for cluster_id, seq_ids in clusters.items():
        for seq_id in seq_ids:
            seq = seq_dict.get(seq_id, "")
            if seq == "":
                print(f"Warning: Sequence ID '{seq_id}' not found in FASTA file.")
            data.append({'Cluster': cluster_id, 'Sequence ID': seq_id})
            sequences.append(seq)
    df = pd.DataFrame(data)
    return df, sequences


if __name__ == "__main__":
    cluster_file = "cluster_both.txt"
    fasta_file = "unique_phage_proteins1.fasta"

    clusters = parse_cdhit_output(cluster_file)
    seq_dict = load_fasta_sequences(fasta_file)
    df, sequences = build_sequences_list(clusters, seq_dict)

    print(f"Total sequences parsed: {len(sequences)}")
    df.to_csv("clustered_sequences1.csv", index=False)
```

```python
from scipy.cluster.hierarchy import dendrogram, linkage
import numpy as np
from sklearn.feature_extraction.text import TfidfVectorizer

vectorizer = TfidfVectorizer(analyzer='char', ngram_range=(3, 4))
X = vectorizer.fit_transform(sequences)
Z = linkage(X.toarray(), method='ward')

plt.figure(figsize=(30, 10))
labels = df.apply(lambda row: f"Cl{row['Cluster']}_{row['Sequence ID']}", axis=1).values
dendrogram(Z, labels=labels, truncate_mode='level', p=10, leaf_rotation=90)
plt.xticks(rotation=90, size=5)
plt.title('Dendrogram of Clusters Both Datasets')
plt.xlabel('Cluster ID')
plt.ylabel('Distance')
plt.show()
```

# 3rd Dataset – Both RBPs CD-HIT



Dendrogram of Clusters Both Datasets

Many sequences <1 → implies high sequence identity
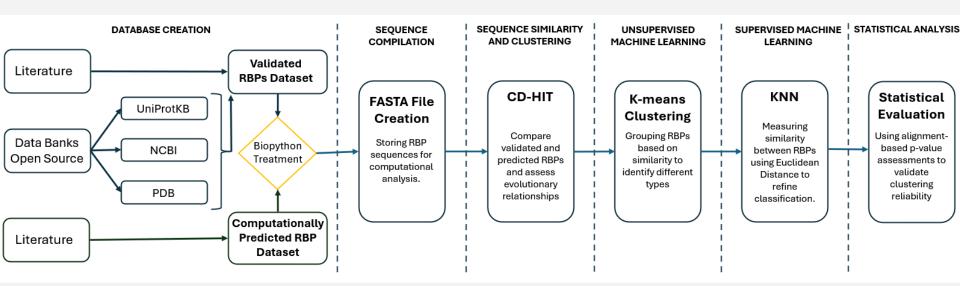
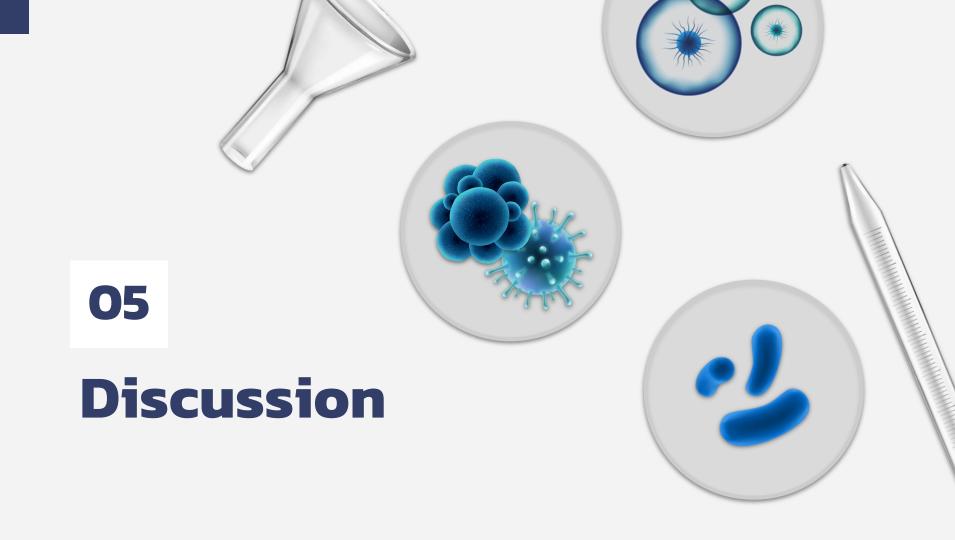Groups with many close sequences may indicate conserved subfamilies.

- Cluster shape: Compact, narrow (short branches within)
- Separation between clusters: High (clear distinction)
- Intra-cluster similarity: High (identity ≥ 40%)
- Inter-cluster difference: High (hard cutoff)

Each cluster could correspond to:
- One phage RBP subtype,
- A known structural domain or tail fiber class,
- Evolutionary lineages with strong sequence conservation.

PROJECT WORKFLOW

**05**

# Discussion

# Discussion

## Challenges

**Creation of the provisory datasets:**
- Filter and selection of RBPs from the 3 Databases
- Automatization of morphotype retrieval
- Datasets Uniformization

**Validation of RBPs:**
- Literature key word identification
- Explicit experimental data

## Success

**Curated RBP Database:**
- Two datasets with Validated RBP and RBP candidates

**Implemented dual clustering approaches**:

- CD-HIT (identity-based): Captured conserved structural/evolutionary groups

- K-Means (motif-based): Revealed functional similarity based on k-mer features

**Visualized clustering via dendrograms:**
- Assess sequence relationships
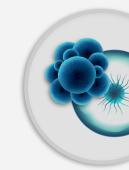
# Conclusion

## 🔍 Next Steps

- **Propythia -** platform for the classification of biological sequences (proteins and DNA) using machine and deep learning

- Train classifiers using validated clusters as labels
- Deploy the RBP-ID web Based tool
- Enable user input and visualization, returning predicted RBP status, cluster info, and motifs
- Expand the validated dataset via literature mining

## 🌐 Vision

RBP-ID aims to become a go-to platform for phage researchers, accelerating discovery of receptor-binding proteins for synthetic biology, diagnostics, and phage therapy

# Thanks