

# Workshop III – Bioinformatics and Viral Genomics

***Untargeted Viromics Session***

***Durban, South Africa***

***February 27, 2025***

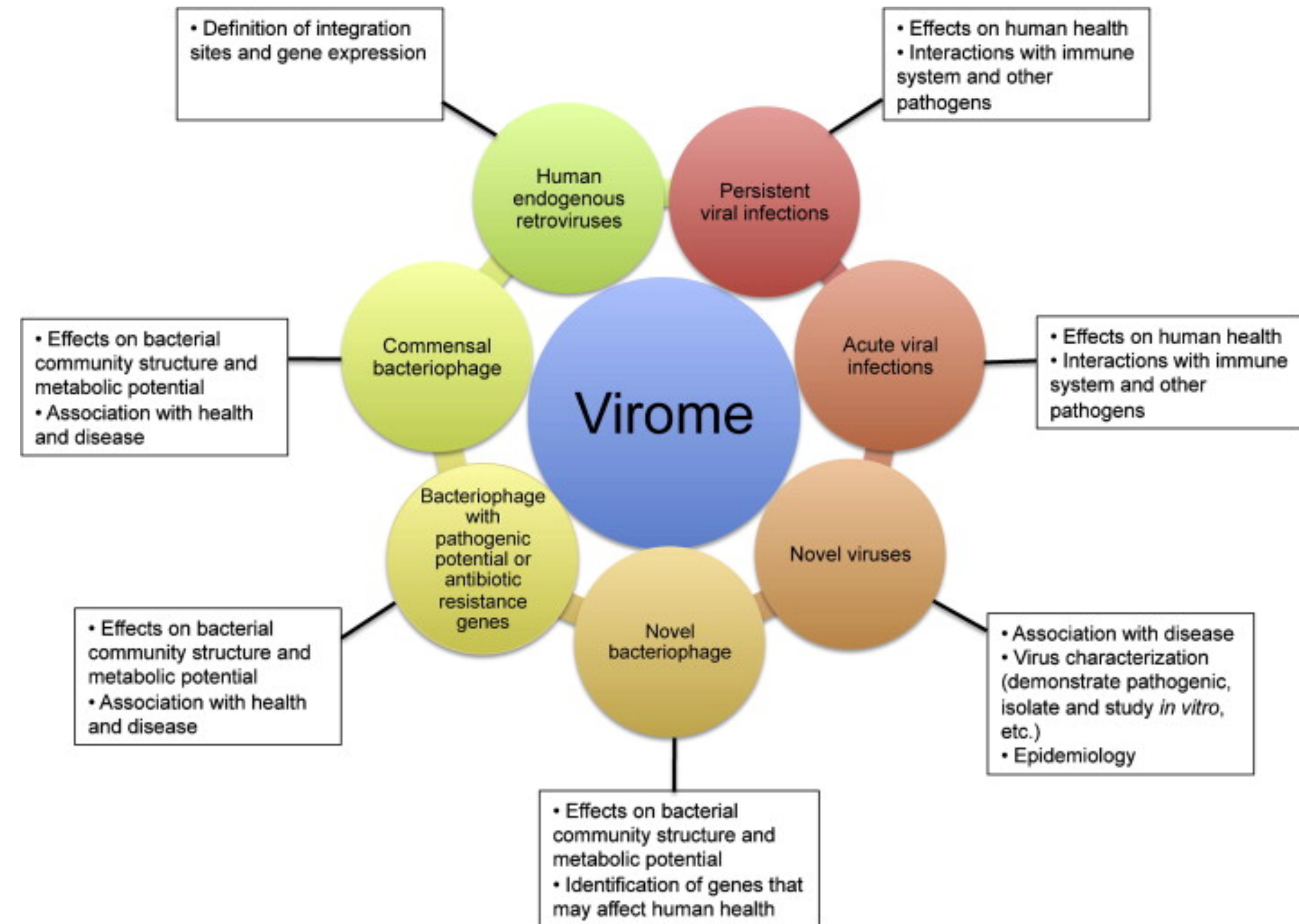
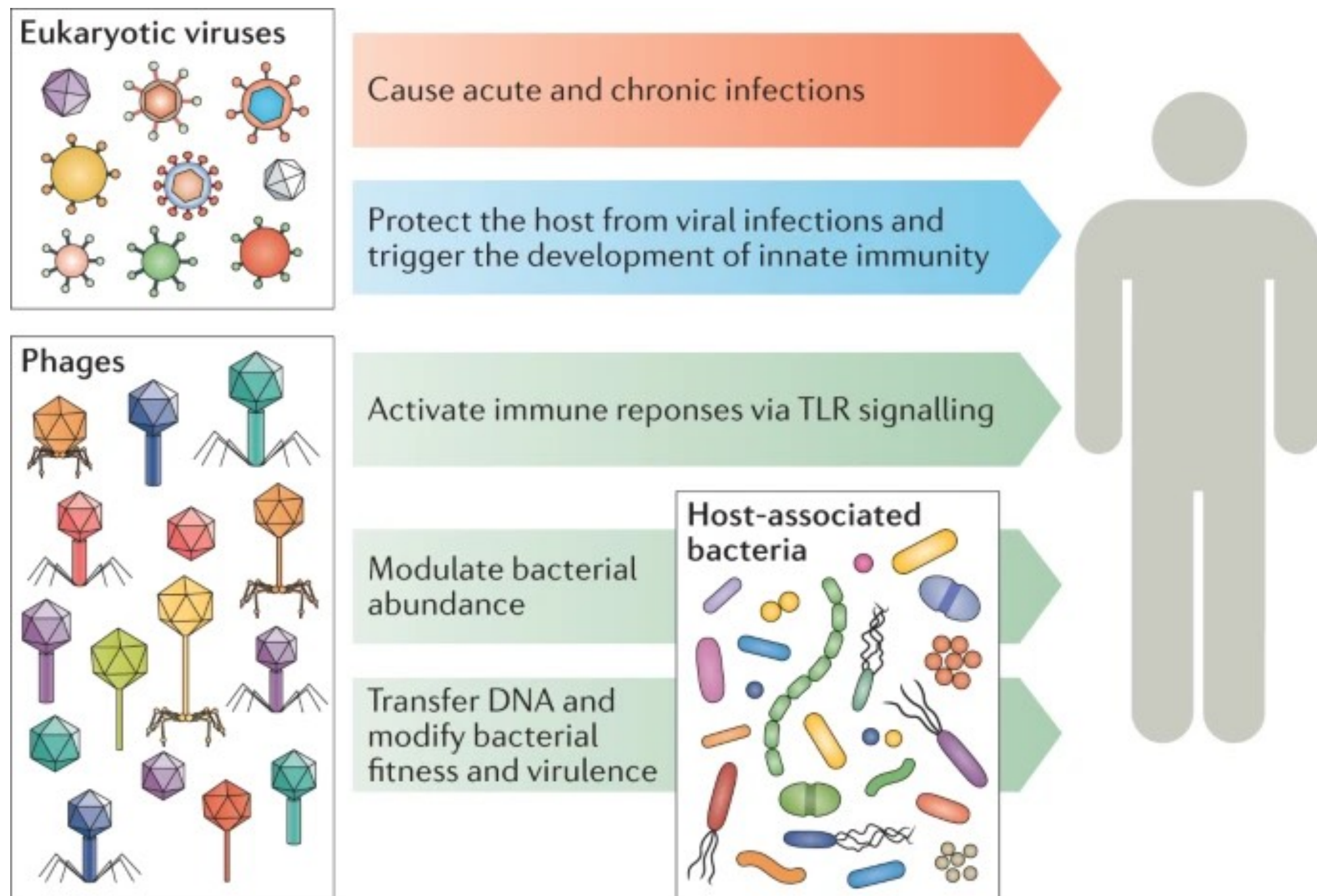
***Luis Chica***

# Outline

- Introductory Presentation
  - Virome Background
  - Uncultivated Viral Genomes (UViGs)
  - Key methodologies for identifying and analyzing UViGs
  - Bulk Metagenomics vs. Virus-Like Particle (VLP) Enrichment (Review)
  - Main Approaches for Viral Prediction
  - Pipelines for Analyzing the Virome
- Hands-On Session
- Hands-On Review

# Definitions and role of the virome

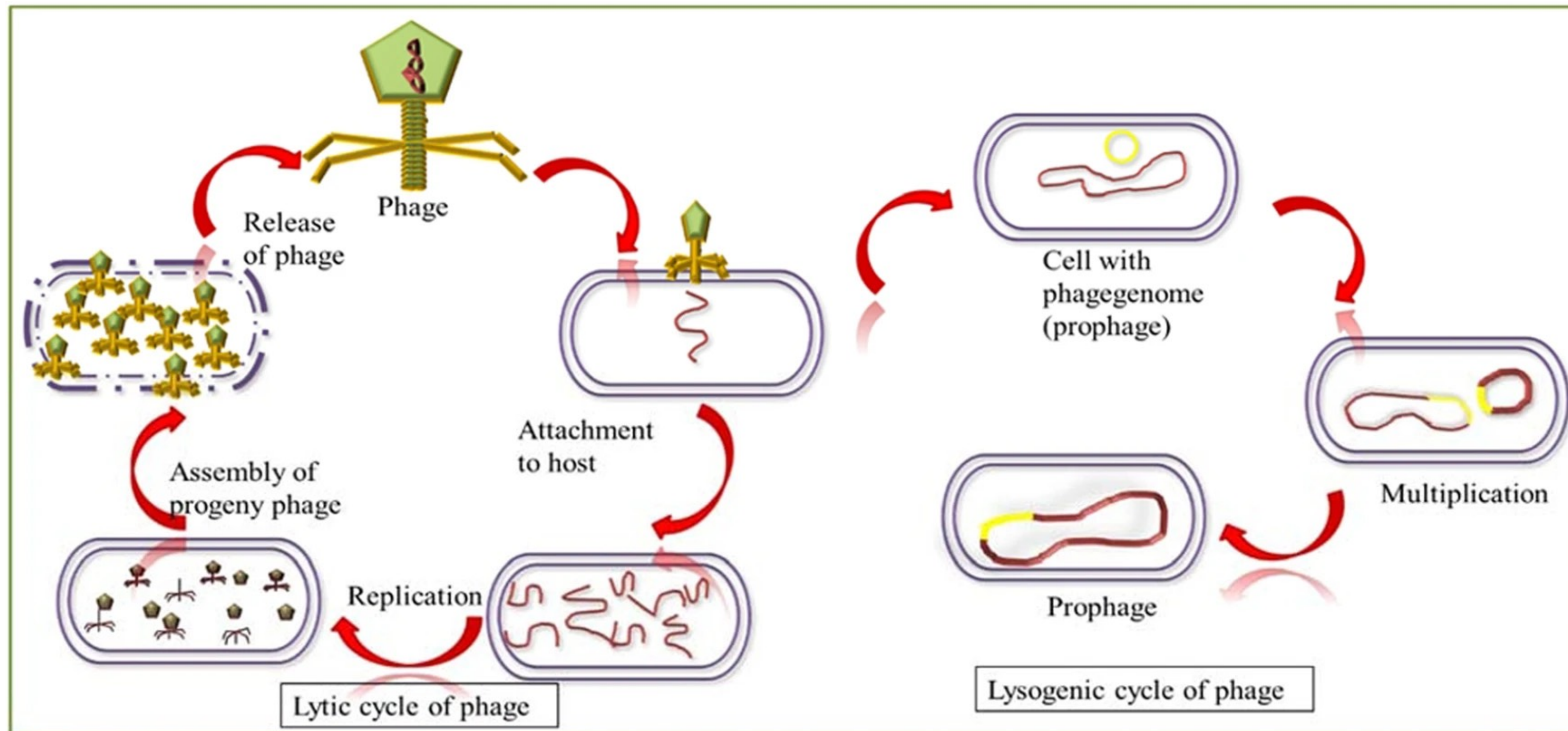
*Spanning eukaryotic viruses to bacteriophages*





# Definitions and role of the virome

## *Concept of bacteriophage lifecycles*

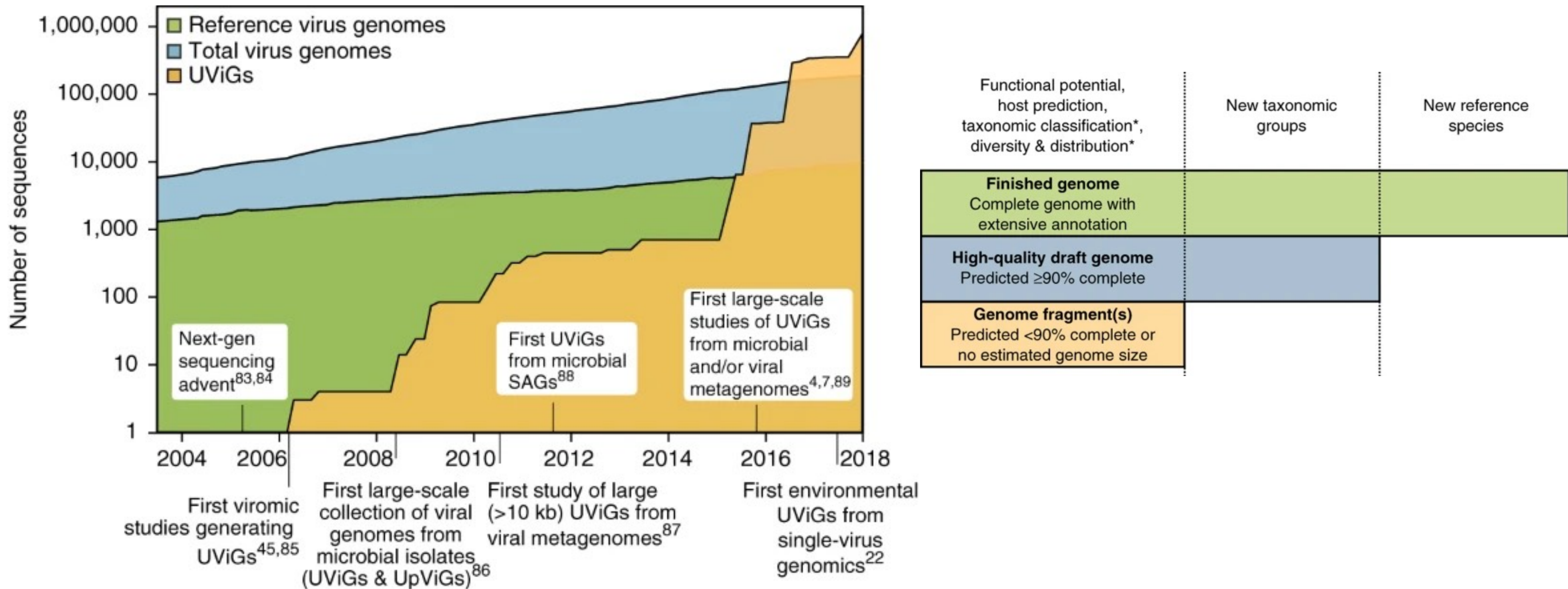


- Lysogenic phages can integrate to the bacterial genome and replicate as long as the bacteria replicates.
- Lysogenic phages = temperate phages.
- Prophages: Stage in which the phage is integrated in the genome.



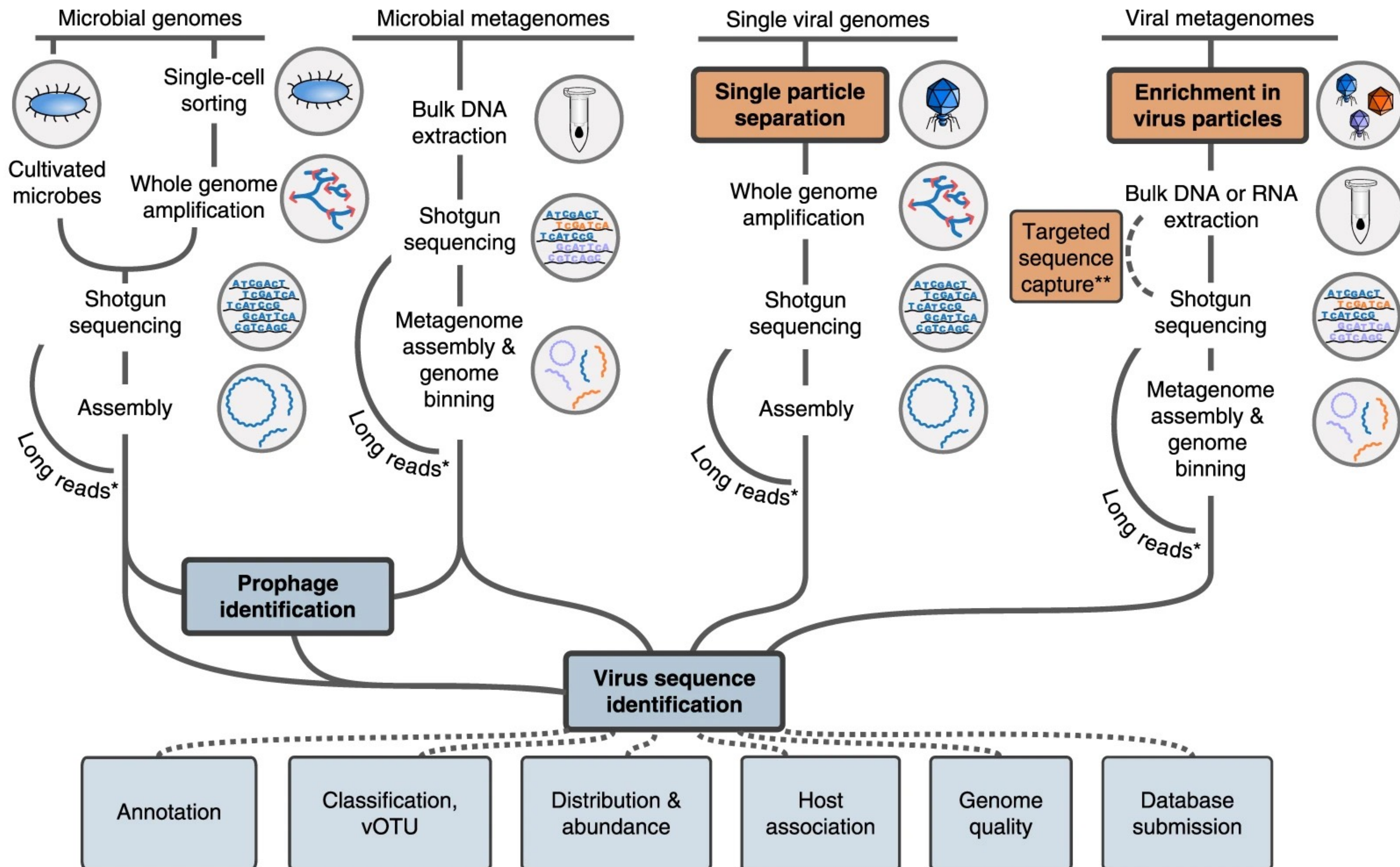
# Uncultivated viruses (uViGs)

*Drastic increase of the number of uViGs deposited in databases over the past years.*



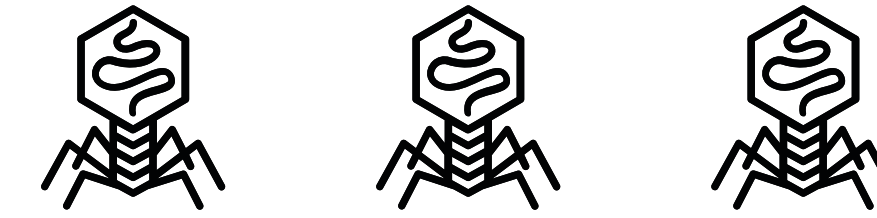
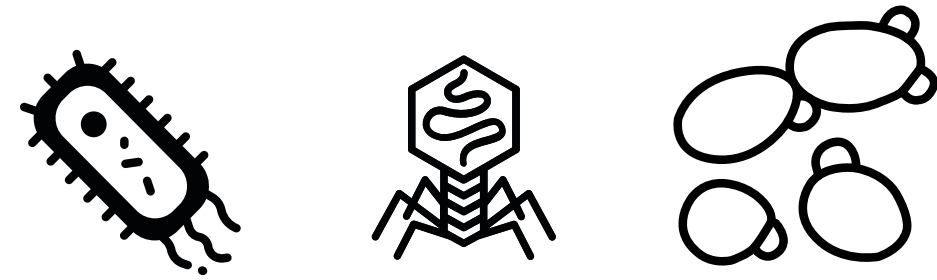



# Main techniques for analysing uncultivated viruses




# Bulk metagenomics vs VLP enrichment

*Depending on the approach your viral results change*



- 
- Comprehensive Sampling: Allow us to relate the different microbial communities within the environment.
  - More accurate identification of prophages and their host.
  - Lower Viral Specificity: hard to detect low-abundance viruses.

- Enhanced Viral Detection: More sensitive method for detecting viruses.
- Reduced Background Noise: Eliminates most non-viral genetic material.
- Possibility of detecting RNA viruses

- 
- Complex Data Analysis
  - Higher Background Noise: potentially mask of viral signals.

- Misses prophages and latent viruses
- Complex and costly sample preparation.



# Main approaches for viral prediction

## COMPUTATIONAL TOOLS FOR PREDICTING VIRAL SEQUENCES

Homology

Comparing sequences using viral databases and local alignment

**Pros:** High accuracy for known viruses. Allows distant homologous detection  
**Cons:** Dependent on the quality and completeness of reference databases. Slow

K-mers

Dividing sequences into subsequences and comparing against DBs

**Pros:** Fast and scalable.  
**Cons:** Detection is limited to high identity relatives within databases

Machine learning

Training model with viral genomic features

**Pros:** Can detect novel viruses. High accuracy with well-trained models.  
**Cons:** computational intensive



# Main approaches for viral prediction

## COMPUTATIONAL TOOLS FOR PREDICTING VIRAL SEQUENCES

Phylogenetic approaches

Using phylogenetic defined references to located the query sequences

**Pros:** Provides evolutionary context. Useful for novel virus discovery.  
**Cons:** Computationally intensive. Requires high-quality alignments.

Hybrid

Normally, it combines machine learning and homology approaches

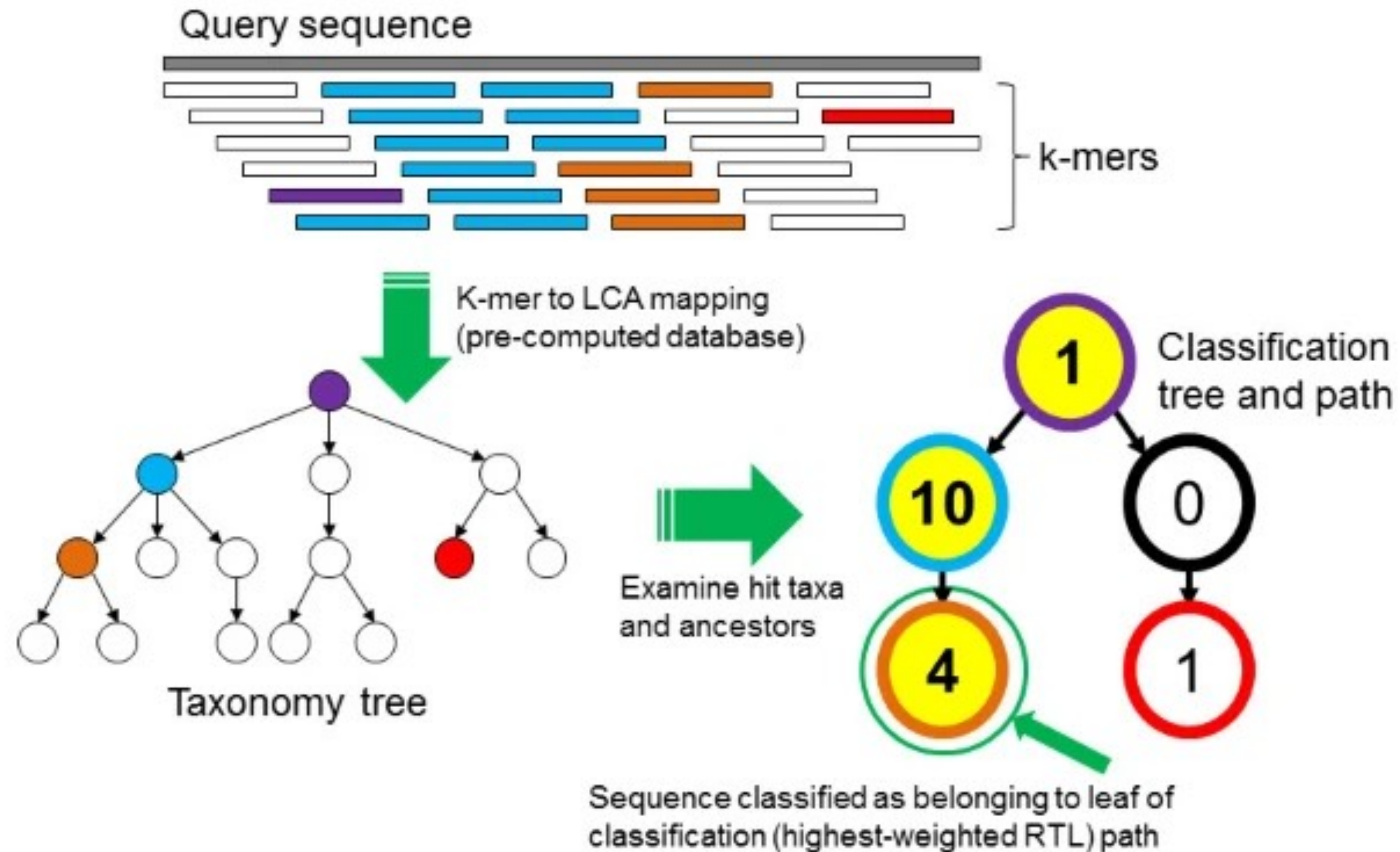
- k-mer analysis typically relies on reads as input.
- Homology-based approaches can use either reads or contigs as input.
- Machine learning and hybrid approaches usually require contigs as input, along with several genomic features for accurate prediction.
- Phylogenetic analysis can use contigs or specific genes extracted after assembly.





# Main approaches for viral prediction

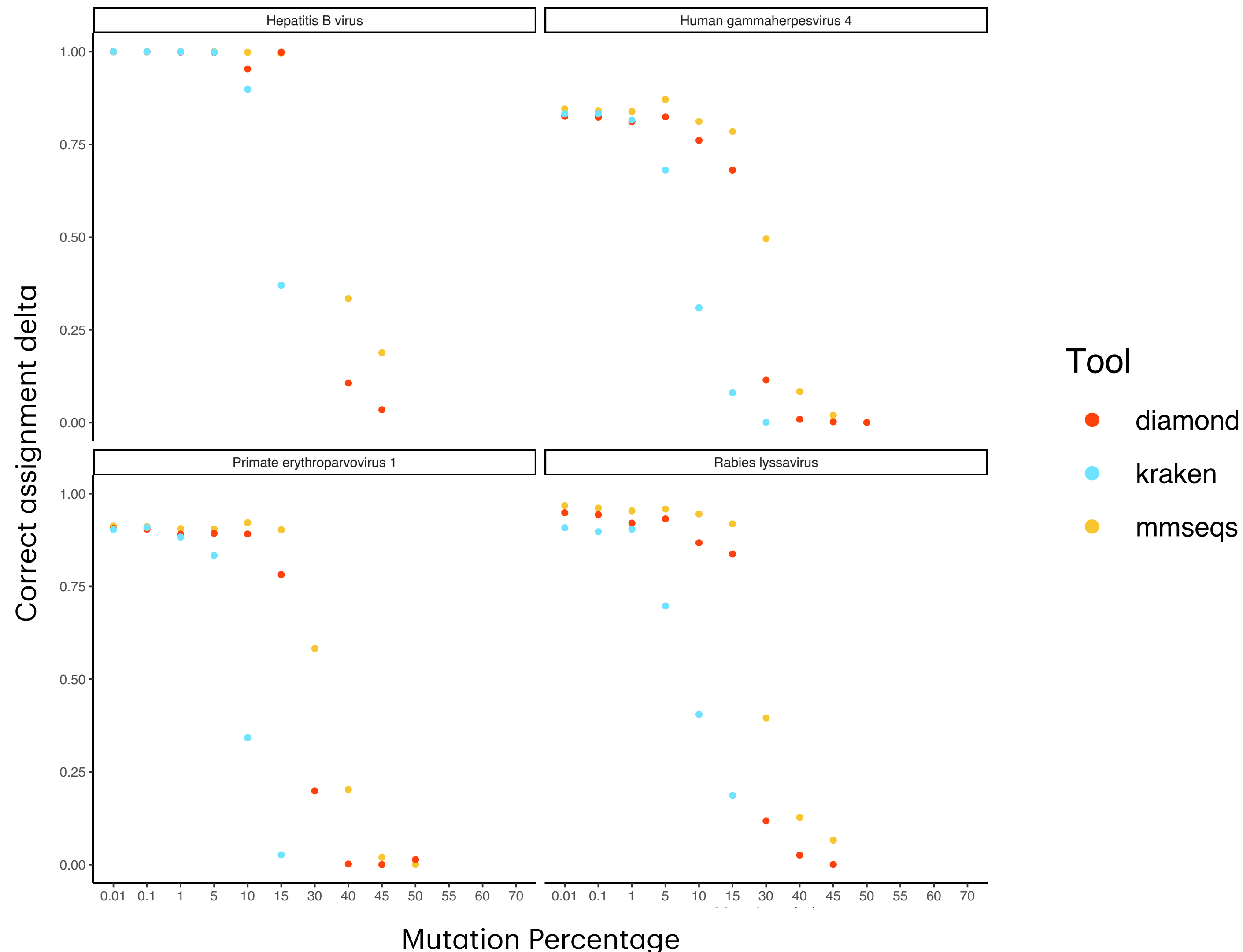
## *K-mer based approaches*



- Use of k-mer frequency analysis to identify unique compositional patterns in viral sequences.
- Comparison of k-mer profiles against curated databases or reference genomes for viral classification.

# When to use homology or k-mers for read based analysis

***How well-known is the system you are working on and data set size as key factors***



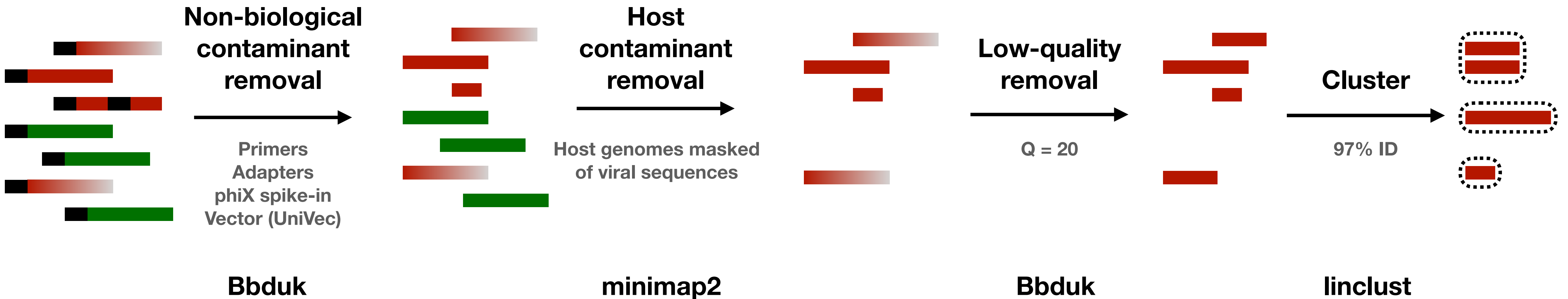
- An increase in mutations can significantly affect the performance of k-mer-based tools.
- This is particularly problematic for novel viruses or viruses without close representatives in the reference database, as they may be missed or misclassified.



***Complete workflows for viral analysis: Read  
and contig-based approaches***

# Quality control for viral genomics

*Key step for avoiding false positives and reducing data sizes*



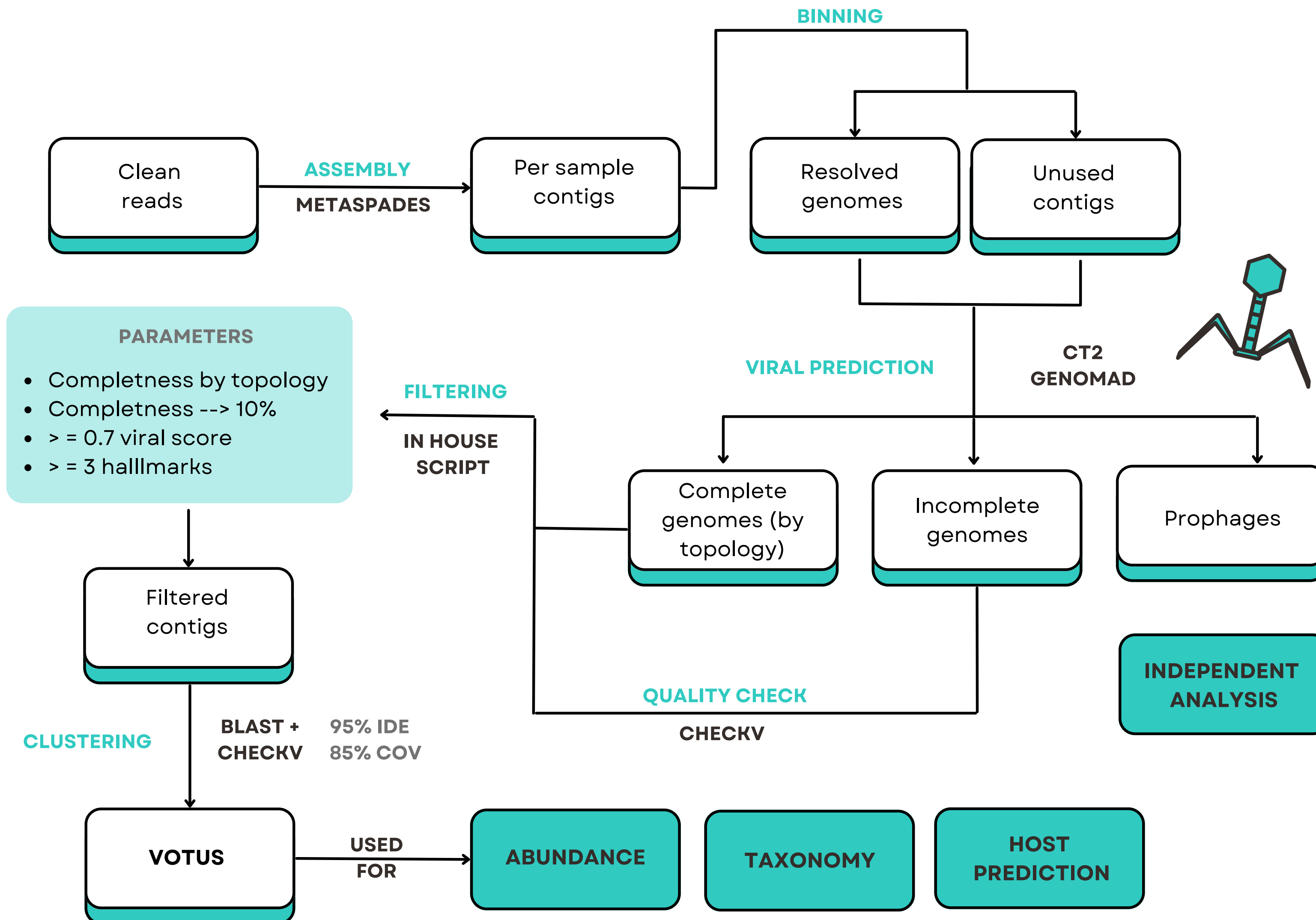
Virus masked reference genomes:

- Human
- Mouse
- Rat
- Macaque
- Pig
- Bat
- Cow
- Camel
- Cat
- Dog
- Mosquito
- Tick
- C. elegans



# Workflow for viral analysis using contigs

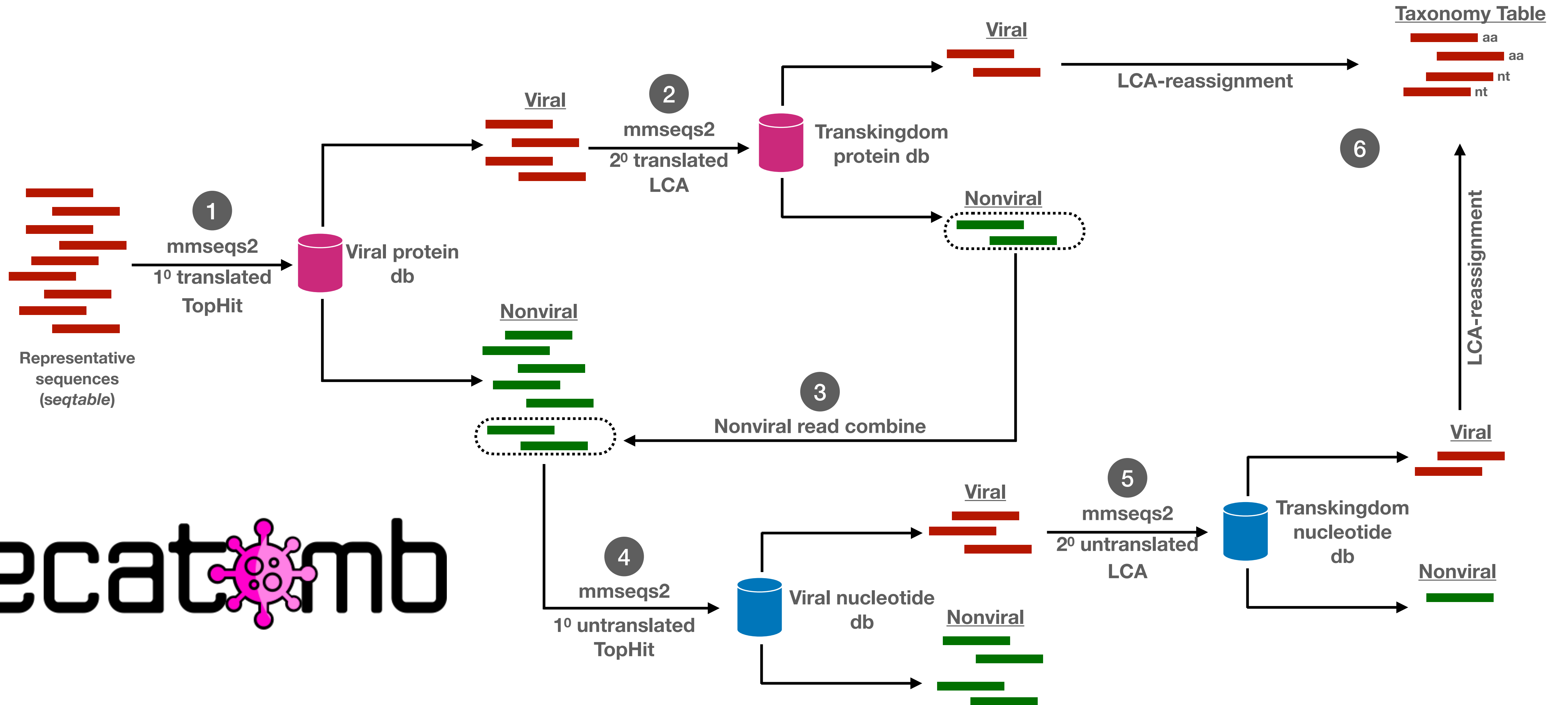
*No gold standard available*



• Optional steps:

- Binning
- Clustering

# Workflow for viral analysis using reads







# Summary

- **Sequencing Methodology:**

- The choice of sequencing approach (e.g., bulk metagenomics, VLP enrichment) depends on the main scope of the project.
- Different methodologies impact the detection of specific viruses and the estimation of viral abundance.

- **Viral Detection:**

- Viral detection should be performed carefully, applying multiple tools to increase confidence in predictions.
- Combining complementary approaches enhances robustness.

- **Tool Selection:**

- For read-based analysis, tool selection should consider the novelty of the system.
- k-mer-based approaches may have limitations in understudied systems due to the lack of reference data.
- Homology-based approaches are time- and resource-intensive but are particularly useful for unknown or novel systems.



# Acknowledgements



## **Handley Lab**

Scott Handley

Leran Wang

Kathie Mihindukulasuriya

Lindsay Droit

Dhoha Abid

Megan Johnson

Martina Moore

Dave Wang

## **Center for Genome Sciences**

Jessica Hoisington-Lopez

MariaLynn Crosby



## **Kwon Lab**

Doug Kwon

Joseph Elsherbini

Sarah Eisa

Cameron Reitan