

# Research Paper Review

Huiming Zhou 1327380

February 25, 2025

## 1. Introduction

### Article Information

**Title:** VirHunter: A Deep Learning-Based Method for Detection of Novel RNA Viruses in Plant Sequencing Data

**Authors:** Grigori Sukhorukov, Maryam Khalili, Olivier Gascuel, Thierry Candresse, Armelle Marais-Colombel and Macha Nikolski

**Publication Year:** 2022

**Journal:** Frontier

**DOI/Link:** [10.3389/fbinf.2022.867111](https://doi.org/10.3389/fbinf.2022.867111)

### Reason for Selection

As a student in the MDS program, I have acquired a solid foundation of theoretical knowledge during my undergraduate and general graduate studies. At the same time, I have always been deeply curious about the field of biology and its interdisciplinary applications. This article, which explores the application of deep learning in the biological domain, perfectly aligns with my interests and expectations. Not only does it present the authors' own model, but it also provides a comprehensive overview of the principles behind other models. Through this article, I can gain insights into the modern detection background of the virome in plant RNA, as well as recent research achievements in the field. This makes it an invaluable resource for expanding my knowledge and understanding of cutting-edge developments in bioinformatics.

## 2. Summary of the Article

### Research Problem/Question

With the adoption of High-Throughput Sequencing (HTS), a DNA sequencing technology, the sequencing of plant viromes has become increasingly common. These datasets are rapidly advancing our understanding of viral diversity. Based on these developments, Stobbe and Rossinck classify viruses in HTS datasets into three groups: (1) viruses previously unknown to infect a given host, (2) viruses with limited sequence similarity to known viruses, and (3) viruses sharing little to no sequence similarity with any known viruses in existing databases. This highlights the need for efficient virus detection methods, especially given the complexity of plant virus analysis, including high mutation rates, co-infection by multiple unrelated viral species, and unavoidable background contamination.

Additionally, the dominance of host material in plant sequencing samples often results in very low abundance of viral genomic material. Current methods for discovering novel viruses from assembled contigs including VirFinder, VirSorter2 and DeepVirFinder, with the first two being computationally intensive and requiring significant expertise for filtering and interpreting results. The third category, which relies on machine learning, offers a more promising approach. However, there is still a need for improved methods to address these challenges effectively.

### Methodology

The proposed VirHunter model combines a multi-path neural network CNN module with a downstream random forest classifier module. The first part of the model employs a k-mer-based approach, where k-mer sizes of 5, 7, and 10 are used as input. Three separate CNN models are utilized, each sharing the same architecture but differing in the number of filters in the convolutional and dense layers. The second dense layer consists of three units and uses a softmax activation function to classify sequences into three categories: virus, plant, or bacteria. The

second part of the model, the random forest classifier, takes the nine outputs from the CNN module (three outputs per k-mer size) and predicts one of the three classes. The model is trained on datasets downloaded from NCBI, including viral sequences with host taxonomic IDs belonging to Viridiplantae (plants), as well as plant sequences from species such as peach, grapevine, sugar beet, and rice.

## **Findings/Results/Impact**

The evaluation of VirHunter demonstrated its superior performance compared to state-of-the-art tools such as DeepVirFinder, VirSorter2, and tBLASTx in terms of True Positives (TPs) when tested on randomly sampled datasets. Across all leave-out experiments, VirHunter consistently achieved a high TP rate in correctly classifying both bacterial and plant fragments. Notably, the model exhibited high accuracy in classifying plant fragments when trained on phylogenetically close plant species.

Furthermore, VirHunter maintained a high TP rate for viral and bacterial fragments even when the hosts in the training and testing datasets were phylogenetically distant. Additionally, the model showed robustness to increasing mutation rates, with the TP rate declining only gradually as the mutation rate increased. These results highlight VirHunter's potential as a reliable and versatile tool for virome analysis, particularly in complex scenarios involving diverse hosts and high mutation rates.

## **3. Critical Analysis**

### **Strengths**

This article demonstrates significant scientific rigor and innovation by building upon the structure of previous models and enhancing them through the integration of various deep learning methods. The proposed model, VirHunter, was rigorously tested on multiple datasets, showcasing its superior performance in detecting novel viruses in plant RNAs. Compared to traditional methods, VirHunter is more user-friendly and time-efficient, while also delivering improved accuracy in identifying novel viruses. This makes it a highly suitable tool for analyzing HTS-acquired viromes.

The clarity of the article is commendable, as it provides a detailed explanation of the methodology, including the combination of multi-path CNN modules and random forest classifiers. The results have a strong potential impact on bioinformatics and related fields, particularly in the study of plant viromes, where efficient and accurate virus detection is critical. By addressing the limitations of existing tools, VirHunter sets a new standard for virome analysis and opens avenues for further research in viral diversity and host-pathogen interactions.

### **Weaknesses**

One limitation of the study is that the proposed method, VirHunter, was trained and tested on datasets from only four specific plant species (peach, grapevine, sugar beet, and rice). While the results are promising, the model's generalizability to other plant species, including non-agricultural plants such as wildflowers, corals, or cacti, remains unverified. Expanding the evaluation to a more diverse range of species would strengthen the validity of the conclusions and demonstrate the model's broader applicability.

Additionally, the study does not thoroughly investigate cases where BLAST-based methods or other state-of-the-art tools outperform VirHunter. A deeper analysis of these scenarios could provide valuable insights into the limitations of the proposed method and guide future improvements. While the conclusions are generally well-supported by the evidence provided, addressing these weaknesses would further enhance the robustness and impact of the research.

### **Clarity of Communication**

The article is well-organized and presents its findings in a logical manner. However, it uses specialized terminology such as HTS (High-Throughput Sequencing), k-mer, CNN (Convolutional Neural Networks), and random forest, which may be familiar to experts in bioinformatics or machine learning but could pose a challenge for readers without a background in these fields. The figures and visual aids in the paper are of high quality, effectively illustrating the experimental results and providing clear, concise representations of the data. Overall, the article is best suited for machine learning experts with some background in biology or bioinformatics professionals.

## **4. Reflection and Suggestions**

### **Reflection**

From this article, I learned about frontier methods for detecting RNA and virus sequences in plants, as well as how to prepare, train, and test models using bioinformatics datasets. It also provided insights into the background and challenges of current technologies, such as high mutation rates and host material dominance. This knowledge allows me to explore applying these techniques to detect viruses in species beyond the plants mentioned in the study. Additionally, I gained a better understanding of how machine learning can be applied in various aspects of bioinformatics.

### **Suggestions**

I suggest testing the methods on additional species not included in the current dataset to evaluate their generalizability. Additionally, comparing the results with a broader range of alternative methods could help better demonstrate the model's performance. Incorporating transfer learning could also be beneficial in improving predictions for unrelated plant species, enhancing the model's adaptability.