

Article

# Deep learning based on stacked sparse autoencoder applied to viral genome classification of SARS-CoV-2 virus

Maria G. F. Coutinho<sup>1,†</sup>, Gabriel B. M. Câmara<sup>1,†</sup>, Raquel de M. Barbosa<sup>1,2,†</sup> and Marcelo A. C. Fernandes<sup>1,3,†,\*</sup>

<sup>1</sup> Laboratory of Machine Learning and Intelligent Instrumentation, Federal University of Rio Grande do Norte, Natal 59078-970, Brazil.;

<sup>2</sup> Laboratory of Drug Development, Department of Pharmacy, Federal University of Rio Grande do Norte, Natal 59078-970, Brazil.; m.g.barbosafernandes@gmail.com

<sup>3</sup> Department of Computer Engineering and Automation, Federal University of Rio Grande do Norte, Natal, RN, 59078-970, Brazil.; mfernandes@dca.ufrn.br

\* Correspondence: mfernandes@dca.ufrn.br

† These authors contributed equally to this work.

**Abstract:** Since December 2019, the world has been intensely affected by the COVID-19 pandemic, caused by the SARS-CoV-2 virus, first identified in Wuhan, China. In the case of a novel virus identification, the early elucidation of taxonomic classification and origin of the virus genomic sequence is essential for strategic planning, containment, and treatments. Deep learning techniques have been successfully used in many viral classification problems associated with viral infections diagnosis, metagenomics, phylogenetic, and analysis. This work proposes to generate an efficient viral genome classifier for the SARS-CoV-2 virus using the deep neural network (DNN) based on stacked sparse autoencoder (SSAE) technique. We performed four different experiments to provide different levels of taxonomic classification of the SARS-CoV-2 virus. The confusion matrix presented the validation and test sets and the ROC curve for the validation set. In all experiments, the SSAE technique provided great performance results. In this work, we explored the utilization of image representations of the complete genome sequences as the SSAE input to provide a viral classification of the SARS-CoV-2. For that, a dataset based on  $k$ -mers image representation, with  $k = 6$ , was applied. The results indicated the applicability of using this deep learning technique in genome classification problems.

**Keywords:** SARS-CoV-2; COVID-19; Deep Learning; Stacked Sparse Autoencoder; Viral classification

## 1. Introduction

Since the emergence of the SARS-CoV-2 virus at the end of 2019, many works are been developed aiming to provide more comprehension about this novel virus. In March 2020, the World Health Organization (WHO) raised the level of contamination to the COVID-19 pandemic, due to its geographical spread across several countries. On July 9, 2021, the disease had registered more than 185 million confirmed cases, and more than 4 million confirmed deaths. In the case of a novel virus identification, the early elucidation of taxonomic classification and origin of the virus genomic sequence is essential for strategic planning, containment, and treatments of the disease [1–3].

One of the fields of research in the bioinformatics area is the analysis of genomic sequences. In the last years, many strategies based on alignment-free methods have been explored as an alternative for the alignment-based methods, considering the limitations of the second approach. Alignment-based programs assume that homologous sequences comprise a series of linearly arranged and more or less conserved sequence stretches, which is not always the case in the real world [4].

Among the alignment-free methodologies, there are some models based on deep learning (DL) techniques, that can provide significant performance in applications of genome analysis [5–7]. Deep neural networks (DNN) can improve prediction accuracy by discovering relevant features of high complexity [7].

**Citation:** Coutinho, M. G. F.; Câmara, G. B. M.; Barbosa, R. de M.; Fernandes, M. A. C.; Deep learning based on stacked sparse autoencoder applied to viral genome classification of SARS-CoV-2 virus. *Preprints* 2021, 1, 0. <https://doi.org/>

Received:

Accepted:

Published:

**Publisher's Note:** MDPI stays neutral with regard to jurisdictional claims in published maps and institutional affiliations.

Figure 1 presents the genome analysis stages and how deep learning integrates this process. The genome analysis stages include the primary analysis, the secondary analysis, and the tertiary analysis. The primary and secondary analysis compose the genome sequencing. The primary analysis receives the biological sample and generates genomic data information, called “reads”, after the processing by the sequencer machine. Then, the secondary analysis processes the reads and produces the complete genome sequence. Lastly, the tertiary analysis provides the genome interpretation, which can be performed for many algorithms and techniques [8–10]. The deep learning techniques have been successful used for the tertiary analysis in many viral classification problems associated with the diagnosis of viral infections, metagenomics, pharmacogenomics, and others [11–15].

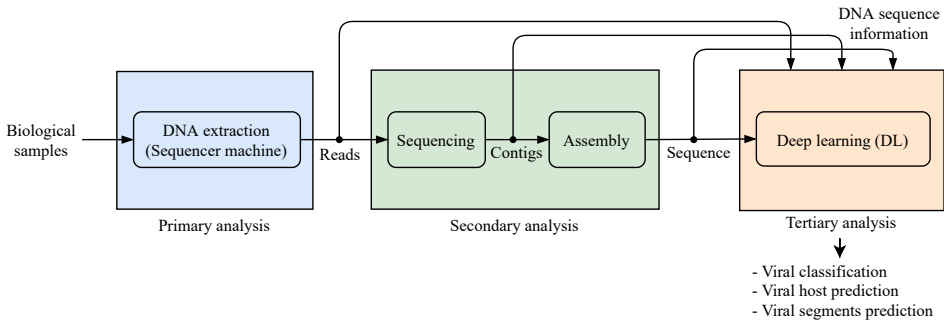


Figure 1. Genome analysis stages with deep learning.

Figure 2 shows the steps of the tertiary analysis using DL, that are the mapping and processing stages. The mapping stage receives the DNA sequence information, that can be the reads, contigs, or the whole genome sequence, and maps this data into a feature space. Various mapping strategies have been present in the works from the state of the art, such as one-hot encoding [13,16–18], number representation [11,12], digital signal processing [19], and other strategies, including multiple mapping strategies applied sequentially [20,21]. The processing stage consists of the utilization of a DNN to perform classification, prediction, and other assumptions about the genome information.

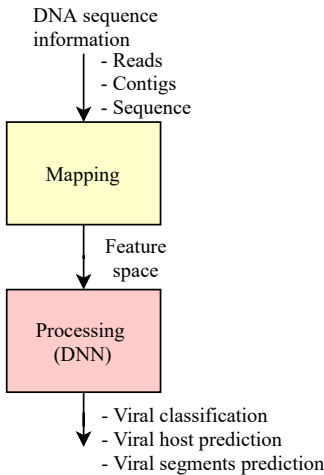


Figure 2. Stages of viral genome analysis using deep learning.

The mapping stage is crucial for the performance of the processing stage. The genome sequence length varies by the type of the virus. Since the DNN only receive a fixed-size input, some researchers have not been using the whole or long sequence length. Nevertheless, longer sequences contain more information and thus are more

convenient to make predictions [17]. In this work, we will explore the utilization of the whole genome sequence mapped by image representation for the use as the DNN input in order to provide viral classification.

Recently works in literature have been applying deep learning as tertiary analysis such as viral prediction, viral host prediction, and viral segments prediction [11–19,22–30].

Tables 1 and 2 present some works from the state of the art that applied DNNs in order to analyse viral genome sequences. Table 1 details the focus of each work as the biology name, the group, the aim, indicates if the proposal was or was not applied for the COVID-19 and present the DNN used. The DNNs applied in those references are divide into 5 groups (CNN+FC, LSTM+FC, BLSTM+FC, BLSTM+CNN+FC, CNN+BLSTM+FC), as we show in the last column of Table 1. Table 2 shows the details about the input and the output of the DNN, besides the biology fields and the bioinformatics area.

In the work presented in [11] was proposed a viral genome deep classifier (VGDC), the first viral genome subtyping based on deep learning techniques found in the literature. Their approach uses a Convolutional Neural Network (CNN) with 25 layers to classify several groups of viruses in subtypes. For the tests, were used five different datasets, each one containing genomes sequences of a specific type of virus. The whole virus genome sequence was used as the input to the network, where the corresponding ASCII code represented each nucleotide. The results indicated that the VGDC was able to achieve better results in comparison with previous works from the state of the art.

In [12] was proposed an approach to assist the tests in the detection of SARS-CoV-2, based on the use of DL techniques. For this, a CNN architecture with 4 layers was used to extract characteristics of the virus genomes, as well as to classify SARS-CoV-2 among Coronavirus type viruses. As presented in [11], the CNN received as input the whole virus genome sequences. The nucleotides were mapped in numerical values (C = 0.25, T = 0.50, G = 0.75, A = 1.0). Missing entries received a value of 0.0. The experiments showed that the CNN was able to correctly identify the sequences even in cases where the noise was added to the genome, reaching accuracies between 0.9674 (with noise) and 0.9875 (without noise). Through the results, the authors also identified a sequence as exclusive for the SARS-CoV-2 virus. They proposed the use of this sequence as a primer for PCR tests.

In [13], was proposed an approach to provide viral classification using the contigs (fragments of the genome sequence) and two different reverse-complement (RC) neural networks architectures: a RC-CNN and a RC-LSTM. These models were also applied to the SARS-CoV-2 virus.

In works presented in [14] and [15], a taxonomic classification for metagenomics applications is proposed. Both works used segments of genome (reads) with DL input (see Figure 1), and the output is the number of the classes. In [14], it was proposed two DL models, one to classify species, and another to classify genus. In [15], a hierarchical taxonomic classification for viral metagenomic data via DL, called CHEER, was proposed. Similar to the work proposed in [14], the CHEER framework classifies the genus, family, and genus.

Proposals presented in [16], [17] and [23] used the contigs with DL input for viral prediction, and classification. In [16], and [17] a DL virus identification framework was proposed and both cases try to recognize if the input is a virus or not.

In work from [16], called ViraMiner, was proposed an approach to detect the presence of viruses on raw metagenomic contigs from different human samples. They used a CNN architecture with two different convolutional branches (pattern and frequency branch) in order to extract relevant features. The outputs of these branches are concatenated and inserted into the fully connected (FC) layer. The ViraMiner output produces a single value that indicates the likelihood of the sequence belonging to the virus class.

Table 1: State of the art references - Part 1.

Biology name	Group		Aim	Ref.	COVID-19	DNN
Genome prediction or sequence classification	Genome classification (taxonomic classification)	Viral classification	Viral Subtyping	[11]	-	CNN+FC
			Primer design	[12]	Yes	CNN+FC
			Identified virus sequence	[13]	Yes	LSTM+FC CNN+FC
			Taxonomic classification	[14] [15]	- -	CNN+FC BLSTM+FC
			Identified virus sequence	[17] [16]	- -	CNN+FC CNN+FC
	Genome prediction	Viral prediction	Identified phage, chromossomes, plasmid	[23]	-	CNN+FC
			Predicting viruses among several hosts	[18] [22]	- Yes	BLSTM+CNN+FC CNN+FC
	Host prediction	Viral host classification				
Genome segments prediction	Genome segments classification	Viral segments classification	Prediction specific regions in the genome	[19] [24] [25] [26]	- - - -	CNN+FC CNN+BLSTM+FC CNN+BLSTM+FC CNN+BLSTM+FC

Table 2: State of the art references - Part 2.

Input	Output	Ref.	Biology fields	Bioinformatics
The DNA or cDNA (RNA virus) of the virus. The whole or part of the genome is used.	Number of the classes	[11] [12] [13] [14] [15]	Metagenomics Diagnosis of viral infections Pharmacogenomics	Free alignments techniques
		[17] [16] [23]	Metagenomics Phylogenetic analysis	
		[18]	Metagenomics Phylogenetic analysis	
		[22]	Metagenomics	
		[19] [24] [25] [26]	Transcriptome Analysis Gene expression analysis	
	Score			
	Binary output			
	Score			
	Number of the classes			
	Score			

In the proposal presented in [17], called DeepVirFinder, the output is a score between 0 and 1 for a binary classification between virus and prokaryote. They fragmented the genomes into non-overlapping sequences of different sizes (150, 300, 500, 1000, and 3000 bp). The sequences were mapped for the network input using the one-hot encoding method. Since they increase the length of the input, i.e. the sequence fragment, they achieve better performance results, which was measured by the area under the receiver operating characteristic curve (AUROC). The maximum AUROC achieved was 0.98 for the 3000 bp fragment.

The work presented in [23] identifies metagenomic fragments as phages, chromosomes or plasmids using the CNN technique. The experiments were performed using artificial contigs and real metagenomic data. The network output, provided by a softmax layer, consists of 3 scores that indicate the probability that each fragment belongs to a specific class.

In the works from [22] and [18] are present DL architectures for host prediction and classification. [22] used a CNN to provide host and infectivity prediction of SARS-CoV-2 virus. In [18] was proposed an approach to predict viral host from three different virus species (influenza A virus, rabies lyssavirus and rotavirus A) from the whole or only fractions of a given viral genome.

In the works from [19], [24], [25] and [26] were proposed methodologies to predict or classify specific regions in the genome sequence. [19] presented a methodology for the classification of three different functional genome types: coding regions, long noncoding regions, and pseudogenes in genomic data. They used a digital signal processing (DSP) methods, called Genomic signal processing (GSP), that converts the nucleotide sequence into a graphical representation of the information contained in the sequence. A CNN with 19 layers was used to perform the classification results.

The authors in [24] proposed a DL framework to identify similar patterns in DNA N6-methyladenine (6mA) sites prediction. This framework, called Deep6mA, is composed of a CNN to extract high-level features in the sequence and a Bi-directional LSTM (BLSTM) to learn dependence structure along the sequence, besides a fully connected layer that determines whether the site is a 6mA site.

In [25] was provided a method based on CNN and BLSTM for exploring the RNA recognition patterns of the CCCTC-binding factor (CTCF) and identify candidate lncRNAs binding. The experiments conducted with two different datasets (human U2OS and mouse ESC) were able to predict CTCF-binding RNA sites from nucleotide sequences. Moreover, [26] propose a computational prediction approach for DNA-protein binding based on CNN and BLSTM.

We intend to provide viral classification using the whole genome sequences, as presented in [11] and [12]. However, in these works were used the length of the longest genome sequence of the dataset as the input of the DNN. So, it was necessary to add some padding for the missing entries. In this work, we will explore the utilization of  $k$ -mers image representation of the complete genome sequences as the DNN input, which will feasibly the use of genome sequences of any length and enable the use of smaller network inputs. The  $k$ -mers representation was used in many works that provide genome sequence classification, as presented in [31], which explores the spectral sequence representation based on  $k$ -mers occurrences. However, that work doesn't explore the  $k$ -mers image representation.

We also explore the utilization of the stacked sparse autoencoder (SSAE) technique as an efficient viral genome classifier. The SSAE has been successfully applied in many biomedical works from the state of the art [6,32–34]. We performed some experiments to provide various levels of taxonomic classification of the SARS-CoV-2 virus, similar to the proposed experiments in [35], using the SSAE technique with a dataset of  $k$ -mers images representations, available on [36].

## 2. Materials and Methods

### 2.1. Dataset

For the experiments, we used a  $k$ -mers representation dataset of SARS-CoV-2 genome, available on [36]. This dataset is composed of 1557 virus instances of SARS-CoV-2, as also, a data stream of 11540 viruses from the Virus-Host DB dataset and the other three Riboviria viruses from NCBI (Betacoronavirus RaTG13, bat-SL-CoVZC45, and bat-SL-CoVZXC21). It also provides  $k$ -mers image representation of all data. The  $k$ -mers images were used to perform the experiments for this work.

Each  $d$ -th sequence, stored in dataset, is expressed by

$$\mathbf{s}_d = [s_{d,1}, \dots, s_{d,n}, \dots, s_{d,N_d}] \quad (1)$$

where  $N_d$  is the length of  $d$ -th sequence and  $s_{d,n}$  is the  $n$ -th nucleotide of the sequence. Each  $n$ -th  $s_{d,n}$  can be characterized as a symbol belonging to an alphabet of 4 possible symbols expressed by set  $\{A, T, C, G\}$  for DNA or by set  $\{A, U, C, G\}$  for RNA, that is,

$$s_{d,n} \in (\{A, T, C, G\} \cup \{A, U, C, G\}). \quad (2)$$

In  $k$ -mers representation, each  $d$ -th nucleotide sequence,  $\mathbf{s}_d$ , is grouped in  $k$ -mers sub-sequences [37,38] that can be expressed as

$$\mathbf{H}_d = \begin{bmatrix} \mathbf{h}_{d,1} \\ \mathbf{h}_{d,2} \\ \vdots \\ \mathbf{h}_{d,i} \\ \vdots \\ \mathbf{h}_{d,N_d-k} \\ \mathbf{h}_{d,N_d-k+1} \end{bmatrix} = \begin{bmatrix} s_{d,1} & \cdots & s_{d,k} \\ s_{d,2} & \cdots & s_{d,k+1} \\ \vdots & \ddots & \vdots \\ s_{d,i} & \cdots & s_{d,i+k} \\ \vdots & \ddots & \vdots \\ s_{d,N_d-k} & \cdots & s_{d,N_d-1} \\ s_{d,N_d-k+1} & \cdots & s_{d,N_d} \end{bmatrix} \quad (3)$$

where the matrix  $\mathbf{H}_d$  stores the  $k$ -mers associated with each  $d$ -th sequence  $\mathbf{s}_d$ . The  $k$ -mers representations are based in each  $d$ -th matrix  $\mathbf{H}_d$  and the matrix  $\mathbf{\Gamma}$ , call here as symbol matrix. The symbol matrix is expressed as

$$\mathbf{\Gamma} = \begin{bmatrix} \gamma_1 \\ \vdots \\ \gamma_i \\ \vdots \\ \gamma_M \end{bmatrix} = \begin{bmatrix} \gamma_{1,1} & \cdots & \gamma_{1,k} \\ \vdots & \ddots & \vdots \\ \gamma_{i,1} & \cdots & \gamma_{i,k} \\ \vdots & \ddots & \vdots \\ \gamma_{M,1} & \cdots & \gamma_{M,k} \end{bmatrix} \quad (4)$$

where each element  $\gamma_{i,j} \in (\{A, T, C, G\} \cup \{A, U, C, G\})$ . The symbol matrix,  $\mathbf{\Gamma}$ , stores all  $M$  possibilities of the  $k$ -mers, where

$$M = 4^k. \quad (5)$$

The  $k$ -mers count 1D representation can be expressed as

$$\mathbf{c}_d = [c_{d,1}, \dots, c_{d,i}, \dots, c_{d,M}] \quad (6)$$

where

$$c_{d,i} = \sum_{j=1}^M \sum_{v=1}^{N-k+1} B_{d,j,v} \quad (7)$$



and

$$B_{d,j,v} = \begin{cases} 0 & \text{for } \gamma_j \neq \mathbf{h}_{d,v} (\exists u = 1, \dots, k : \gamma_{j,u} \neq h_{d,v,u}) \\ 1 & \text{for } \gamma_j = \mathbf{h}_{d,v} (\forall u = 1, \dots, k : \gamma_{j,u} = h_{d,v,u}) \end{cases} \quad (8)$$

Table 3 shows a example of the  $k$ -mers count 1D representation values (with  $k = 2$ ) for SARS-CoV-2 from China-Wuhan (ID: LR757995), USA-MA (ID: MT039888), Brazil (ID: MT126808), and Italy (ID: MT066156). The dataset provide in [36] has  $k$ -mers count 1D representation for  $k = 2, \dots, 6$ .

The  $k$ -mers count 2D representation for each  $d$ -th sequence,  $\mathbf{s}_d$ , is described by

$$\mathbf{\Lambda}_d = \begin{bmatrix} \lambda_{d,1,1} & \cdots & \lambda_{d,1,L} \\ \vdots & \ddots & \vdots \\ \lambda_{d,i,1} & \cdots & \lambda_{d,i,L} \\ \vdots & \ddots & \vdots \\ \lambda_{d,L,1} & \cdots & \lambda_{d,L,L} \end{bmatrix} = \begin{bmatrix} c_{d,1} & \cdots & c_{d,L} \\ \vdots & \ddots & \vdots \\ c_{d,(i-1) \times L + 1} & \cdots & c_{d,i \times L} \\ \vdots & \ddots & \vdots \\ c_{d,M-L+1} & \cdots & c_{d,M} \end{bmatrix} \quad (9)$$

where

$$L = \sqrt{M} = \sqrt{2^k}. \quad (10)$$

Finally, the  $k$ -mers image representation, for each  $d$ -th sequence, can be represented as

$$\mathbf{\Phi}_d = \begin{bmatrix} \phi_{d,1,1} & \cdots & \phi_{d,1,L} \\ \vdots & \ddots & \vdots \\ \phi_{d,i,1} & \cdots & \phi_{d,i,L} \\ \vdots & \ddots & \vdots \\ \phi_{d,L,1} & \cdots & \phi_{d,L,L} \end{bmatrix} \quad (11)$$

where  $\phi_{d,i,j}$  represents each pixel associated with  $d$ -th image  $\mathbf{\Phi}_d$ . Each pixel,  $\phi_{d,i,j}$ , is be expressed as

$$\phi_{d,i,j} = \left\lfloor \frac{2^b - 1}{\max\{\mathbf{\Lambda}_d\}} \times \lambda_{d,i,j} \right\rfloor \quad (12)$$

where  $\max\{\cdot\}$  is the maximum value in  $d$ -th matrix  $\mathbf{\Lambda}_d$ ,  $\lfloor \cdot \rfloor$  is the greatest integer less than or equal, and  $b$  is number of bits associated with the image pixels. Figure 3 show the  $k$ -mers image representation, matrix  $\mathbf{\Phi}$ , (with  $k = 6$  and  $b = 8$ ) for Geminiviridae (ID: HE616777), Alphacoronavirus (ID: JQ410000), and SARS-CoV-2 (Betacoronavirus) from China-Wuhan (ID: LR757995) and Brazil (ID: MT126808).

In this work, we used  $k$ -mers image representation with  $k = 6$ . In the work presented in [16], the 6-mers reached the best performance in comparisons with other values of  $k$  (3, 4, 5 and 7). The data of each experiment was partitioned using the holdout method, which splits the data into a training set and a validation set at random. We used the proportion of 80% for the training set and 20% for the validation set. Each class data was split respecting these percentages. The SARS-CoV-2  $k$ -mers images were used only for the test set.

## 2.2. DNN Architecture

All experiments were performed using the SSAE technique. In these models each hidden layer is composed of an individually trained sparse autoencoder in an unsupervised way. A sparse autoencoder is an autoencoder whose training involves a sparse penalty, which functions as a regularizing term added to the loss function [39]. The autoencoder (AE) is a DL technique specialized in dimensionality reduction and feature extraction. The AE output can provide the reconstruction of the input information. These networks are composed of three layers: an input, a hidden and an output. The encoder is formed by the input and hidden layers, and the decoder is formed by the hidden and

Table 3: Examples of  $k$ -mers count 1D representation values (with  $k = 2$ ) for SARS-CoV-2.

$k$ -mers ( $k = 2$ )	China-Wuhan (ID: LR757995)	USA-MA (ID: MT039888)	Brazil (ID: MT126808)	Italy (ID: MT066156)
AA	2862	2859	2853	2847
AC	2022	2022	2022	2022
AG	1741	1741	1742	1742
AT	2306	2309	2309	2308
CA	2085	2082	2084	2082
CC	886	888	888	888
CG	439	439	440	439
CT	2080	2081	2080	2082
GA	1612	1612	1612	1611
GC	1167	1167	1169	1168
GG	1092	1093	1092	1092
GT	1990	1990	1988	1989
TA	2373	2378	2377	2378
TC	1415	1412	1413	1413
TG	2589	2589	2587	2587
TT	3212	3217	3219	3216

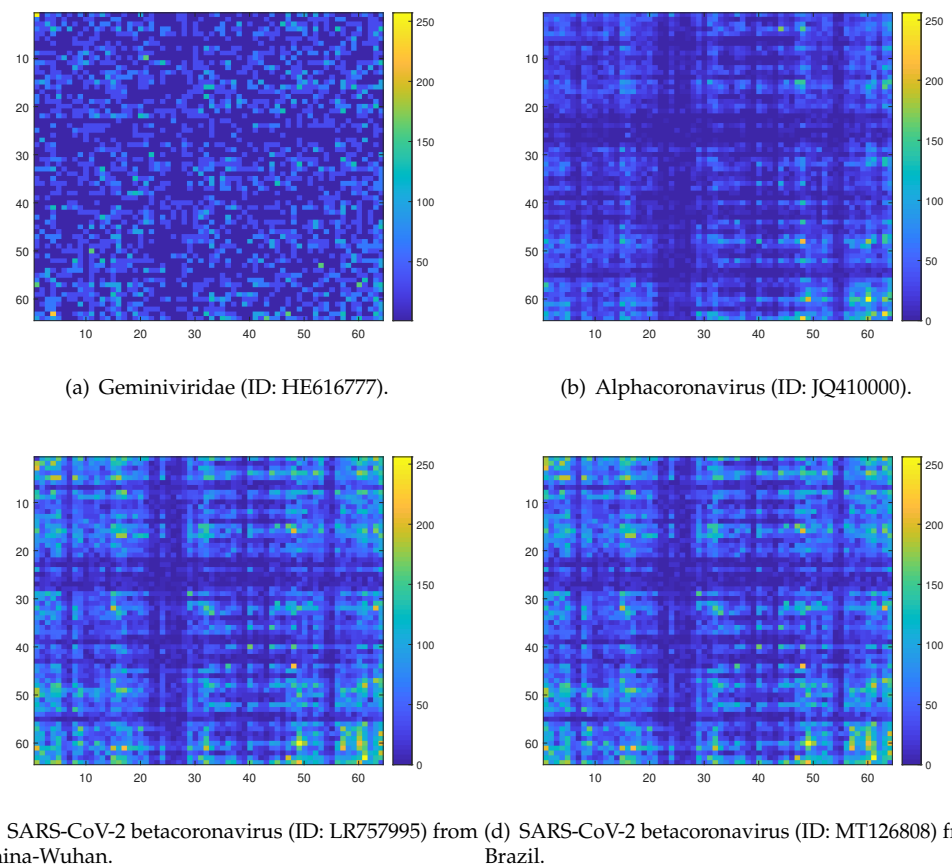
Table 4: Examples of  $k$ -mers count 2D representation values (with  $k = 2$ ) for SARS-CoV-2.

China-Wuhan (ID: LR757995)					USA-MA (ID: MT039888)				
$\Lambda_{17} =$	2862	2022	1741	1741	$\Lambda_{32} =$	2859	2022	1741	1741
	2085	886	439	439		2082	888	439	439
	1612	1167	1092	1092		1612	1167	1093	1093
	2373	1415	2589	2589		2378	1412	2589	2589
Brazil (ID: MT126808)					Italy (ID: MT066156)				
$\Lambda_{52} =$	2853	2022	1742	1742	$\Lambda_{79} =$	2853	2022	1742	1742
	2084	888	440	440		2084	888	440	440
	1612	1169	1092	1092		1612	1169	1092	1092
	2377	1413	2587	2587		2377	1413	2587	2587

output layers [39]. For the output layer, we used a softmax layer, where the number of neurons consists of the number of classes of the experiment. Figure 4 illustrates the DL SSAE with  $P$  inputs,  $K$  hidden layers, and a output layer. Each  $i$ -th hidden layer has  $Q_i$  neurons and the output layer has  $U$  neurons. Functions  $\varphi(\cdot)$  and  $f(\cdot)$  are the action functions in each  $p$ -th neuron (in each  $i$ -th hidden layer) and each  $u$ -th neuron in output layer, respectively.

For all experiments, the network architecture used three hidden layers ( $K = 3$ ), containing 3000 neurons in the first hidden layer,  $Q_1$ , 1000 in the second hidden layer,  $Q_2$ , and 500 in the third hidden layer  $Q_3$ . For input of the SSAE, it was used  $k$ -mers images, with  $k = 6$ , generating images, matrix  $\Phi$ , with  $64 \times 64$  pixels (based on Equation

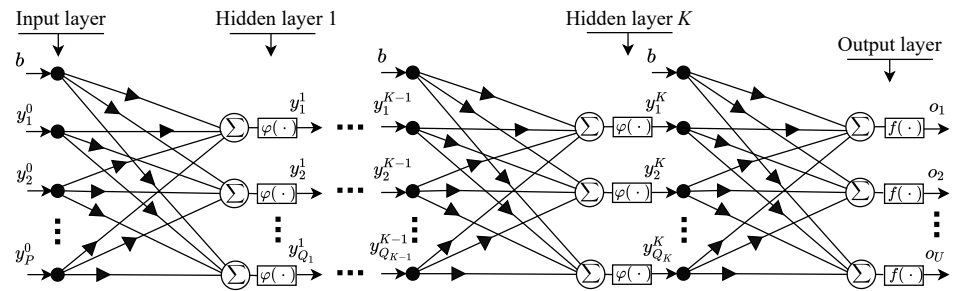




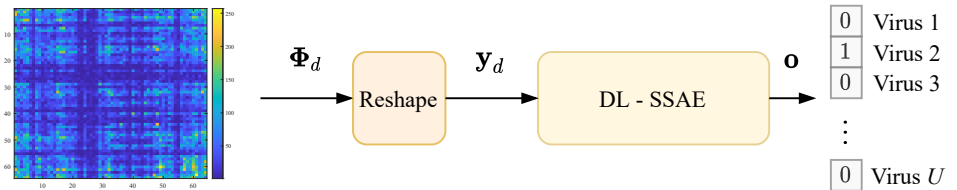
**Figure 3.** Examples of  $k$ -mers images representation with  $k = 6$ . Based on Equation 10,  $L = 64$  and each image, matrix  $\Phi$  (see Equation 11), is composed by  $64 \times 64$  pixels with  $b = 8$  (see Equation 12).

10,  $L = \sqrt{4^6} = 64$ ). Each  $d$ -th image,  $\Phi_d$ , associated with a  $d$ -th viral genome sequence is reshaped into a vector expressed by

$$\mathbf{y}_d = \begin{bmatrix} y_{d,1}^0 \\ y_{d,2}^0 \\ \vdots \\ y_{d,i-1}^0 \\ y_{d,i}^0 \\ y_{d,i+1}^0 \\ \vdots \\ y_{d,p-1}^0 \\ y_{d,p}^0 \end{bmatrix} = \begin{bmatrix} \phi_{d,1,1} \\ \vdots \\ \phi_{d,L,1} \\ \phi_{d,1,2} \\ \vdots \\ \phi_{d,L,2} \\ \vdots \\ \phi_{d,1,L} \\ \vdots \\ \phi_{d,L,L} \end{bmatrix} \quad (13)$$



**Figure 4.** Deep learning stacked sparse autoencoder architecture (DL-SSAE).



**Figure 5.** Viral classification process using  $k$ -mers images representation with the DL-SSAE.

with  $P = 64 \times 64 = 4096$  values and applied to the SSAE. The number of neurons in output layer,  $U$ , is defined by the number of different viruses in a specific taxonomic level such as family, genus, realm and other. The output can be expressed by

$$\mathbf{o} = \begin{bmatrix} o_1 \\ \vdots \\ o_u \\ \vdots \\ o_U \end{bmatrix} \quad (14)$$

where each  $u$ -th output,  $o_u$ , represents a specific virus in a taxonomic level classification and is defined by

$$o_u = \begin{cases} 1 & \text{if } \mathbf{y}_d \text{ is the } u\text{-th virus} \\ 0 & \text{otherwise} \end{cases} \quad (15)$$

Figure 5 illustrates how the sequence information is passed through the DL-SSAE to perform the viral classification. The DL-SSAE input was normalized in the range of 0 to 1. First, the SSAE receives the training set as input to perform the training phase. Then, the validation set, which only contains samples that were not applied in the training phase, is used to identify the capacity of generalization of the DNN. After the network validation, the SSAE was applied for the test set, which only contains SARS-CoV-2 sequences. The SARS-CoV-2  $k$ -mers images were not used for the training phase of the SSAE.

The SSAE was implemented in the Matlab platform (License 596681) [40], adopting the deep learning toolbox. All network was trained with the Scaled Conjugate Gradient (SCG) algorithm. The loss function used for the training in each AE was the Mean Squared Error with L2 and Sparsity Regularizers, that can be expressed as

$$E = \frac{1}{I} \sum_{i=1}^I \sum_{u=1}^U (o_{ui}^{ref} - o_{ui})^2 + \lambda \times \Omega_{weights} + \beta \times \Omega_{sparsity}, \quad (16)$$

where  $I$  is the number of training examples,  $U$  is the number of classes,  $\Omega_{weights}$  is the L2 regularization term,  $\lambda$  is the coefficient for the L2 regularization term,  $\Omega_{sparsity}$  is the sparsity regularization term, and  $\beta$  is the coefficient for the sparsity regularization term.

The loss function applied for the softmax layer was the Cross-Entropy. In this work, after the training in each layer, the fine-tuning was performed, which retrained all the stacked network in a supervised way in order to improve the classification results. The fine-tuning process also used the Cross-Entropy as the loss function, as in the softmax layer.

### 3. Results and discussion

We performed four different experiments to provide different levels of taxonomic classification of the SARS-CoV-2 virus, similar to the experimental methodology present in [35]. The details about the data and the network architecture used in each experiment are shown in Table 5. The SSAE architecture was chosen by the observation of the MSE obtained with the reconstruction of the validation set in each AE. In order to validate the proposed idea of this work, the results are present by the confusion matrix for the validation and test sets. We also measured the performance of the viral classifier proposed with some popular classification metrics, as precision, recall, F1-score, and specificity. The precision value measure the percentages of all the examples predicted to belong to each class that are correctly classified, which corresponds to the positive predictive value. The recall, also called sensibility, corresponds to the percentages of all the examples belonging to each class that are correctly classified, which is the true positive rate. The F1-score can be interpreted as a weighted average of the precision and recall, and the specificity indicates the true negative rate. The column on the far right of each confusion matrix shows the percentages of precision per class, and the row at the bottom of each confusion matrix shows the percentages of recall per class. The cell in the bottom right of the plot of each confusion matrix shows the overall accuracy. Besides, for the validation set we also present the receiver operating characteristic (ROC) curve. The ROC curve measures the classification performance, that is the true positive rate and the false positive rate of each class, at various thresholds settings.

In Experiment 1, we intended to classify the viruses in 14 different classes, as presented in Table 5, which consists of 10 families (Adenoviridae, Anelloviridae, Circoviridae, Geminiviridae, Genomoviridae, Microviridae, Papillomaviridae, Parvoviridae, Polyomaviridae and Tolecusatellitidae), three orders (Caudovirales, Herpesvirales and Ortervirales) and Riboviria realm. The Riboviria class contains various families that belong to the realm Riboviria, including the Coronaviridae family. To ensure data balance, only the classes with at least 100 sequences from the original dataset were considered. For the classes with more than 500 sequences, only 500 sequences were selected at random, except for the Riboviria class, in which was prioritized the Coronaviridae family sequences, to guarantee the correct classification of the test data (SARS-CoV-2 sequences), which is the focus of this work. In this particular case, were selected all Coronaviridae family sequences available in the dataset (206 samples), and the other 294 sequences were select from the rest of the Riboviria data at random. After this balancing, Experiment 1 comprised 3433 samples of virus sequences.

The SSAE architecture used in Experiment 1 was the 4096 – 3000 – 1000 – 500 – 14 architecture. The three AEs were trained for 400 epochs. The softmax layer was trained for 3000 epochs or until reach the minimum gradient ( $< 1 \times 10^{-6}$ ). Lastly, the fine-tuning was performed. For each experiment, the fine-tuning phase uses the same stopping condition as the softmax layer.

The confusion matrix and the ROC curve from the validation set of Experiment 1 are present in Figures 6 and 7, respectively. In Experiment 1, the classification accuracy from the validation set reached 92%. This result is promising, especially considering the challenges of the classification in high-level taxonomies because of the high diversity of the viruses sequences. It is essential to mention that the balancing process may have caused the classification more complicated because some crucial sequences may have been excluded from the dataset. However, this result can be improved in many ways that will be discussed following.

Table 5: Experiments data.

Experiments	Classes	Number of sequences	SSAE architecture $P - Q_1 - Q_2 - Q_3 - U$
Experiment 1	Adenoviridae	195	4096 – 3000 – 1000 – 500 – 14
	Anelloviridae	114	
	Caudovirales	500	
	Circoviridae	243	
	Geminiviridae	500	
	Genomoviridae	115	
	Herpesvirales	136	
	Microviridae	102	
	Ortervirales	214	
	Papillomaviridae	354	
	Parvoviridae	168	
	Polyomaviridae	142	
	Riboviria	500	
	Tolecusatellitidae	150	
Experiment 2	Picornaviridae	423	4096 – 3000 – 1000 – 500 – 8
	Caliciviridae	392	
	Coronaviridae	206	
	Potyviridae	232	
	Flaviviridae	217	
	Rhabdoviridae	186	
	Betaflexiviridae	129	
	Reoviridae	111	
Experiment 3	Alphacoronavirus	52	4096 – 3000 – 1000 – 500 – 4
	Betacoronavirus	123	
	Deltacoronavirus	20	
	Gammacoronavirus	9	
Experiment 4	Embecovirus	47	4096 – 3000 – 1000 – 500 – 4
	Merbecovirus	17	
	Nobecovirus	9	
	Sarbecovirus	46	

268 Regarded to the classification performance per class, the precision value presented  
269 in the last column shows that the worse result was obtained from an order class (71.4%  
270 from the Herpesvirales). Among the five worst classification results, two are from order  
271 classes (71.4% and 83.3% from Herpesvirales and Ortervirales, respectively). Since  
272 these classes can contain viruses from many different realms and families, they can  
273 difficult the training process. The Riboviria realm, which is the focus of this work,  
274 reached a classification accuracy of 93%. Analyse the results per classes can give more  
275 understanding about the dataset used and the implications of this dataset for the results,  
276 which is important to make decisions for the next experiments.

277 The confusion matrix from the test set of Experiment 1 is present in Figure 8. In  
278 the test phase of this experiment, all the 1557 sequences of SARS-CoV-2 was correctly  
279 classified as belonging to the Riboviria realm, so the classification accuracy reached  
280 100%.

281 Experiment 2 performs the classification of Riboviria families. As in Experiment 1,  
282 only classes with at least 100 sequences were considered. This experiment includes 1896  
283 sequences separated into eight families (Picornaviridae, Caliciviridae, Coronaviridae,  
284 Potyviridae, Flaviviridae, Rhabdoviridae, Betaflexiviridae and Reoviridae). We used  
285 the 4096 – 3000 – 1000 – 500 – 8 SSAE architecture. The three AEs were trained for 400  
286 epochs each and the softmax layer was trained for 1000 epochs or until reaching the  
287 minimum gradient, as well as the fine-tuning phase.

Output Class	Adenoviridae	39 5.7%	0 0.0%	1 0.1%	1 0.1%	0 0.0%	0 0.0%	2 0.3%	0 0.0%	0 0.0%	0 0.0%	1 0.1%	0 0.0%	0 0.0%	0 0.0%	88.6% 11.4%
	Anelloviridae	0 0.0%	20 2.9%	0 0.0%	0 0.0%	0 0.0%	0 0.0%	0 0.0%	0 0.0%	1 0.1%	0 0.0%	0 0.0%	1 0.1%	0 0.0%	0 0.0%	90.9% 9.1%
	Caudovirales	0 0.0%	0 0.0%	96 14.0%	1 0.1%	0 0.0%	0 0.0%	0 0.0%	0 0.0%	1 0.1%	0 0.0%	1 0.1%	0 0.0%	0 0.0%	0 0.0%	97.0% 3.0%
	Circoviridae	0 0.0%	0 0.0%	0 0.0%	40 5.8%	1 0.1%	2 0.3%	1 0.1%	1 0.1%	0 0.0%	0 0.0%	1 0.1%	0 0.0%	2 0.3%	0 0.0%	83.3% 16.7%
	Geminiviridae	0 0.0%	0 0.0%	0 0.0%	0 0.0%	98 14.3%	1 0.1%	0 0.0%	0 0.0%	0 0.0%	0 0.0%	1 0.1%	1 0.1%	0 0.0%	1 0.1%	96.1% 3.9%
	Genomoviridae	0 0.0%	0 0.0%	0 0.0%	1 0.1%	0 0.0%	20 2.9%	2 0.3%	0 0.0%	0 0.0%	0 0.0%	0 0.0%	0 0.0%	0 0.0%	0 0.0%	87.0% 13.0%
	Herpesvirales	1 0.1%	0 0.0%	0 0.0%	2 0.3%	0 0.0%	0 0.0%	20 2.9%	0 0.0%	0 0.0%	0 0.0%	0 0.0%	1 0.1%	1 0.1%	0 0.0%	80.0% 20.0%
	Microviridae	0 0.0%	0 0.0%	0 0.0%	0 0.0%	0 0.0%	0 0.0%	0 0.0%	19 2.8%	1 0.1%	0 0.0%	0 0.0%	0 0.0%	0 0.0%	0 0.0%	95.0% 5.0%
	Ortervirales	0 0.0%	1 0.1%	0 0.0%	1 0.1%	0 0.0%	0 0.0%	3 0.4%	0 0.0%	35 5.1%	0 0.0%	2 0.3%	0 0.0%	0 0.0%	3 0.4%	77.8% 22.2%
	Papillomaviridae	0 0.0%	0 0.0%	0 0.0%	0 0.0%	0 0.0%	0 0.0%	0 0.0%	0 0.0%	1 0.1%	70 10.2%	1 0.1%	0 0.0%	0 0.0%	0 0.0%	97.2% 2.8%
	Parvoviridae	0 0.0%	1 0.1%	1 0.1%	1 0.1%	0 0.0%	0 0.0%	0 0.0%	0 0.0%	2 0.3%	0 0.0%	27 3.9%	0 0.0%	0 0.0%	1 0.1%	81.8% 18.2%
	Polyomaviridae	0 0.0%	0 0.0%	0 0.0%	0 0.0%	0 0.0%	0 0.0%	0 0.0%	0 0.0%	0 0.0%	0 0.0%	0 0.0%	25 3.6%	0 0.0%	0 0.0%	100% 0.0%
	Riboviria	0 0.0%	0 0.0%	3 0.4%	1 0.1%	1 0.1%	0 0.0%	0 0.0%	0 0.0%	1 0.1%	0 0.0%	0 0.0%	0 0.0%	93 13.6%	0 0.0%	93.9% 6.1%
	Tolecusatellitidae	0 0.0%	0 0.0%	0 0.0%	0 0.0%	0 0.0%	0 0.0%	0 0.0%	0 0.0%	0 0.0%	0 0.0%	0 0.0%	0 0.0%	0 0.0%	29 4.2%	100% 0.0%
		97.5% 2.5%	90.9% 9.1%	95.0% 5.0%	83.3% 16.7%	98.0% 2.0%	87.0% 13.0%	71.4% 28.6%	95.0% 5.0%	83.3% 16.7%	100% 0.0%	79.4% 20.6%	89.3% 10.7%	93.0% 7.0%	96.7% 3.3%	92.0% 8.0%
		Adenoviridae	Anelloviridae	Caudovirales	Circoviridae	Geminiviridae	Genomoviridae	Herpesvirales	Microviridae	Ortervirales	Papillomaviridae	Parvoviridae	Polyomaviridae	Riboviria	Tolecusatellitidae	
		Target Class														

**Figure 6.** Confusion matrix of the validation set from the Experiment 1.

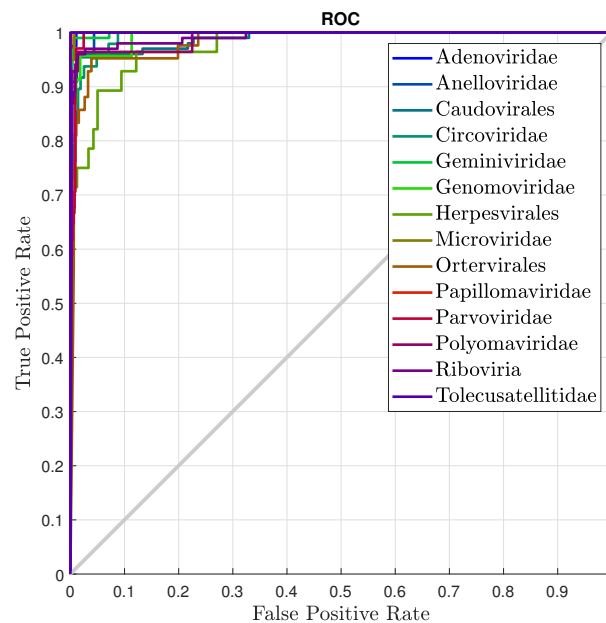
288 The confusion matrix and the ROC curve from the validation set of Experiment 2  
289 are present in Figures 9 and 10, respectively. The classification accuracy from Experiment  
290 2 reached 96.3%. From the 379 sequences applied in this validation, only 11 were  
291 not correctly classified. Besides, the SSAE classified all sequences that belong to the  
292 Coronaviridae family correctly. The ROC curve from Experiment 2 also provides  
293 excellent results.

294 The confusion matrix from the test set of Experiment 2 is present in Figure 11. The  
295 SSAE achieve 100% of classification accuracy, i.e., all SARS-CoV-2 sequences applied in  
296 this experiment were perfectly classified as Coronaviridae family sequences.

297 In Experiment 3 we aim to provide the classification among the Coronaviridae  
298 genera. For this experiment, 204 sequences divided into four genera (Alphacoronavirus,  
299 Betacoronavirus, Deltacoronavirus and Gammacoronavirus) were used. The SSAE  
300 architecture used in this experiment was the 4096 – 3000 – 1000 – 500 – 4 architecture.  
301 The three AEs were trained for 400 epochs each, and the softmax layer was trained for  
302 2000 epochs or until reaching the minimum gradient.

303 Figures 12 and 13 show the resulting confusion matrix and ROC curve from the  
304 Experiment 3, respectively. This experiment achieved 95% of classification accuracy of  
305 the validation set. The classification performance of the model obtained for the Betacoro-  
306 navirus genus was 95.8%. Also, the ROC curve plotted for all classes of Experiment 3  
307 provides satisfactory results.

308 Regarding the test set of Experiment 3, the confusion matrix is present in Figure  
309 14. The test phase of Experiment 3 achieved 98.9% of classification accuracy. In the  
310 validation phase of Experiment 3, the Betacoronavirus genus did not reach the highest  
311 performance, which probably explains these result in the test phase.



**Figure 7.** ROC curve of the validation set from the Experiment 1.

In Experiment 4, we provide the Betacoronaviridae subgenera classification. This test includes 119 genome sequences divided into four classes (Embecovirus, Marbecovirus, Nobecovirus and Sarbecovirus). The SSAE architecture was the same as the architecture used in Experiment 3 (4096 – 3000 – 1000 – 500 – 4), as well as the training parameters.

The confusion matrix and the ROC curve from the validation set of Experiment 4 are present in Figures 15 and 16, respectively. In this experiment, the SSAE achieved the highest classification accuracy (100%), which is reaffirmed for the ROC curve plot.

Figure 15 exposes the confusion matrix from the test set of Experiment 4. In this case, the SSAE achieved 99.9% of classification accuracy, that is equivalent to only one sequence wrong classified.

Table 6 presents the results regarding some popular classification performance metrics obtained from the validation set. The first column of the table indicates the experiment proposed. The second column shows the overall accuracy for each experiment. The precision, recall, F1-score, and specificity are present in the others columns, which were obtained by the average of the values obtained for each class.

Table 6: Classification performance metrics results obtained from the validation set.

Experiment	Accuracy	Precision	Recall	F1 score	Specificity
1	0.920 (92.0%)	0.924 (92.4%)	0.920 (92.0%)	0.931 (93.1%)	0.993 (99.3%)
2	0.963 (96.3%)	0.968 (96.8%)	0.971 (97.1%)	0.962 (96.2%)	0.997 (99.7%)
3	0.950 (95.0%)	0.979 (97.9%)	0.979 (97.9%)	0.955 (95.5%)	0.983 (98.3%)
4	1 (100%)	1 (100%)	1 (100%)	1 (100%)	1 (100%)

All the metrics presented in Table 6 indicate that the viral classifier proposed performs great for all experiments. The highest performance was obtained for the Experiment 4. Besides, Experiments 2 and 3, reached values more than 0.95 for all the metrics evaluated. The classification performance slightly decreased in the Experiment 1, which is acceptable because of the high diversity of the viruses sequences applied.





		Confusion Matrix							
Output Class	Picornaviridae	81 21.4%	1 0.3%	0 0.0%	2 0.5%	1 0.3%	0 0.0%	0 0.0%	95.3% 4.7%
	Caliciviridae	3 0.8%	78 20.6%	0 0.0%	2 0.5%	0 0.0%	0 0.0%	0 0.0%	94.0% 6.0%
	Coronaviridae	0 0.0%	0 0.0%	41 10.8%	0 0.0%	0 0.0%	0 0.0%	0 0.0%	100% 0.0%
	Potyviridae	0 0.0%	0 0.0%	0 0.0%	43 11.3%	0 0.0%	0 0.0%	1 0.3%	97.7% 2.3%
	Flaviviridae	0 0.0%	0 0.0%	0 0.0%	0 0.0%	42 11.1%	0 0.0%	2 0.5%	95.5% 4.5%
	Rhabdoviridae	0 0.0%	0 0.0%	0 0.0%	0 0.0%	0 0.0%	37 9.8%	0 0.3%	97.4% 2.6%
	Betaflexiviridae	0 0.0%	0 0.0%	0 0.0%	0 0.0%	0 0.0%	25 6.6%	1 0.3%	96.2% 3.8%
	Reoviridae	0 0.0%	0 0.0%	0 0.0%	0 0.0%	0 0.0%	0 0.0%	18 4.7%	100% 0.0%
		96.4% 3.6%	98.7% 1.3%	100% 0.0%	91.5% 8.5%	97.7% 2.3%	100% 0.0%	96.2% 3.8%	81.8% 18.2%
		Target Class							
		Picornaviridae	Caliciviridae	Coronaviridae	Potyviridae	Flaviviridae	Rhabdoviridae	Betaflexiviridae	Reoviridae

**Figure 9.** Confusion matrix of the validation set from the Experiment 2.

cross-validation scheme. Besides, we also intend to study data balancing alternatives, based on the analysis of the results presented here.

#### 4. Conclusions

This work presented an alignment-free methodology, based on the stacked sparse autoencoder technique, in order to classify genome sequences of the SARS-CoV-2 virus in various levels of taxonomy (realm, family, genus and subgenus). We explored the utilization of  $k$ -mers image representation of the whole genome sequence, which feasibility the use of genome sequences of any length and enable the use of smaller network inputs. The results were presented by the confusion matrix for the validation and test sets, and the ROC curve for the validation set. All experiments provided great performance results, reaching accuracies between 98.9% and 100% for the test set. These results indicated the applicability of using the stacked sparse autoencoder technique in genome classification problems.

**Author Contributions:** All the authors have contributed in various degrees to ensure the quality of this work. (e.g., Maria G. F. Coutinho, Gabriel B. M. Câmara, Raquel de M. Barbosa and Marcelo A. C. Fernandes conceived the idea and experiments; Maria G. F. Coutinho, Gabriel B. M. Câmara, Raquel de M. Barbosa and Marcelo A. C. Fernandes designed and performed the experiments; Maria G. F. Coutinho, Gabriel B. M. Câmara, Raquel de M. Barbosa and Marcelo A. C. Fernandes analyzed the data; Maria G. F. Coutinho, Gabriel B. M. Câmara, Raquel de M. Barbosa and Marcelo A. C. Fernandes wrote the paper. Marcelo A. C. Fernandes coordinated the project.). All authors have read and agreed to the published version of the manuscript.

**Funding:** This study was financed in part by the Coordenação de Aperfeiçoamento de Pessoal de Nível Superior (CAPES)—Finance Code 001.

**Acknowledgments:** The authors wish to acknowledge the financial support of the Coordenação de Aperfeiçoamento de Pessoal de Nível Superior (CAPES) for their financial support.

**Conflicts of Interest:** The authors declare no conflict of interest.

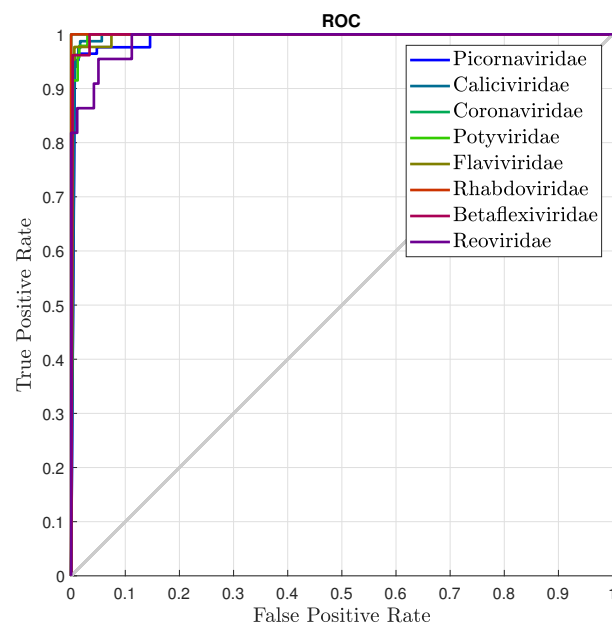
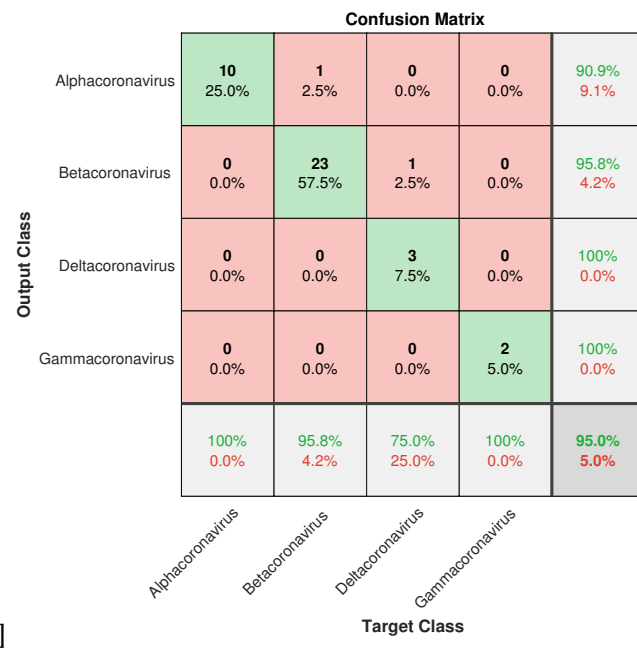


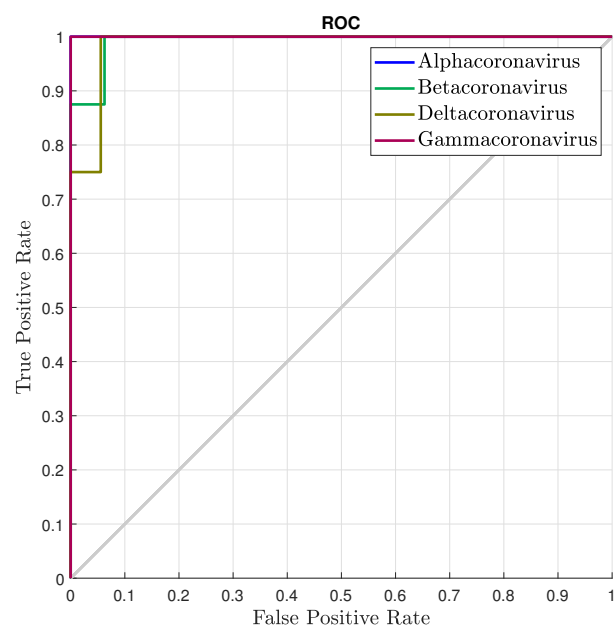
Figure 10. ROC curve of the validation set from the Experiment 2.

		Confusion Matrix							
Output Class	Picornaviridae	0 0.0%	0 0.0%	0 0.0%	0 0.0%	0 0.0%	0 0.0%	0 0.0%	NaN% NaN%
	Caliciviridae	0 0.0%	0 0.0%	0 0.0%	0 0.0%	0 0.0%	0 0.0%	0 0.0%	NaN% NaN%
	Coronaviridae	0 0.0%	0 0.0%	1557 100%	0 0.0%	0 0.0%	0 0.0%	0 0.0%	100% 0.0%
	Potyviridae	0 0.0%	0 0.0%	0 0.0%	0 0.0%	0 0.0%	0 0.0%	0 0.0%	NaN% NaN%
	Flaviviridae	0 0.0%	0 0.0%	0 0.0%	0 0.0%	0 0.0%	0 0.0%	0 0.0%	NaN% NaN%
	Rhabdoviridae	0 0.0%	0 0.0%	0 0.0%	0 0.0%	0 0.0%	0 0.0%	0 0.0%	NaN% NaN%
	Betaflexiviridae	0 0.0%	0 0.0%	0 0.0%	0 0.0%	0 0.0%	0 0.0%	0 0.0%	NaN% NaN%
	Reoviridae	0 0.0%	0 0.0%	0 0.0%	0 0.0%	0 0.0%	0 0.0%	0 0.0%	NaN% NaN%
		NaN% NaN%	NaN% NaN%	100% 0.0%	NaN% NaN%	NaN% NaN%	NaN% NaN%	NaN% NaN%	100% 0.0%
		Picornaviridae	Caliciviridae	Coronaviridae	Potyviridae	Flaviviridae	Rhabdoviridae	Betaflexiviridae	Reoviridae
		Target Class							

Figure 11. Confusion matrix of the test set from the Experiment 2.



**Figure 12.** Confusion matrix of the validation set from the Experiment 3.



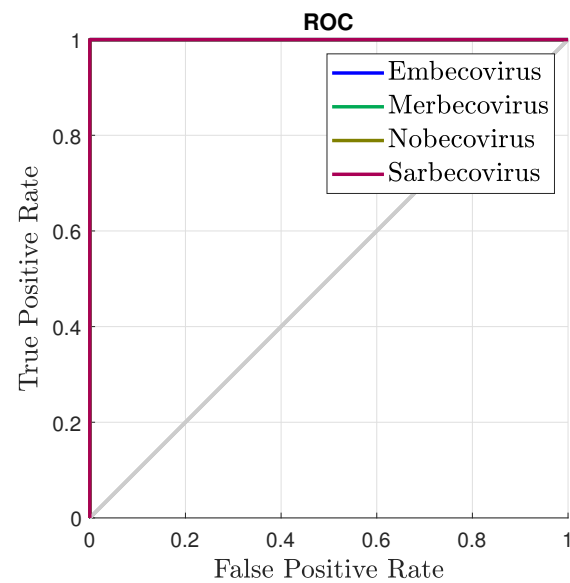
**Figure 13.** ROC curve of the validation set from the Experiment 3.

		Confusion Matrix				
Output Class	Alphacoronavirus	0 0.0%	2 0.1%	0 0.0%	0 0.0%	0.0% 100%
	Betacoronavirus	0 0.0%	1540 98.9%	0 0.0%	0 0.0%	100% 0.0%
	Deltacoronavirus	0 0.0%	0 0.0%	0 0.0%	0 0.0%	NaN% NaN%
	Gammacoronavirus	0 0.0%	15 1.0%	0 0.0%	0 0.0%	0.0% 100%
		NaN% NaN%	98.9% 1.1%	NaN% NaN%	NaN% NaN%	98.9% 1.1%
		Target Class				
		Alphacoronavirus	Betacoronavirus	Deltacoronavirus	Gammacoronavirus	

**Figure 14.** Confusion matrix of the test set from the Experiment 3.

		Confusion Matrix				
Output Class	Embecovirus	9 39.1%	0 0.0%	0 0.0%	0 0.0%	100% 0.0%
	Merbecovirus	0 0.0%	4 17.4%	0 0.0%	0 0.0%	100% 0.0%
	Nobecovirus	0 0.0%	0 0.0%	1 4.3%	0 0.0%	100% 0.0%
	Sarbecovirus	0 0.0%	0 0.0%	0 0.0%	9 39.1%	100% 0.0%
		100% 0.0%	100% 0.0%	100% 0.0%	100% 0.0%	100% 0.0%
		Target Class				
		Embecovirus	Merbecovirus	Nobecovirus	Sarbecovirus	

**Figure 15.** Confusion matrix of the validation set from the Experiment 4.



**Figure 16.** ROC curve of the validation set from the Experiment 4.

		Confusion Matrix				
Output Class	Embecovirus	0 0.0%	0 0.0%	0 0.0%	0 0.0%	NaN% NaN%
	Merbecovirus	0 0.0%	0 0.0%	0 0.0%	1 0.1%	0.0% 100%
	Nobecovirus	0 0.0%	0 0.0%	0 0.0%	0 0.0%	NaN% NaN%
	Sarbecovirus	0 0.0%	0 0.0%	0 0.0%	1556 99.9%	100% 0.0%
		NaN% NaN%	NaN% NaN%	NaN% NaN%	99.9% 0.1%	99.9% 0.1%
		Embecovirus	Merbecovirus	Nobecovirus	Sarbecovirus	
		Target Class				

**Figure 17.** Confusion matrix of the test set from the Experiment 4.



## References

1. Lam, T.T.Y.; Shum, M.H.H.; Zhu, H.C.; Tong, Y.G.; Ni, X.B.; Liao, Y.S.; Wei, W.; Cheung, W.Y.M.; Li, W.J.; Li, L.F.; others. Identifying SARS-CoV-2 related coronaviruses in Malayan pangolins. *Nature* **2020**, pp. 1–6.
2. Andersen, K.G.; Rambaut, A.; Lipkin, W.I.; Holmes, E.C.; Garry, R.F. The proximal origin of SARS-CoV-2. *Nature medicine* **2020**, *26*, 450–452.
3. Graham, R.L.; Baric, R.S. SARS-CoV-2: Combating Coronavirus Emergence. *Immunity* **2020**.
4. Zielezinski, A.; Vinga, S.; Almeida, J.; Karlowski, W.M. Alignment-free sequence comparison: benefits, applications, and tools. *Genome biology* **2017**, *18*, 186.
5. Zou, J.; Huss, M.; Abid, A.; Mohammadi, P.; Torkamani, A.; Telenti, A. A primer on deep learning in genomics. *Nature genetics* **2019**, *51*, 12–18.
6. Tang, B.; Pan, Z.; Yin, K.; Khateeb, A. Recent advances of deep learning in bioinformatics and computational biology. *Frontiers in genetics* **2019**, *10*, 214.
7. Eraslan, G.; Avsec, Ž.; Gagneur, J.; Theis, F.J. Deep learning: new computational modelling techniques for genomics. *Nature Reviews Genetics* **2019**, *20*, 389–403.
8. Pareek, C.S.; Smoczynski, R.; Tretyn, A. Sequencing technologies and genome sequencing. *Journal of applied genetics* **2011**, *52*, 413–435.
9. Pabinger, S.; Dander, A.; Fischer, M.; Snajder, R.; Sperk, M.; Efremova, M.; Krabichler, B.; Speicher, M.R.; Zschocke, J.; Trajanoski, Z. A survey of tools for variant analysis of next-generation genome sequencing data. *Briefings in bioinformatics* **2014**, *15*, 256–278.
10. Posada-Céspedes, S.; Seifert, D.; Beerenwinkel, N. Recent advances in inferring viral diversity from high-throughput sequencing data. *Virus research* **2017**, *239*, 17–32.
11. Fabijańska, A.; Grabowski, S. Viral Genome Deep Classifier. *IEEE Access* **2019**, *7*, 81297–81307.
12. Lopez-Rincon, A.; Tonda, A.; Mendoza-Maldonado, L.; Claassen, E.; Garssen, J.; Kraneveld, A.D. Accurate identification of sars-cov-2 from viral genome sequences using deep learning. *bioRxiv* **2020**.
13. Bartoszewicz, J.M.; Seidel, A.; Renard, B.Y. Interpretable detection of novel human viruses from genome sequencing data. *bioRxiv* **2020**. doi:10.1101/2020.01.29.925354.
14. Liang, Q.; Bible, P.W.; Liu, Y.; Zou, B.; Wei, L. DeepMicrobes: taxonomic classification for metagenomics with deep learning. *NAR Genomics and Bioinformatics* **2020**, *2*, lqaa009.
15. Shang, J.; Sun, Y. CHEER: hierarChical taxonomic classification for viral mEtagEnomic data via deep leaRning. *Methods* **2020**.
16. Tampuu, A.; Bzhalava, Z.; Dillner, J.; Vicente, R. ViraMiner: Deep learning on raw DNA sequences for identifying viral genomes in human samples. *PloS one* **2019**, *14*, e0222271.
17. Ren, J.; Song, K.; Deng, C.; Ahlgren, N.A.; Fuhrman, J.A.; Li, Y.; Xie, X.; Poplin, R.; Sun, F. Identifying viruses from metagenomic data using deep learning. *Quantitative Biology* **2020**, pp. 1–14.
18. Mock, F.; Viehweger, A.; Barth, E.; Marz, M. Viral host prediction with Deep Learning. *bioRxiv* **2019**, p. 575571.
19. Morales, J.A.; Saldaña, R.; Santana-Castolo, M.H.; Torres-Cerna, C.E.; Borrayo, E.; Mendizabal-Ruiz, A.P.; Vélez-Pérez, H.A.; Mendizabal-Ruiz, G. Deep Learning for the Classification of Genomic Signals. *Mathematical Problems in Engineering* **2020**, *2020*.
20. Nguyen, N.G.; Tran, V.A.; Ngo, D.L.; Phan, D.; Lumbanraja, F.R.; Faisal, M.R.; Abapihi, B.; Kubo, M.; Satou, K.; others. DNA sequence classification by convolutional neural network. *Journal of Biomedical Science and Engineering* **2016**, *9*, 280.
21. Guo, Y.; Li, W.; Wang, B.; Liu, H.; Zhou, D. DeepACLSTM: deep asymmetric convolutional long short-term memory neural models for protein secondary structure prediction. *BMC bioinformatics* **2019**, *20*, 1–12.
22. Zhu, H.; Guo, Q.; Li, M.; Wang, C.; Fang, Z.; Wang, P.; Tan, J.; Wu, S.; Xiao, Y. Host and infectivity prediction of Wuhan 2019 novel coronavirus using deep learning algorithm. *BioRxiv* **2020**.
23. Fang, Z.; Tan, J.; Wu, S.; Li, M.; Xu, C.; Xie, Z.; Zhu, H. PPR-Meta: a tool for identifying phages and plasmids from metagenomic fragments using deep learning. *GigaScience* **2019**, *8*, giz066.
24. Pian, C.; Li, Z.; Jiang, H.; Kong, L.; Chen, Y.; Zhang, L. Deep6mA: a deep learning framework for exploring similar patterns in DNA N6-methyladenine sites across different species. *bioRxiv* **2019**.
25. Kuang, S.; Wang, L. Identification and analysis of consensus RNA motifs binding to the genome regulator CTCF. *NAR Genomics and Bioinformatics* **2020**, *2*, lqaa031.
26. Zhang, Y.; Qiao, S.; Ji, S.; Li, Y. DeepSite: bidirectional LSTM and CNN models for predicting DNA–protein binding. *International Journal of Machine Learning and Cybernetics* **2019**, pp. 1–11.
27. Remita, M.A.; Halioui, A.; Daigle, B.; Kiani, G.; Diallo, A.B.; others. A machine learning approach for viral genome classification. *BMC bioinformatics* **2017**, *18*, 208.
28. Ren, J.; Song, K.; Deng, C.; Ahlgren, N.A.; Fuhrman, J.A.; Li, Y.; Xie, X.; Poplin, R.; Sun, F. Identifying viruses from metagenomic data using deep learning. *Quantitative Biology* **2020**, pp. 1–14.
29. Dey, L.; Chakraborty, S.; Mukhopadhyay, A. Machine learning techniques for sequence-based prediction of viral–host interactions between SARS-CoV-2 and human proteins. *Biomedical journal* **2020**.
30. Bzhalava, Z.; Tampuu, A.; Bała, P.; Vicente, R.; Dillner, J. Machine Learning for detection of viral sequences in human metagenomic datasets. *BMC bioinformatics* **2018**, *19*, 336.
31. Rizzo, R.; Fiannaca, A.; La Rosa, M.; Urso, A. A deep learning approach to dna sequence classification. *International Meeting on Computational Intelligence Methods for Bioinformatics and Biostatistics*. Springer, 2015, pp. 129–140.

32. Xu, J.; Xiang, L.; Liu, Q.; Gilmore, H.; Wu, J.; Tang, J.; Madabhushi, A. Stacked sparse autoencoder (SSAE) for nuclei detection on breast cancer histopathology images. *IEEE transactions on medical imaging* **2016**, *35*, 119–130.
33. Pratiher, S.; Chatteraj, S.; Vishwakarma, K. Application of stacked sparse autoencoder in automated detection of glaucoma in fundus images. *Unconventional Optical Imaging*. International Society for Optics and Photonics, 2018, Vol. 10677, p. 106772X.
34. Xiao, Y.; Wu, J.; Lin, Z.; Zhao, X. A semi-supervised deep learning method based on stacked sparse auto-encoder for cancer prediction using RNA-seq data. *Computer methods and programs in biomedicine* **2018**, *166*, 99–105.
35. Randhawa, G.S.; Soltysiak, M.P.; El Roz, H.; de Souza, C.P.; Hill, K.A.; Kari, L. Machine learning using intrinsic genomic signatures for rapid classification of novel pathogens: COVID-19 case study. *Plos one* **2020**, *15*, e0232391.
36. de M. Barbosa, R.; Fernandes, M.A. k-mers 1D and 2D representation dataset of SARS-CoV-2 nucleotide sequences. *Mendeley Data* **2020**, *v2*. doi:<http://dx.doi.org/10.17632/f5y9cggxny.2>.
37. Mapleson, D.; Garcia Accinelli, G.; Kettleborough, G.; Wright, J.; Clavijo, B.J. KAT: a K-mer analysis toolkit to quality control NGS datasets and genome assemblies. *Bioinformatics* **2016**, *33*, 574–576. doi:10.1093/bioinformatics/btw663.
38. Chor, B.; Horn, D.; Goldman, N.; Levy, Y.; Massingham, T. Genomic DNA k-mer spectra: models and modalities. *Genome biology* **2009**, *10*, R108.
39. Goodfellow, I.; Bengio, Y.; Courville, A. *Deep Learning*; MIT press, 2016.
40. The MathWorks. Matlab. <https://www.mathworks.com/>, 2020.