

# Estimated size of the total genome and protein space of viruses

Congyu Lu,<sup>1,2</sup> Yifan Wu,<sup>1</sup> Zheng Zhang,<sup>1,3</sup> Longfei Mao,<sup>1</sup> Xingyi Ge,<sup>1</sup> Aiping Wu,<sup>4,5</sup> Fengzhu Sun,<sup>6</sup> Yongqiang Jiang,<sup>7</sup> Yousong Peng<sup>1</sup>

**AUTHOR AFFILIATIONS** See affiliation list on p. 9.

**ABSTRACT** Recent metagenomic studies have identified a vast number of viruses. However, the systematic assessment of the true genetic diversity of the whole virus community on our planet remains to be investigated. Here, we explored the genome and protein space of viruses by simulating the process of virus discovery in viral metagenomic studies. Among multiple functions, the power function was found to best fit the increasing trends of virus diversity and was, therefore, used to predict the genetic space of viruses. The estimate suggests that there are at least  $8.23 \times 10^8$  viral operational taxonomic units and  $1.62 \times 10^9$  viral protein clusters on Earth when assuming the saturation of the virus genetic space, taking into account the balance of costs and the identification of novel viruses. It is noteworthy that less than 3% of the viral genetic diversity has been uncovered thus far, emphasizing the vastness of the unexplored viral landscape. To saturate the genetic space, a total of  $3.08 \times 10^8$  samples would be required. Analysis of viral genetic diversity by ecosystem yielded estimates consistent with those mentioned above. Furthermore, the estimate of the virus genetic space remained robust when accounting for the redundancy of sampling, sampling time, sequencing platform, and parameters used for protein clustering. This study provides a guide for future sequencing efforts in virus discovery and contributes to a better understanding of viral diversity in nature.

**IMPORTANCE** Viruses are the most abundant and diverse biological entities on Earth. In recent years, a large number of viruses have been discovered based on sequencing technology. However, it is not clear how many kinds of viruses exist on Earth. This study estimates that there are at least 823 million types of viruses and 1.62 billion types of viral proteins. Remarkably, less than 3% of this large diversity has been uncovered to date. These findings highlight the enormous potential for discovering new viruses and reveal a significant gap in our current understanding of the viral world. This study calls for increased attention and resources to be directed toward viral discovery and metagenomics and provides a guide for future sequencing efforts, enhancing our knowledge of viral diversity in nature for ecology, biology, and public health.

**KEYWORDS** virome, viral diversity, viral genetic space, ecosystem

Viruses are the most abundant and diverse biological entities on Earth (1). According to the Baltimore classification system, viruses can be classified into seven groups, including double-stranded DNA viruses (dsDNA), single-stranded DNA viruses, double-stranded RNA viruses, positive-sense single-stranded RNA viruses, negative-sense single-stranded RNA viruses, positive-sense ssRNA viruses with a DNA intermediate, and dsDNA viruses with ssRNA intermediates, based on their genetic materials and replication mode (2). Viruses exhibit wide variations in terms of host range, virion morphology, genome type, and size (3). For example, the size of a virus genome can span from several hundred to millions of bases (4, 5). Beyond their impact on human health and diseases, viruses also play a crucial role in maintaining the balance of global ecosystems.

**Editor** Michael J. Imperiale, University of Michigan, Ann Arbor, Michigan, USA

Address correspondence to Yousong Peng, [pys2013@hnu.edu.cn](mailto:pys2013@hnu.edu.cn), or Yongqiang Jiang, [jiangyq710327@sina.com](mailto:jiangyq710327@sina.com).

**Received** 12 August 2024

**Accepted** 16 December 2024

**Published** 25 February 2025

Copyright © 2025 Lu et al. This is an open-access article distributed under the terms of the [Creative Commons Attribution 4.0 International license](https://creativecommons.org/licenses/by/4.0/).

Previous studies have suggested that viruses are responsible for the daily demise of 20% of bacterial hosts in the oceans (1).

An enormous number of viruses have been identified, thanks to the rapid development of next-generation sequencing technology, in comparison to traditional virus identification methods based on virus isolation (6–9). For example, the IMG/VR database, recognized as the largest virus sequence database, has amassed over 15 million viral genomic sequences (6). Furthermore, novel viruses continue to be discovered at an unprecedented rate. The Global Ocean Viromes 2.0 project, for instance, has unveiled nearly 200,000 marine viral populations through the sequencing of 145 ocean samples (7). This raises a natural question regarding the genetic space of viruses. Previous studies by Rohwer estimated that there were 100 million phage species on Earth (10). The Global Virome Project put the estimate at 1.67 million viruses in birds and mammals (11). Koonin and colleagues estimated the total number of distinct virus species to be between  $10^7$  and  $10^9$  by combining available estimates of prokaryote species and virome size (12, 13). These studies based their estimates on the assumption that a host species harbors, on average, dozens of viruses, without considering the actual process of virus discovery through sampling and sequencing. This study delves into the genome and protein spaces of viruses by simulating the process of virus discovery in viral metagenomic studies.

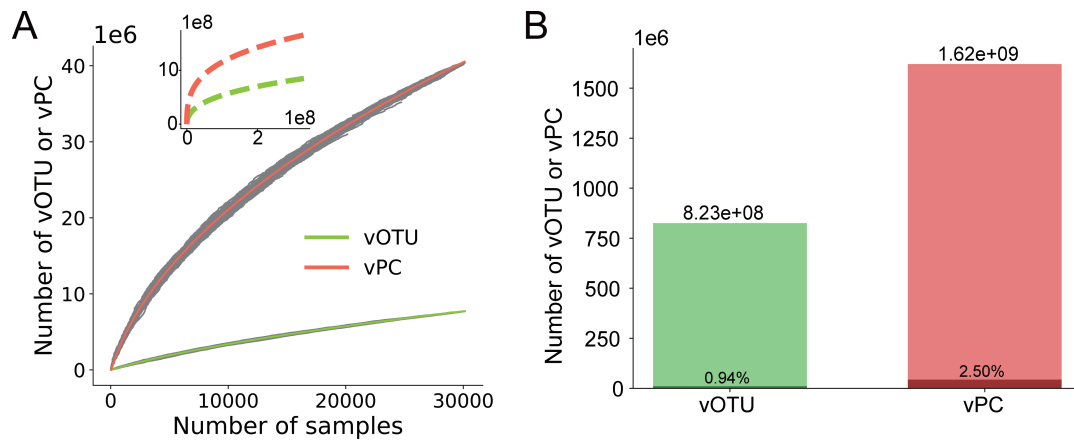
## RESULTS

### The increasing trends of the virus genome and protein space

First, we investigated the increasing trends in the virus genome and protein space as the number of samples increased. Specifically, we randomly selected 10 sequenced samples from the IMG/VR database without replacement and then counted the accumulated number of viral operational taxonomic units (vOTUs) and viral protein clusters (vPCs) (the definitions of vOTU and vPC are described in Materials and Methods) identified in the accumulated samples selected up to that point. This process was repeated until no samples remained (see Materials and Methods). The entire selection process was repeated 100 times. As depicted in Fig. 1A, both the number of vOTUs and vPCs increased rapidly. With all the samples in the IMG/VR database, a total of 7,721,789 vOTUs and 40,464,268 vPCs were identified.

To capture the increasing trends in the viral genetic space, we employed several common mathematical functions, including power, sigmoid, second-order polynomial, triple-order polynomial, logarithmic, exponential, and inverse proportional, to fit the relationship between the number of vOTUs or vPCs and the number of samples selected. To determine the best function, 3,000 samples were randomly selected as the training data, while the remaining 27,158 samples were used to test these functions. As illustrated in Fig. S1, only the power function fitted well for both vOTU and vPC, achieving a much lower mean squared error than other functions on the testing samples, even though all functions fit well on the training samples. Similar results were obtained when changing the data set size for training and testing, with the power function consistently fitting best among all functions (see Table S1). Therefore, the power function was used to predict the increasing trends of vOTU and vPC, assuming that more sequenced samples would be added in future studies (refer to the top-left sub-figure in Fig. 1A). According to the predictions, vOTUs and vPCs would continue to increase rapidly and then begin to expand with a slowing rate.

While the generation and disappearance of viruses occur daily due to their rapid evolution, there exists a balance between the birth and death of viruses, suggesting that the number of viruses on Earth should be finite. Consequently, it is anticipated that the identification of novel viruses will decrease as sampling continues. However, given the vast diversity and rapid evolution of viruses, it may be challenging to identify all viruses on Earth. To estimate the total number of viruses, the virus genetic space was assumed to be saturated when less than one novel virus was identified in a sample. Based on this assumption, the total numbers of vOTUs and vPCs were estimated to be  $8.23 \times 10^8$  and



**FIG 1** Estimating the virus genome and protein space. (A) Fitting the increasing trend of vOTUs and vPCs versus the number of selected samples. The gray points represent the number of vOTUs and vPCs identified in the selected samples in 100 simulations of sequencing studies (see Materials and Methods). The solid line depicts the fitted curve of the power function for vOTUs (green) and vPCs (red). The top-left sub-figure illustrates the prediction of the power function. (B) The estimated total number of vOTUs and vPCs when the virus genetic space was saturated, along with the proportion of explored viral genetic space at the levels of vOTUs and vPCs.

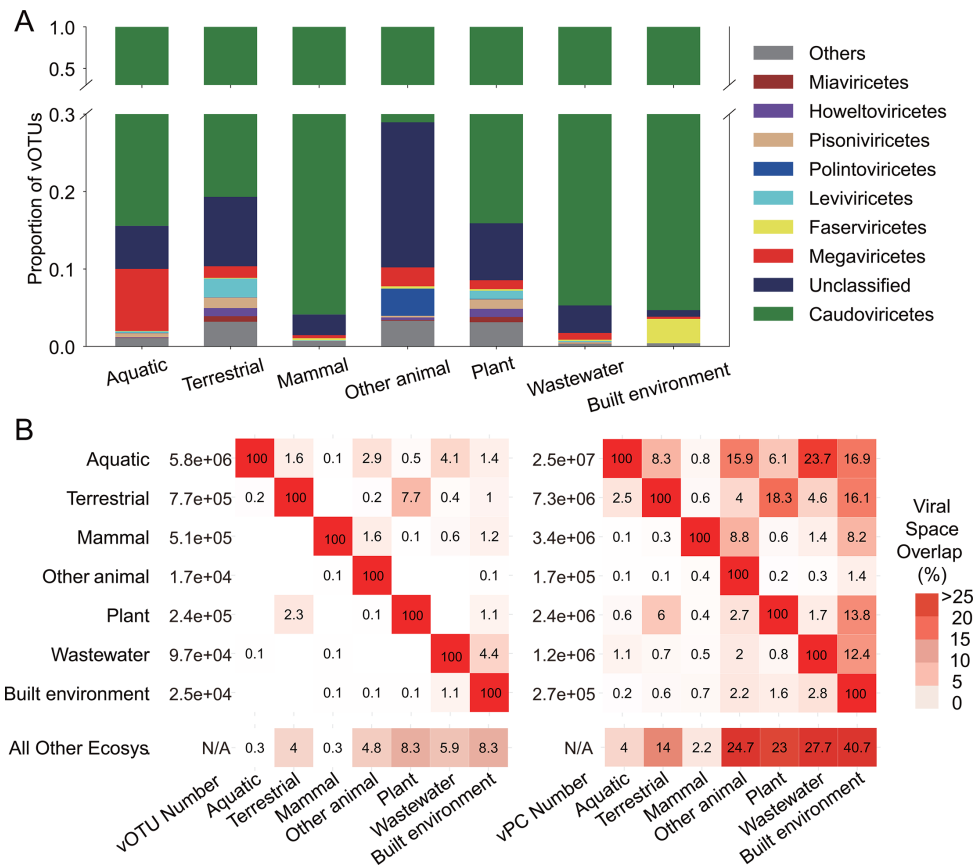
1.62e+09, respectively, when the virus genetic space was saturated (see Fig. 1B). The Proportion of Explored Viral Genetic Space (PEVGS), defined as the total number of vOTUs or vPCs identified in the IMG/VR database divided by the estimated numbers when the virus genetic space was saturated, was calculated to be 0.94% for vOTU and 2.50% for vPC. As depicted in Fig. 1A, the number of new viruses identified in new samples decreases as sampling continues. It was estimated that 3.08e+08 samples were required to saturate the genetic space. Unfortunately, only less than 40,000 samples had been sequenced so far according to the IMG/VR database.

### Taxonomy composition of viruses in different ecosystems

According to the Genomes OnLine Database ecological classification system (14), the samples used in this study primarily originated from seven ecosystems, including two environmental ecosystems (Aquatic and Terrestrial ecosystems accounting for 44% and 23% of all samples, respectively), three host-associated ecosystems (Mammal, Other-animal, and Plant ecosystems accounting for 13%, 2%, and 5% of all samples, respectively), and two engineered ecosystems (Wastewater and Built-environment ecosystems accounting for 1% and 4% of all samples, respectively) (refer to Table S2). The investigation into the taxonomic composition of viruses in different ecosystems involved calculating the proportion of vOTUs of different classes in each ecosystem (15). As depicted in Fig. 2A, the taxonomic composition of viruses across different ecosystems exhibited similarity, with the class of *Caudoviricetes* and unclassified viruses collectively constituting approximately 90% of all vOTUs. The class of *Caudoviricetes* held the largest proportion among all viral classes in every ecosystem, ranging from 71% to 96%. The taxonomic composition of the remaining vOTUs (approximately 10%) varied significantly among ecosystems. For instance, the class of *Faserviricetes* dominated the remaining vOTUs in the Built-environment ecosystem, where few other viruses were detected. Conversely, in the Other-animal ecosystem, the class of *Polintoviricetes* prevailed among the remaining vOTUs. Interestingly, the Aquatic ecosystem exhibited an enrichment of the class *Megaviricetes*, aligning with the predominance of their hosts (including algae and amoeba) in this ecosystem (15).

### Shared genetic diversity among different ecosystems

We then proceeded to analyze the extent of shared viral genetic diversity among different ecosystems. As illustrated in Fig. 2B, various ecosystems exhibited a very small



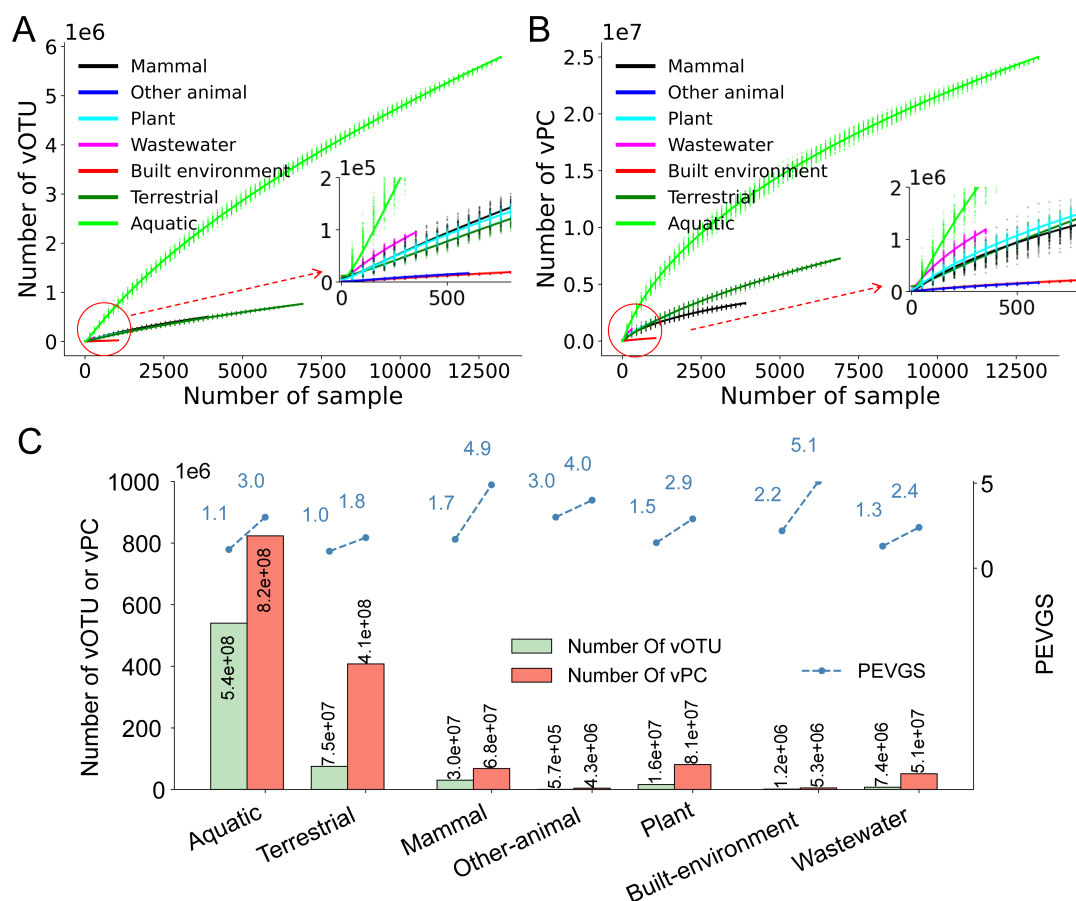
**FIG 2** Analysis of the viral genome and protein space by ecosystem. (A) The taxonomic composition of viruses identified in different ecosystems. (B) The proportion of vOTUs and vPCs shared between different ecosystems. The numbers presented in the heatmap indicate the percentage of vOTUs or vPCs shared between two ecosystems (indicated on the left and bottom) relative to the total vOTUs or vPCs in the ecosystem shown at the bottom. The blank squares in the heatmap represent values less than 0.1%. The number in the bottom row “All Other Ecosys.” refers to the proportion of shared vOTU or vPC between all other ecosystems and this particular ecosystem.

proportion of shared vOTUs. The proportion of shared vOTUs between any given ecosystem and other ecosystems ranged from 0.3% to 8.3%. Notably, the Aquatic ecosystem, although possessing the largest number of vOTUs, only covered a very small proportion of vOTUs in other ecosystems. For example, it covered 0.1% and 0.5% of vOTUs in Mammal and Plant ecosystems, respectively. Interestingly, the Plant ecosystem exhibited the highest degree of sharing with the Terrestrial ecosystem.

In contrast, the proportions of shared vPCs between different ecosystems were considerably larger than those of the shared vOTUs, ranging from 2.2% to 40.7%. The Mammal ecosystem had the smallest proportion of shared vPCs with other ecosystems (2.2%). Moreover, more than one-fifth of vPCs in three ecosystems, namely, Other-animal, Plant, and Wastewater, were shared with other ecosystems. Remarkably, for the Built-environment ecosystem, over 40% of the vPCs were shared with other ecosystems.

### The virus genome and protein space estimated by ecosystem

To analyze the increasing trend of virus genetic diversity by ecosystem as the number of samples increases, random sampling was conducted for each ecosystem, utilizing only samples from the corresponding ecosystem (refer Table S2 for the number of samples in different ecosystems). The increasing trends of vOTUs and vPCs in each ecosystem were analyzed and fitted using power functions. As depicted in Fig. 3A, the increasing trends of vOTUs in all ecosystems exhibited similarity, showing a slowing rate as the number



**FIG 3** Analysis of the viral genome and protein space by ecosystem. Panels A and B refer to the increasing trends of the numbers of vOTUs and vPCs, respectively, as the number of selected samples in different ecosystems. The simulation process was conducted based on the ecosystem. (C) The histogram displays the estimated total numbers of vOTUs and vPCs in each ecosystem when the virus genetic space was saturated. Additionally, the Proportion of Explored Viral Genetic Space of viruses in each ecosystem is illustrated at the levels of vOTU and vPC (indicated by blue dashed lines).

of selected samples increased. Notably, the number of vOTUs in the Aquatic ecosystem increased most rapidly with the growing number of selected samples, while the increase in other ecosystems was considerably slower. At the vPC level (Fig. 3B), the overall increasing trends mirrored those observed at the vOTU level. However, the differences in increasing trends of vPCs between the Aquatic and other ecosystems, especially the Aquatic and Terrestrial ecosystems, were smaller than those observed at the vOTU level.

We then estimated the total numbers of vOTUs and vPCs in different ecosystems when the viral genetic space was saturated (refer Fig. 3C). Predictions indicated that the Aquatic ecosystem was expected to have the largest number of vOTUs ( $5.4 \times 10^8$ ) and vPCs ( $8.2 \times 10^8$ ), followed by the Terrestrial ecosystem ( $7.5 \times 10^7$  vOTUs and  $4.1 \times 10^8$  vPCs) and the Mammal ecosystem ( $3.0 \times 10^7$  vOTUs and  $6.8 \times 10^7$  vPCs). In contrast, the four remaining ecosystems were predicted to contain considerably less viral diversity, with the Built-environment and Wastewater ecosystems projected to have only  $1.2 \times 10^6$  and  $7.4 \times 10^6$  vOTUs, respectively.

The PEVGS was calculated for each ecosystem at levels of vOTU and vPC (depicted in Fig. 3C). Across all ecosystems, PEVGS values increased from the level of vOTU (ranging from 1.0% to 3.0%) to vPC (ranging from 1.8% to 5.1%). Notably, host-associated ecosystems exhibited higher PEVGS values compared to environmental and engineered ecosystems. For instance, the Mammal ecosystem had PEVGS values of 1.7% and 4.9% at levels of vOTU and vPC, respectively, while those for the Terrestrial ecosystem were 1.0% and 1.8%, respectively.

**TABLE 1** Influence of different factors on the estimate of virus genetic space

Factor	Number of samples used	Estimated number of vOTUs
Redundancy of samples		
Unique IMG sample	30,158	8.23e+08
Unique BioSample	25,836	9.44e+08
Unique BioProject	20,260	9.40e+08
Sampling time		
Before 2018	15,349	6.15e+08
2018 and after	14,809	8.38e+08
Sequencing platform		
Illumina NovaSeq	7,426	5.38e+08
Illumina HiSeq	15,546	5.60e+08

### Stability analysis of the estimate of the virus genetic space

Finally, we delved into the influence of various factors on the estimate of the virus genetic space. The first factor considered was the redundancy of samples, where one sample might undergo sequencing multiple times or the sequencing data of the sample might be utilized in multiple studies. Upon eliminating redundancy using the NCBI BioSample, 25,836 biosamples were retained, and the estimated number of virus vOTUs was 9.44e+08. Similarly, when redundancy was removed using the NCBI BioProject (ensuring only one sample was used in each BioProject), the estimated number of virus vOTUs was 9.40e+08. Notably, both estimates, after eliminating redundancy, closely aligned with the original estimate (8.23e+08) (refer Table 1), suggesting that sample redundancy had a minor influence on the estimate of the virus genetic space.

The number of viruses identified in samples varied significantly over different periods, experiencing a substantial increase in 2018 (see Fig. S2). To investigate the impact of the year of virus identification on the estimate, all samples used in the study were categorized into two groups: one group comprised samples obtained before 2018, totaling 15,349 samples, while the other group included samples obtained in 2018 and later, totaling 14,809 samples. The virus genetic space was then estimated for each group of samples. As outlined in Table 1, the estimated numbers of viral vOTUs for the two groups were 6.15e+08 and 8.38e+08, respectively. Despite the former showing a nearly 30% decrease compared to the original estimate, both estimates were of the same order of magnitude as the original estimate.

The number of viruses identified from sequencing data generated by different sequencing platforms exhibited considerable variation (see Fig. S3). For instance, a median of 140 and 93 vOTUs were identified per sample sequenced by Illumina NovaSeq and HiSeq, respectively, while only 17 vOTUs were identified per sample sequenced by Illumina MiSeq. Consequently, we investigated the influence of the sequencing platform on the estimate of virus genetic space. As the majority of samples were sequenced by Illumina NovaSeq (7,426, 25%) and HiSeq (15,546, 52%), we estimated the virus genetic space based on samples sequenced exclusively by Illumina NovaSeq or HiSeq, resulting in estimated values of 5.38e+08 and 5.60e+08 vOTUs, respectively (refer Table 1). Once again, both estimates were of the same order of magnitude as the original estimate.

Virus protein clusters were obtained by clustering protein sequences with the threshold of identity and coverage set at 0.7. We then investigated the influence of

**TABLE 2** The influence of thresholds of identity and coverage on the estimate of vPCs

Threshold for identity and coverage	Number of vPCs obtained	Estimate of the number of vPCs
0.5	3.24e+07	1.10e+09
0.6	3.59e+07	1.30e+09
0.7	4.05e+07	1.62e+09
0.8	4.66e+07	2.01e+09



these thresholds on the estimate of viral genetic space. Increasing the threshold from 0.5 to 0.8 resulted in an increase in the number of vPCs from  $3.24 \times 10^7$  to  $4.66 \times 10^7$ . The estimated total number of vPCs also increased from  $1.10 \times 10^9$  to  $2.01 \times 10^9$ , remaining within the same order of magnitude as the original estimate when the threshold was set to 0.7 (refer Table 2).

## DISCUSSION

In recent years, numerous large-scale projects targeting the entire virome in various specific ecological environments have uncovered an abundance of novel viruses. Despite these efforts, a comprehensive understanding of the ultimate extent of virus genetic diversity on Earth has remained elusive. This study represents the first attempt to estimate the potential virus genome and protein space based on current virus genetic diversity, revealing a total of  $8.23 \times 10^8$  vOTUs and  $1.62 \times 10^9$  vPCs on Earth. This is in line with the numbers presented by Koonin and colleagues (11, 12). Notably, the estimated virus genetic space aligned with the sum of viral genetic diversity estimated individually within each ecosystem. Furthermore, the estimates of virus genetic space remained relatively stable when considering multiple factors, highlighting the robustness of the estimations regarding viral genome and protein space.

However, it is crucial to note that less than 3% of the viral genetic space has been unveiled thus far, emphasizing the necessity for additional virome projects to comprehensively capture viral genetic diversity. Unfortunately, as more viruses are identified, the number of novel viruses in new samples diminishes. This study estimates that a substantial  $3.08 \times 10^8$  samples would be required to encompass the entire viral genetic diversity. Notably, samples from the Aquatic ecosystem emerged as key contributors, capturing a significantly larger portion of viral genetic diversity than other ecosystems and should thus be prioritized in virus discovery efforts. Moreover, certain viral orders exhibited preferences for specific ecosystems, exemplified by the enrichment of Megaviricetes in the Aquatic ecosystem. Sequencing studies tailored to such preferences hold promise for the precise identification of specific kinds of viruses.

Despite utilizing the ICTV and Baltimore systems for viral taxonomy, most viruses analyzed in this study remained unclassified at the order and family levels. A critical need exists to develop an expandable classification framework for all viruses. The ideal classification framework should have the ability to classify all viruses, be easy to use, and be scalable. Previous studies by Jang and colleagues have developed the vConTACT for the classification of prokaryotic viruses based on genomic sequences (16). In the study, a gene-sharing network of viruses based on the shared protein clusters among viral genomes was built, and distance-based clustering and metric were integrated to provide the measures of clustering confidence. vConTACT is a scalable and automated classification tool as it can build classifications for more than 10,000 sequences, which were previously unclassified in oceanic viromics studies. Moreover, Jang's study provided a novel and important framework for organizing viral genetic diversity and demonstrated that the gene content-based methods are suitable for a unified classification of all viruses. Unfortunately, it was not widely used in virome studies. There was still a lack of a feasible classification framework for the whole virome on Earth.

Several limitations must be acknowledged. First, the estimated virus protein space may be influenced by parameter settings in protein clustering. Fortunately, the estimated total number of vPCs exhibited stability across varying parameters, ensuring the robustness of the estimate. Second, the study predominantly sampled from seven ecosystems, with potential bias in their representation. Some ecosystems, such as Other-animal and extreme environments, might be underrepresented, impacting the estimate of viral genetic diversity. Additionally, regional and country-based sampling bias, favoring developed countries, may further affect the accuracy of the estimate. Third, the RNA viruses were underrepresented as approximately 80% of the samples collected in the IMG/VR database were sequenced using the metagenomic strategy that favors DNA viruses. In combination, it is essential to view the estimate of virus genetic space

as a conservative lower bound. Finally, the estimate of the number of samples required to encompass the entire viral genetic diversity may be influenced by the sequencing method or protocols. Using high-throughput sequencers along with a virus-enriched method could potentially reduce the number of samples required.

## Conclusion

In summary, this study explores the virus genome and protein space, estimating the total number of vOTUs and vPCs on Earth when the virus genetic space is saturated. It serves as a guiding framework for future virus discovery sequencing efforts and significantly contributes to our understanding of viral diversity in nature.

## MATERIALS AND METHODS

### Data retrieval

The viral genome and protein sequences were retrieved from the IMG/VR database (version 4) (6). Meta information for metagenomic samples, from which viruses were identified, was obtained from the IMG/M database (17). Samples derived from experimental cultures were excluded as they do not adequately reflect natural virus diversity. Restricted samples were removed according to the JGI data utilization policy (6). A total of 30,158 samples were retained for further analysis, encompassing 13,746,160 viral genomic sequences, 7,721,789 vOTUs (defined as clusters of genomic sequences with an average of 95% or higher sequence identity on more than 85% of genomic regions according to the IMG/VR database [6]), and 179,131,383 protein sequences.

### Generation of protein clusters

To reduce computational costs in generating protein clusters, all viral protein sequences were initially clustered using the linclust algorithm of MMseqs2, with both coverage and identity set to 0.7 (18). This process yielded a total of 40,464,268 vPCs. The influence of parameter selection on the generation of vPCs is detailed in Table 2.

### Simulation of sampling in viral metagenomic studies

To simulate the virus discovery process in viral metagenomic studies, 10 sequenced samples were randomly selected from all samples without replacement. The number of vOTUs and vPCs identified in the accumulated samples selected up to that point was recorded. This sampling process was iterated until no samples remained. The increasing trend of the accumulated vOTUs and vPCs versus the number of accumulated samples selected was analyzed. The simulation process was repeated 100 times to mitigate randomness in sampling.

### Predicting the size of vOTUs and vPCs when the genetic space was saturated

Various mathematical functions, including sigmoid, power, exponential, second-order polynomial, triple-order polynomial, logarithmic, and inverse proportional functions, were employed to fit the increasing trend of the accumulated vOTUs versus the common logarithm of the accumulated samples selected in the simulation process. This was achieved using the `scipy.optimize.curve_fit` package in Python (formulas for these functions are listed in the supplemental material) (19).

### The calculation of vOTUs and vPCs shared between different ecosystems

If a vOTU includes sequences from samples of different ecosystems, it is considered to be shared between these ecosystems. The proportion of vOTUs shared between different



ecosystems was calculated by dividing the number of vOTUs that are shared between ecosystems by the total number of vOTUs in the ecosystem. In more detail, for two ecosystems A and B, if the number of shared vOTUs between A and B is N, the proportion of N to the total number of vOTUs in ecosystem A represents the proportion of vOTUs shared in A with B. Similarly, the proportion of N to the total number of vOTUs in ecosystem B represents the proportion of vOTUs shared in B with A. Similar methods were used for calculating the proportion of vPCs shared between different ecosystems.

## ACKNOWLEDGMENTS

We thank Prof. Simon Roux for guidance with the usage of the IMG/VR database and members in Peng's lab for helpful suggestions.

This work was supported by the National Natural Science Foundation of China (32170651 and 32370700), Hunan Provincial Natural Science Foundation of China (2024JJ2015), R&D Program of Guangzhou National Laboratory (GZNL2024A01002), and Scientific Research Program of the Educational Department of Hunan Province, China (22C0097).

All authors gave their consent for publication. Y.P. conceived and designed the analysis. Y.P., C.L., and Y.W. performed the analysis and prepared all figures. Y.P., C.L., Z.Z., Y.J., F.S., A.W., X.G., and L.M. wrote the main manuscript text. All authors reviewed the manuscript.

## AUTHOR AFFILIATIONS

<sup>1</sup>Hunan Provincial Key Laboratory of Medical Virology, Bioinformatics Center, College of Biology, Hunan University, Changsha, Hunan, China

<sup>2</sup>Hunan Provincial People's Hospital (The First Affiliated Hospital of Hunan Normal University), Hunan Normal University, Changsha, Hunan, China

<sup>3</sup>Hunan Engineering and Technology Research Center for Agricultural Big Data Analysis & Decision-making, College of Plant Protection, Hunan Agricultural University, Changsha, Hunan, China

<sup>4</sup>Institute of Systems Medicine, Peking Union Medical College, Chinese Academy of Medical Sciences, Beijing, China

<sup>5</sup>Suzhou Institute of Systems Medicine, Suzhou, China

<sup>6</sup>Department of Quantitative and Computational Biology, University of Southern California, Los Angeles, California, USA

<sup>7</sup>State Key Laboratory of Pathogen and Biosecurity, Institute of Microbiology and Epidemiology, Academy of Military Medical Sciences, Academy of Military Sciences, Beijing, China

## AUTHOR ORCIDs

Xingyi Ge  <http://orcid.org/0000-0003-3964-5140>

Aiping Wu  <http://orcid.org/0000-0002-5869-651X>

Yousong Peng  <http://orcid.org/0000-0002-5482-9506>

## DATA AVAILABILITY

All data used in the study are derived from the IMG/VR and IMG/M databases.

## ADDITIONAL FILES

The following material is available [online](#).

## Supplemental Material

**Supplemental material (mSphere00683-24-S0001.docx).** Mathematical functions used in the study, Fig. S1 to S3, and Tables S1 and S2.

## REFERENCES

1. Suttle CA. 2007. Marine viruses—major players in the global ecosystem. *Nat Rev Microbiol* 5:801–812. <https://doi.org/10.1038/nrmicro1750>
2. Lefkowitz EJ, Dempsey DM, Hendrickson RC, Orton RJ, Siddell SG, Smith DB. 2018. Virus taxonomy: the database of the International Committee on Taxonomy of Viruses (ICTV). *Nucleic Acids Res* 46:D708–D717. <https://doi.org/10.1093/nar/gkx932>
3. Masson P, Hulo C, De Castro E, Bitter H, Gruenbaum L, Essioux L, Bougueleret L, Xenarios I, Le Mercier P. 2013. ViralZone: recent updates to the virus knowledge resource. *Nucleic Acids Res* 41:D579–D583. <https://doi.org/10.1093/nar/gks1220>
4. Philippe N, Legendre M, Doutre G, Couté Y, Poirot O, Lescot M, Arslan D, Seltzer V, Bertaux L, Bruley C, Garin J, Claverie J-M, Abergel C. 2013. Pandoraviruses: amoeba viruses with genomes up to 2.5 Mb reaching that of parasitic eukaryotes. *Science* 341:281–286. <https://doi.org/10.1126/science.1239181>
5. Collins RF, Gellatly DL, Sehgal OP, Abouhaidar MG. 1998. Self-cleaving circular RNA associated with rice yellow mottle virus is the smallest viroid-like RNA. *Virology (Auckl)* 241:269–275. <https://doi.org/10.1006/viro.1997.8962>
6. Camargo AP, Nayfach S, Chen I-M, Palaniappan K, Ratner A, Chu K, Ritter SJ, Reddy TBK, Mukherjee S, Schulz F, Call L, Neches RY, Woyke T, Ivanova NN, Elie-Fadrosh EA, Kyrpides NC, Roux S. 2023. IMG/VR v4: an expanded database of uncultivated virus genomes within a framework of extensive functional, taxonomic, and ecological metadata. *Nucleic Acids Res* 51:D733–D743. <https://doi.org/10.1093/nar/gkac1037>
7. Gregory AC, Zayed AA, Conceição-Neto N, Temperton B, Bolduc B, Alberti A, Ardyna M, Arkhipova K, Carmichael M, Cruaud C, et al. 2019. Marine DNA viral macro- and microdiversity from pole to pole. *Cell* 177:1109–1123. <https://doi.org/10.1016/j.cell.2019.03.040>
8. Lu C, Peng Y. 2021. Computational viromics: applications of the computational biology in viromics studies. *Virol Sin* 36:1256–1260. <https://doi.org/10.1007/s12250-021-00395-7>
9. Li J-C, Zhao J, Li H, Fang L-Q, Liu W. 2022. Epidemiology, clinical characteristics, and treatment of severe fever with thrombocytopenia syndrome. *Infect Med* 1:40–49. <https://doi.org/10.1016/j.imj.2021.10.001>
10. Rohwer F. 2003. Global phage diversity. *Cell* 113:141. [https://doi.org/10.1016/S0092-8674\(03\)00276-9](https://doi.org/10.1016/S0092-8674(03)00276-9)
11. Carroll D, Daszak P, Wolfe ND, Gao GF, Morel CM, Morzaria S, Pablos-Méndez A, Tomori O, Mazet JAK. 2018. The Global Virome Project. *Science* 359:872–874. <https://doi.org/10.1126/science.aap7463>
12. Koonin E.V, Krupovic M, Dolja VV. 2023. The global virome: How much diversity and how many independent origins? *Environ Microbiol* 25:40–44. <https://doi.org/10.1111/1462-2920.16207>
13. Koonin Eugene V, Kuhn JH, Dolja VV, Krupovic M. 2024. Megataxonomy and global ecology of the virosphere. *ISME J* 18:wrad042. <https://doi.org/10.1093/ismejo/wrad042>
14. Mukherjee S, Stamatis D, Bertsch J, Ovchinnikova G, Sundaramurthi JC, Lee J, Kandimalla M, Chen I-M, Kyrpides NC, Reddy T. 2021. Genomes OnLine Database (GOLD) v. 8: overview and updates. *Nucleic Acids Res* 49:D723–D733. <https://doi.org/10.1093/nar/gkaa983>
15. Schoch CL, Ciufo S, Domrachev M, Hottel CL, Kannan S, Khovanskaya R, Leipe D, Mcveigh R, O'Neill K, Robertse B, Sharma S, Soussov V, Sullivan JP, Sun L, Turner S, Karsch-Mizrachi I. 2020. NCBI taxonomy: a comprehensive update on curation, resources and tools. *Database* 2020:baaa062. <https://doi.org/10.1093/database/baaa062>
16. Bin Jang H, Bolduc B, Zablocki O, Kuhn JH, Roux S, Adriaenssens EM, Brister JR, Kropinski AM, Krupovic M, Lavigne R, Turner D, Sullivan MB. 2019. Taxonomic assignment of uncultivated prokaryotic virus genomes is enabled by gene-sharing networks. *Nat Biotechnol* 37:632–639. <https://doi.org/10.1038/s41587-019-0100-8>
17. Chen I-MA, Chu K, Palaniappan K, Ratner A, Huang J, Huntemann M, Hajek P, Ritter S, Varghese N, Seshadri R, Roux S, Woyke T, Elie-Fadrosh EA, Ivanova NN, Kyrpides NC. 2021. The IMG/M data management and analysis system v.6.0: new tools and advanced capabilities. *Nucleic Acids Res* 49:D751–D763. <https://doi.org/10.1093/nar/gkaa939>
18. Steinegger M, Söding J. 2017. MMseqs2 enables sensitive protein sequence searching for the analysis of massive data sets. *Nat Biotechnol* 35:1026–1028. <https://doi.org/10.1038/nbt.3988>
19. Virtanen P, Gommers R, Oliphant TE, Haberland M, Reddy T, Cournapeau D, Burovski E, Peterson P, Weckesser W, Bright J, et al. 2020. SciPy 1.0: fundamental algorithms for scientific computing in Python. *Nat Methods* 17:261–272. <https://doi.org/10.1038/s41592-019-0686-2>