

1 **Protein Set Transformer: A protein-based genome language**
2 **model to power high diversity viromics**

3

4 Cody Martin^{1,2}, Anthony Gitter^{3,4,5}, and Karthik Anantharaman^{1,6,*}

5

6 ¹ Department of Bacteriology, University of Wisconsin-Madison, Madison, WI, USA

7 ² Microbiology Doctoral Training Program, University of Wisconsin-Madison, Madison,
8 WI, USA

9 ³ Department of Biostatistics and Medical Informatics, University of Wisconsin-Madison,
10 Madison, WI, USA

11 ⁴ Morgridge Institute for Research, Madison, WI, USA

12 ⁵ Department of Computer Sciences, University of Wisconsin-Madison, Madison, WI,
13 USA

14 ⁶ Department of Integrative Biology, University of Wisconsin-Madison, Madison, WI,
15 USA

16

17 * Correspondence: karthik@bact.wisc.edu (Karthik Anantharaman)

18

19 Keywords: viromics, viral metagenomics, genome language model, protein language
20 model, transformer

21 **Abstract**

22

23 Exponential increases in microbial and viral genomic data demand transformational
24 advances in scalable, generalizable frameworks for their interpretation. Standard
25 homology-based functional analyses are hindered by the rapid divergence of microbial
26 and especially viral genomes and proteins that significantly decreases the volume of
27 usable data. Here, we present Protein Set Transformer (PST), a protein-based genome
28 language model that models genomes as sets of proteins without considering sparsely
29 available functional labels. Trained on >100k viruses, PST outperformed other
30 homology- and language model-based approaches for relating viral genomes based on
31 shared protein content. Further, PST demonstrated protein structural and functional
32 awareness by clustering capsid-fold-containing proteins with known capsid proteins and
33 uniquely clustering late gene proteins within related viruses. Our data establish PST as
34 a valuable method for diverse viral genomics, ecology, and evolutionary applications.
35 We posit that the PST framework can be a foundation model for microbial genomics
36 when trained on suitable data.

37 Introduction

38 Viruses are the most abundant biological entity on the planet and inhabit every
39 ecosystem. Understanding how viruses modulate microbiome community dynamics and
40 functional outputs is an active area of research that spans various scales from global
41 biogeochemistry¹ to human health and disease². Despite the sheer abundance and
42 influence of viruses, comprehensive large-scale viral metagenomics (viromics) studies
43 are severely hindered by the enormous genetic diversity of viruses as most genomics
44 tools rely on sequence similarity to existing reference databases. These problems are
45 compounded by the lack of universal genes in viruses, complicating phylogenetic and
46 comparative analyses across diverse groups of viruses. Overall, these challenges have
47 impeded the development of viromics software that is both accurate and scalable to
48 increasingly diverse viral datasets. Thus, there is a clear need to develop data-driven
49 frameworks to study viruses at-scale using generalizable genomic principles instead of
50 simple sequence homology-based methods.

51
52 Protein language models (pLMs) are promising deep learning frameworks for
53 generalizable genomics. Trained on corpuses of millions of proteins³⁻⁵, pLMs have been
54 shown to model amino acids patterns in protein sequences akin to reading words in
55 sentences, capturing biochemical, functional, and structural features of proteins using
56 contextual information of the amino acids within a protein^{4,5}. Applications of pLMs to
57 viral datasets have demonstrated increased capacity for protein function annotation^{6,7}
58 and host prediction⁸. However, these studies only focused on specific tasks without
59 considering that pLMs could be universally beneficial for a variety of viromics tasks⁹,
60 thus missing out on the true potential of foundation pLMs. An additional shortcoming of
61 the pLMs themselves is that they do not account for evolution-driven genome
62 organization. Recent work has addressed this issue by contextualizing pLM
63 embeddings across short genomic fragments¹⁰ and even representing the entire
64 genome as an aggregation of the pLM embeddings⁸. However, each of these models
65 only targets one specific kind of representation: the former represents proteins with
66 added genome context, while the latter represents genomes as a weighted sum of
67 protein embeddings subject to a specific classification task. Thus, none of these
68 approaches are truly generalizable to a variety of viromics tasks that require both
69 protein- and genome-level reasoning.

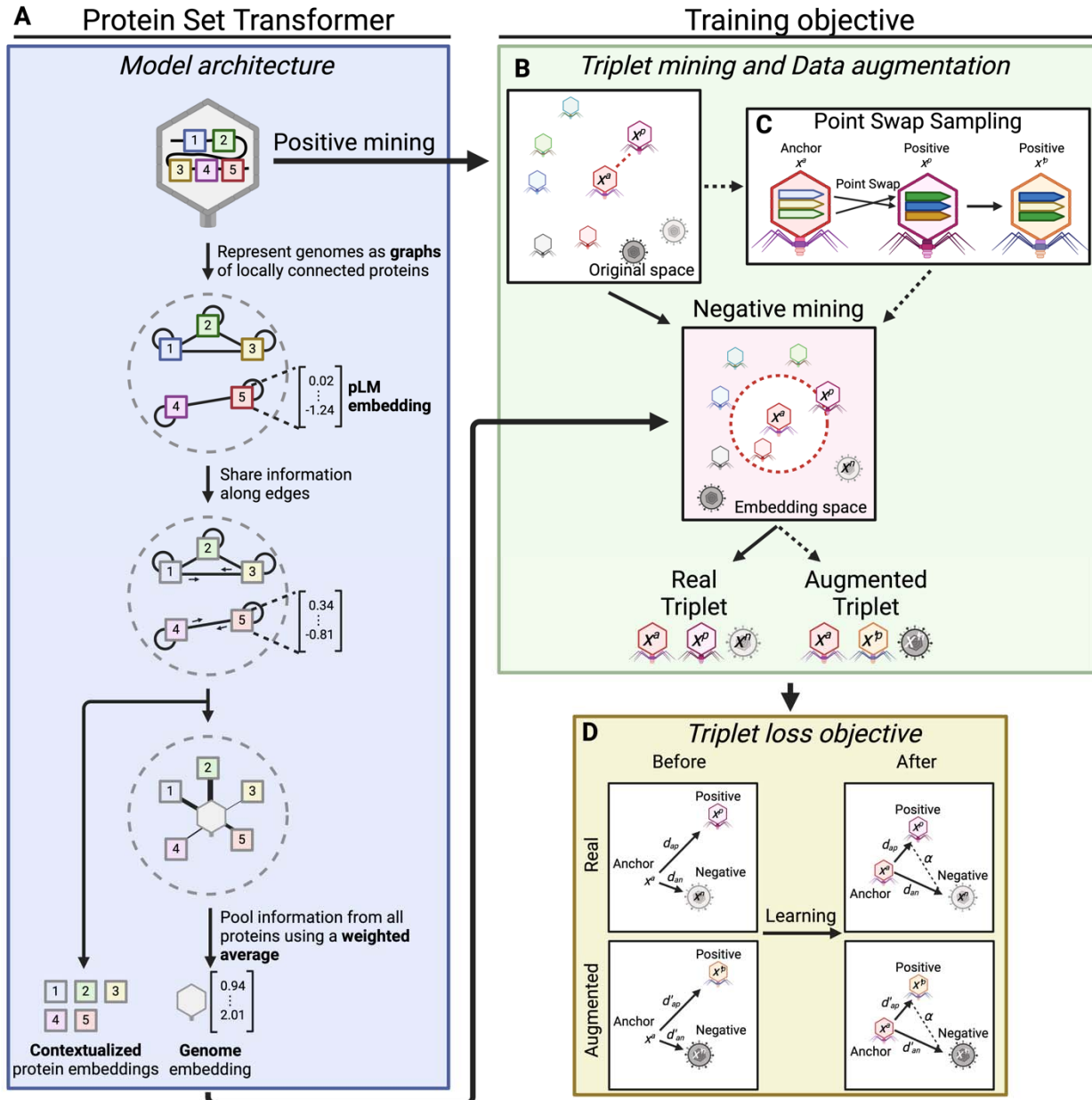
70
71 Here, we present our Protein Set Transformer (PST), a protein-based genome
72 language model that uses an encoder-decoder paradigm to simultaneously produce
73 both genome-contextualized protein embeddings and genome-level embeddings within
74 a single end-to-end model. We pretrained a foundation viral PST (vPST) model on
75 >100k high-quality dereplicated viral genomes encoding >6M proteins and evaluated on
76 a distinct test dataset of >150k high-quality viral genomes encoding >7M proteins from
77 IMG/VR v4¹¹. We demonstrate that vPST better relates viral genome-genome
78 relationships based on shared protein content. Further, we observe that only vPST can
79 consistently cluster operationally related proteins like late gene proteins, indicating the
80 importance of genome context-aware training. Additionally, vPST protein embeddings
81 are associated with protein structure relationships, as demonstrated by clustering

82 capsid-fold-containing proteins with no annotated function with annotated capsid
83 proteins.

84
85 Notably, neither the genome-contextualized vPST protein embeddings nor the
86 genome embeddings were learned with respect to any external labels, meaning that
87 they will be useful for a wide range of applications. Due to this flexibility of the vPST, we
88 propose that the vPST can be used for transfer learning to model other viral-centric
89 tasks such as viral gene and genome identification, genome quality control, genome
90 binning, taxonomy, and host prediction, which are major components of viromics
91 research⁹. Thus, we expect that the vPST will be foundational to future viromics studies.
92 Further, we posit that the PST architecture can be a general-purpose model for
93 microbial genomics when trained on microbial instead of or in addition to viral genomes.

94 95 **Results**

96
97 *Developing the Protein Set Transformer (PST) as a genome language model*
98



99
100
101
102
103
104
105
106
107
108
109
110
111
112
113

Figure 1. The Protein Set Transformer (PST) architecture and training regime. **A)** General overview of the graph-based PST for learning genome representations from contextualized protein embeddings. Each protein is represented by an ESM2 protein embedding. The PST internally represents each genome as a graph, consisting of multiple subgraphs of fully connected, locally adjacent proteins. The size of each subgraph is a tuned hyperparameter. The PST uses multi-head attention both to contextualize protein embeddings within each genome and to learn per-protein weights for a weighted averaged over each genome. See **Extended Data Fig. 1** for a modeling-centric view of the PST. Both protein and genome representations can be used for the appropriate downstream task. **B)** Triplet mining workflow that includes the data augmentation technique **C)** PointSwap sampling. For each training genome, a positive genome is identified from the ESM2 embedding space defined as the minimum Chamfer distance. Then, a negative, less related, genome is chosen from the PST embedding space that is the next farther genome after the positive. We augment our training data by creating hybrid genomes that swap similar protein vectors between each genome and its positive genome. **D)** Pictorial representation of the triplet loss objective function used to train the viral PST (vPST). The operational objective of triplet loss is to

114 embed each genome and its positive genome closer in embedding space than each genome and its
115 negative genome, within a tunable distance margin.

116
117 The PST (**Fig. 1A, Extended Data Fig. 1**) models genomes as sets of proteins
118 using principles from the natural language processing and set¹² and pointset¹³
119 transformer¹⁴ fields. We thus refer to the PST as a protein-based genome language
120 model, since it contextualizes protein information at genome-scale. In PST, all proteins
121 from each genome are embedded using the well-established ESM2 pLM^{3,4}. Unlike the
122 Set Transformer¹², the PST concatenates small vectors onto the pLM embeddings to
123 model both the protein genome position and coding strand. This set of updated protein
124 embeddings from each genome are fed to the PST encoder, which uses multi-head
125 attention¹⁴ to contextualize the protein representations within each genome (referred to
126 as simply “PST protein embeddings” from here out). These PST protein embeddings
127 can be used for protein-level tasks like protein classification and functional annotation.
128 In the end-to-end PST, the PST protein embeddings are further passed to the PST
129 decoder, which also uses multi-head attention to weigh the relative importance of each
130 protein in a genome. These weights are used for a weighted average of the
131 contextualized proteins to produce a genome representation.

132
133 A common training objective for language models that is used by a similar protein-
134 based genome language model¹⁰ is masked language modeling¹⁵, which involves
135 predicting masked tokens (words) in sentences from the rest of the sentence. In the
136 case of genome sentences composed of protein words represented as dense vectors,
137 masked language modeling is less intuitive and likely overcomplicates training. We
138 instead opt to mirror relationship-guided genomics to better understand patterns of
139 genetic diversity using the triplet loss function^{13,16} (**Fig. 1B, 1D**). During self-supervised
140 pretraining of the vPST foundation model, triplet loss uses the distance in vPST
141 embedding space as a measure of genome-genome relatedness. In vPST, genome-
142 genome relationships are implicitly conditioned on protein-protein relatedness. Briefly,
143 triplet loss involves the formation of genome triplets, consisting of one as an anchor, the
144 genome most related to the anchor as a positive example, and a genome less similar
145 than the positive genome as a negative example^{13,16} (**Fig. 1B, 1D**). Positive examples
146 are defined using the Chamfer distance in the input embedding space among genomes
147 within a training minibatch, while negative examples are sampled in the vPST
148 embedding space. Chamfer distance is computed as the average minimum of protein-
149 protein distances for pairs of genomes, meaning that the positive genome has the most
150 similar proteins to the anchor genome. The objective of triplet loss is to embed the
151 positive genome closer to the anchor than the negative within a tunable margin (**Fig.**
152 **1D**).

153
154 To help the vPST learn more generalizable representations, we used the data
155 augmentation technique PointSwap¹³ (**Fig. 1C**) for each genome and its most related
156 genome defined by Chamfer distance above (**Fig. 1B**). Each genome pair swaps
157 protein vectors that are most similar at a defined, tunable rate, analogous to
158 homologous recombination. We then update the triplet loss objective to include
159 maximizing the similarity between the anchor genome and its corresponding augmented
160 hybrid genome produced by PointSwap.

161
162 *Tuning the vPST using a modified Leave-One-Taxon-Out cross validation strategy*
163

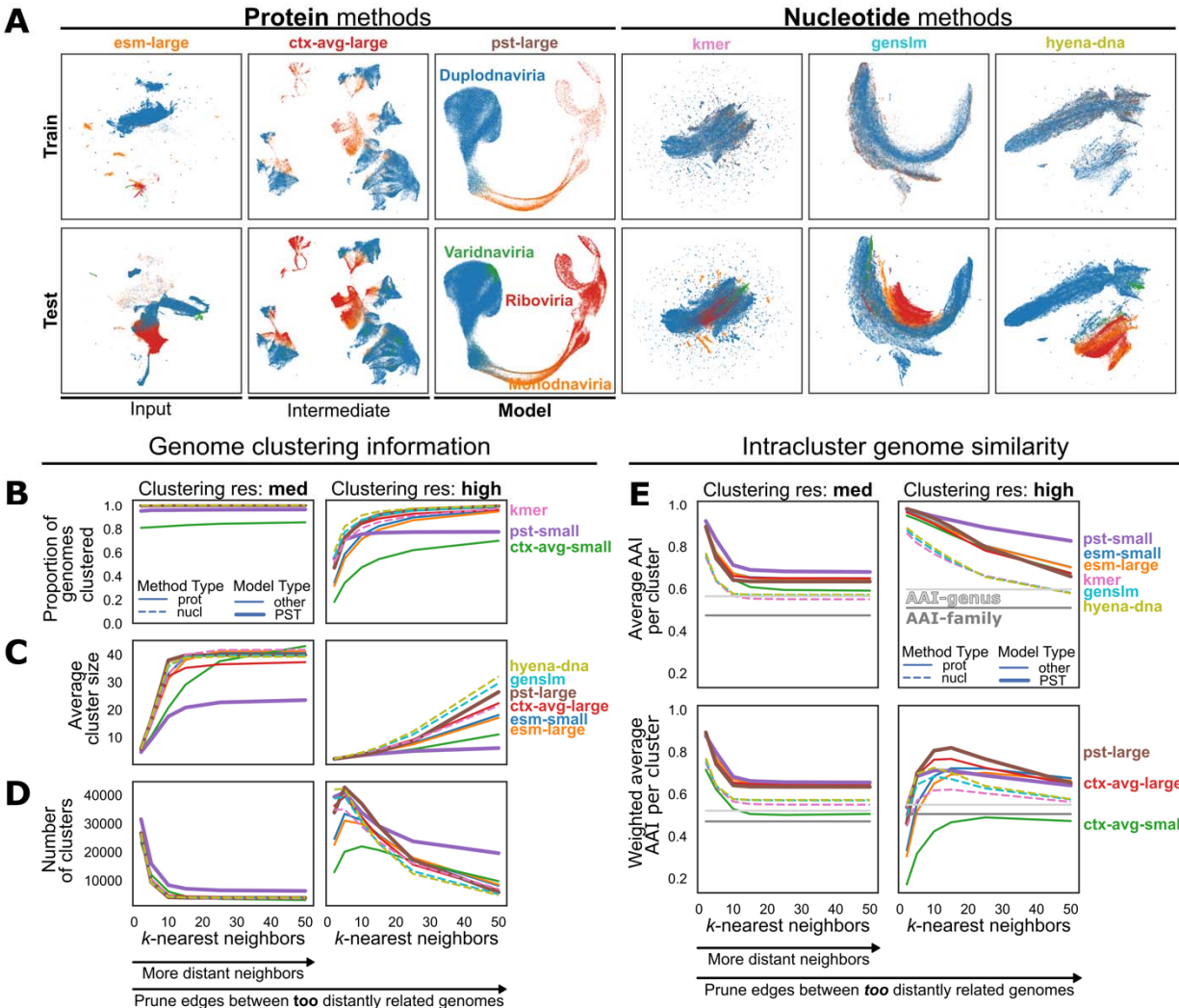
164 To train the vPST viral foundation model, we collected 103,589 viruses from 12
165 different publicly available sources^{1,17–27} as a training dataset (**Extended Data Fig. 2B**).
166 151,255 viruses from IMG/VR v4¹¹ that were distinct at the nucleotide level (<95%
167 average nucleotide identity over 85% of either genome) from the training viruses were
168 used as the test dataset. Dereplication at the nucleotide level was sufficient to reduce
169 train-test genome similarity at the protein level (**Extended Data Fig. 2A**). Most viruses
170 in either set were predicted to encode between 2–100 proteins (**Extended Data Fig.**
171 **2C**) and to be Duplodnaviria, Monodnaviria, or Riboviria (**Extended Data Fig. 2D**).
172 Additionally, most were from environmental sources not associated with a holobiont
173 (**Extended Data Fig. 2D**), especially marine systems. Among the viruses with a known
174 or predictable host, most are bacterial viruses (**Extended Data Fig. 2D**).
175

176 We tuned 2 different vPSTs with small (6-layer, 8M param) and large (30-layer, 150M
177 param) ESM2 protein embeddings, respectively, using a variant of leave-one-group-out
178 cross validation (CV), where the group is the viral taxonomic realm. In our variation, the
179 Duplodnaviria group is always included as a CV training fold since this group composes
180 a significant fraction of our training set (65.4%, **Extended Data Fig. 2D**). This CV setup
181 notably helps choosing model hyperparameters optimal for all viruses rather than just
182 the most abundant. The best models were chosen based on the lowest triplet loss
183 averaged among all folds at the end of tuning (**Extended Data. 3AB**). Using this
184 strategy, we tuned training-specific (dropout, layer dropout, learning rate, weight decay,
185 batch size), model-specific (number of attention heads and encoder layers), PST-
186 specific (chunk size), PointSwap-specific (rate), and triplet loss-specific (distance
187 margin, scale factor) hyperparameters (**Extended Data Fig. 3CD**). For the small vPST,
188 fewer attention heads and encoder layers led to optimal performance, while the reverse
189 is true for the large vPST, likely reflecting the increased information capacity of larger
190 pLM embeddings. Increasing values of the AdamW optimizer (PyTorch v2.0.0) weight
191 decay parameter, increasing values of the PointSwap rate, and decreasing values of the
192 triplet loss distance margin led to decreased loss (better performance) for both vPST
193 sizes. After hyperparameter tuning, we trained a final model for each ESM2 input using
194 the best hyperparameters (**Extended Data Fig. 3E**). The remaining results are based
195 on these 2 models that we refer to as “pst-small” (5M parameters) and “pst-large” (178M
196 parameters), respectively, when discussing both the learned genome and protein
197 representations.
198

199 *The vPST captures biologically relevant information about viral genomes*
200

201 To evaluate if the vPST learned biologically meaningful representations of viral
202 genomes, we compared the genome embeddings produced by the vPST against other
203 protein- and nucleotide-based methods with a quantitative clustering assessment on the
204 vPST test dataset. For protein-based methods, we performed an ablation study
205 comparing unweighted averages of the input ESM2 embeddings and of the vPST
206 protein embeddings over each genome (“ctx-avg” methods). For nucleotide-based

207 methods, we used 4-mer nucleotide frequency vectors, GenSLM²⁸ embeddings, and
 208 HyenaDNA²⁹ embeddings. The latter methods were chosen for both their availability
 209 relative to the course of this study and their relevance to genome language modeling,
 210 as there have been several recently described nucleotide-based models^{30,31}. Notably,
 211 both GenSLM and HyenaDNA have also been referred to as genome language models,
 212 so we explicitly refer to these as nucleotide language models to distinguish them from
 213 our protein-based genome language model. GenSLM was trained to focus on codon-
 214 level words in a genome sentence. Thus, to produce GenSLM genome embeddings, we
 215 embedded each nucleotide open reading frame (ORF) and then averaged these over
 216 each genome. Meanwhile, HyenaDNA is a long-context nucleotide language model that
 217 can contextualize up to 1Mb fragments, which is well above the size of most viral
 218 genomes. Protein-based methods and HyenaDNA appeared to better reflect the
 219 evolutionary relationships among viruses in both the vPST training and test datasets in
 220 a qualitative analysis of the genome embeddings in which there are visually distinct
 221 clusters of the 4 viral taxonomic realms (**Fig. 2A**).
 222



223 **Figure 2.** The vPST learns biologically meaningful genome representations for diverse sets of viruses. **A)**
 224 UMAP dimensionality reduction plots for the genome embeddings produced by each method, color coded
 225

226 by the viral realm. “Kmer” represents 4-mer nucleotide frequency vectors. “Ctx-avg” methods are
227 averages of the vPST protein embeddings over each genome. **B–D**) Statistics of genome clusters
228 detected by the Leiden algorithm on a k -nearest neighbor graph of the genome embeddings from the test
229 dataset (see Methods): **B**) proportion of genomes clustered, **C**) average number of genomes per cluster,
230 and **D**) total number of clusters. A cluster is only counted if there are at least 2 genomes. **E**) *Top*: Pairwise
231 amino acid identity (AAI) was computed for all pairs of viruses in a cluster and averaged for the entire
232 cluster. Then, the AAI for each cluster was averaged for each method, weighting the clusters by their size.
233 *Bottom*: The data in the top row were scaled by the proportion of genomes clustered from the test
234 dataset. All analyses were performed with the vPST test dataset.

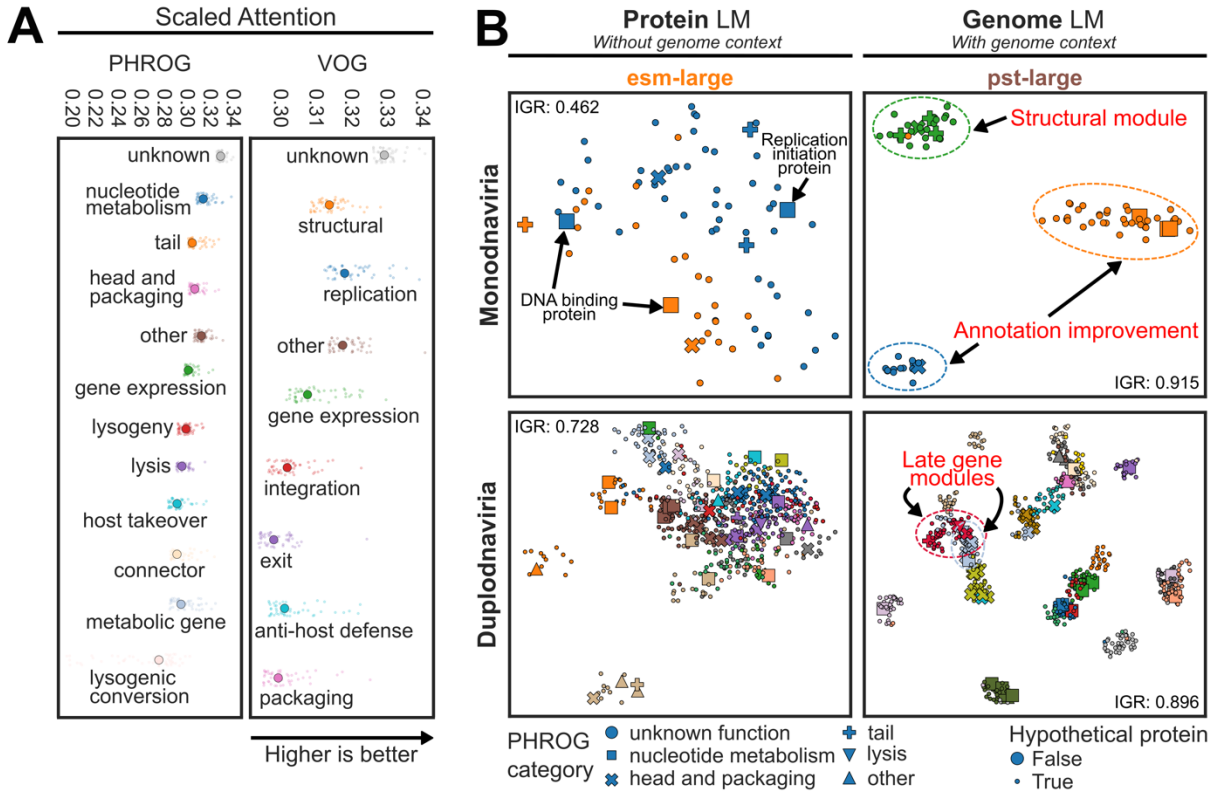
235
236 To quantitatively evaluate each genome representation, a similarity-weighted k -
237 nearest neighbor (k NN) graph was constructed from each of the genome embeddings
238 from the vPST test dataset and then clustered using the Leiden algorithm³². We
239 considered a range of values for k , the number of genome neighbors, and for the
240 clustering resolution, which sets a threshold for how distant connections can be, to
241 better understand the clustering trends with each genome representation (**Fig. 2B–D**).
242 As expected, increasing k from 2 to 50 leads to a greater proportion of viruses clustered
243 with at least 1 other genome (**Fig. 2B**), increases the average cluster size (**Fig. 2C**),
244 and decreases the total number of clusters (**Fig. 2D**). Likewise, increasing the clustering
245 resolution in the Leiden algorithm has the opposite effect when k is constant, since
246 more distant connections are pruned in the k NN graph (**Fig. 2B–D**, right column).

247
248 We then computed average amino acid identity (AAI) between pairs of genomes in
249 each genome cluster and aggregated the AAI over all genome clusters to assess the
250 quality of the genome clusters. As expected, protein-based methods lead to genome
251 clusters that have higher intra-cluster AAI than nucleotide-based methods (**Fig. 2E**),
252 suggesting that these methods use overall protein similarity to understand viral genome
253 relationships. Notably, pst-small genome clusters have the highest AAI among all
254 methods (**Fig. 2E**). However, when penalizing high rates of genome singletons, pst-
255 large genome clusters have the highest AAI (**Fig. 2E**). Importantly, this implies that pst-
256 large not only clusters viral genomes based on protein similarity but also can relate the
257 largest proportion of genomes. Additionally, most methods also outperform the baseline
258 of clustering viruses specifically using AAI at the genus or family level (**Fig. 2E**, “AAI-”
259 lines). Further, evaluating the taxonomic purity of both the viruses and their hosts across
260 the genome clusters does not strongly separate any method (**Extended Data Fig. 4**).
261 This may suggest that current viral taxonomy is not as informative for understanding
262 viral-viral relationships across diverse sets of viruses compared to AAI, which is based
263 on more intrinsic information to the viral genomes. Further, the proportion of viruses with
264 a predicted host is low (**Extended Data Fig. 2D**), which may also skew this analysis.

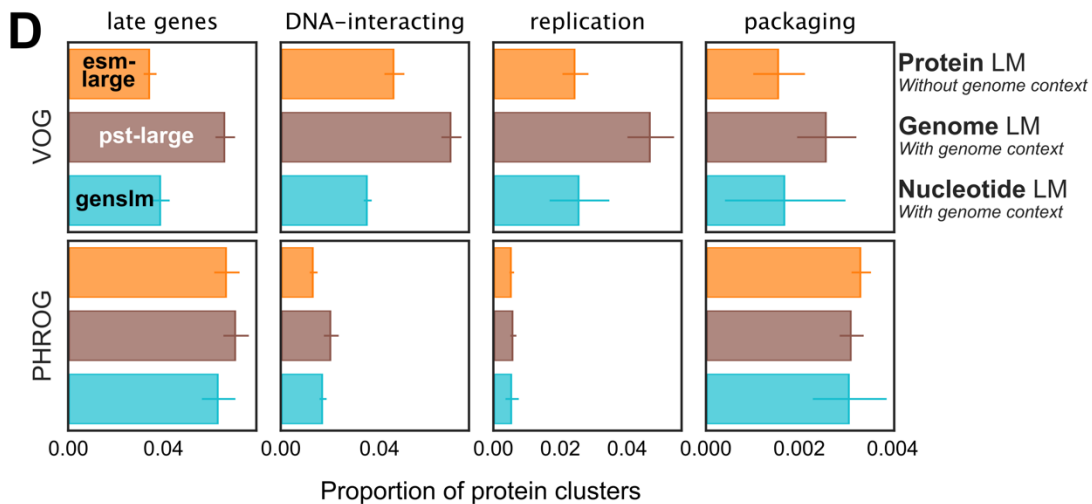
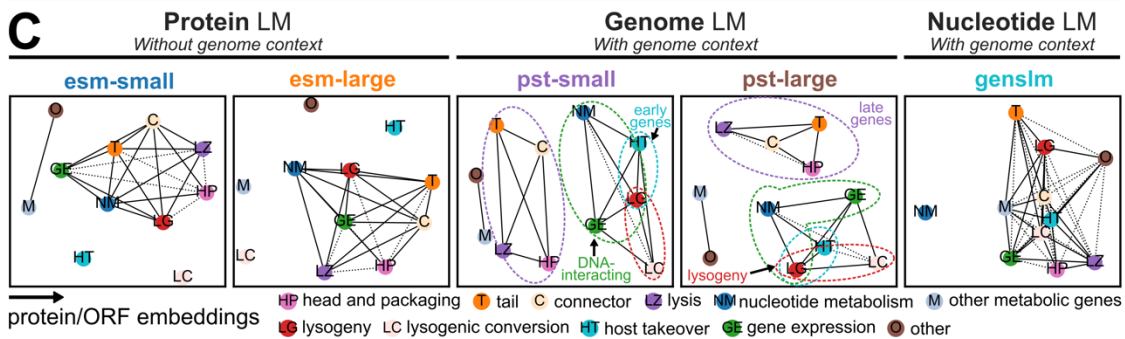
265
266 *The vPST detects important viral protein functions, including identifying new potential*
267 *hallmark proteins*

268
269 The vPST genome representations are produced as a function of the input protein
270 embeddings that get contextualized by the intermediate PST encoder. Thus, we
271 expected that the biologically meaningful genome embeddings of the vPST should be
272 generated from meaningful protein representations. We first analyzed the attention
273 scores of each protein per genome from pst-large, which are used as importance scores

274 when pooling the vPST protein embeddings for the final genome representation. We
275 considered that the general function of each protein was likely associated with high
276 attention. Indeed, structural proteins (head, packaging, tail) and replication or nucleotide
277 metabolism proteins were generally most attended to by the model (**Fig. 3A**). This is
278 intuitive since these proteins are essential to viruses and likely reflects their relatively
279 greater abundance in the dataset (**Extended Data Fig. 5B**). Further, we found a subtle
280 association between the attention scores with the number of proteins belonging to the
281 same sequence identity-based cluster (**Extended Data Fig. 5A**). This reflects the model
282 assigning a higher weight to proteins seen more frequently.



Functional co-occurrence



284 **Figure 3.** The vPST leverages genomic context to learn protein function relationships. **A)** Scaled attention
285 from pst-large normalized to compare across genomes with differing numbers of proteins (see Methods)
286 with respect to protein function. Scaled attention is the max scaled attention of all proteins in each of the
287 top 50 sequence identity-based protein clusters (mmseqs2). **B)** UMAP dimensionality reduction plots for 2
288 genome clusters that were primarily ($\geq 85\%$ of genomes) composed of Monodnaviria (top, 13 genomes, 80
289 proteins) or Duplodnaviria (bottom, 4 genomes, 682 proteins). Colors indicate protein cluster membership
290 defined by clustering the k -nearest neighbor graph of the indicated protein embedding with the Leiden
291 algorithm. Here, pst-large refers to the vPST protein embeddings. “IGR” refers to the average weighted
292 information gain ratio for all protein clusters within each of the two genome clusters as a measure of
293 protein cluster functional purity (see Methods). Shapes indicate the PHROG functional category. **C)**
294 Summary of functional co-clustering based on PHROG annotations. Each connected component was
295 clustered in a co-occurrence graph using the Leiden algorithm with resolution of 1.0. The edges indicate
296 pairs of functional categories that were more enriched in protein clusters defined by clustering the k -
297 nearest neighbor graph of the corresponding protein/ORF embeddings (columns) relative to the
298 background distribution of annotation profiles. The length of the edges reflects the degree of enrichment
299 since the networks were visualized using a spring force algorithm. Dotted lines indicate connections that
300 were less enriched than or equal to expected, while solid lines were more enriched than expected. **D)** The
301 proportion of protein clusters that correspond to one of the indicated function modules (columns) when
302 using either the VOG (top) or PHROG (bottom) annotation databases. For B and C, genomes were
303 clustered using pst-large genome embeddings ($k=15$, clustering resolution=“high”). Proteins were
304 clustered within each genome cluster with $k=15$ and clustering resolution=“med”. All analyses were
305 generated using the vPST test dataset.

306
307 To quantitatively assess the ability of the vPST to understand protein relationships,
308 we conducted a similar analysis as with the genome clusters. The embedding-based
309 protein clusters were generated using the Leiden algorithm on a similarity-weighted k NN
310 graph. To reduce potential noise when clustering, we restricted the set of nearest
311 neighbors to only include proteins from genomes in the same genome cluster,
312 specifically using the hyperparameters that maximized intra-genome-cluster AAI ($k=15$,
313 resolution=“high”, **Fig. 2E**). We performed a similar purity analysis of the protein clusters
314 with respect to VOG and PHROG functional categories that did not strongly indicate
315 which protein or genome clustering methods produced the most functionally pure
316 genome clusters (**Extended Data Fig. 6**). However, clustering the genomes with the
317 ctx-avg-large embeddings tended to perform best for protein cluster functional purity
318 (**Extended Data Fig. 6B**). This result makes sense because the vPST protein
319 embeddings used for the ctx-avg-large genome embeddings are the last time the vPST
320 directly considers protein information. Additionally, vPST protein embeddings led to the
321 overall highest protein functional purity.

322
323 To identify cases where the vPST outperforms the input ESM2, we visualized the
324 protein embeddings from 2 representative genome clusters primarily ($\geq 85\%$ of
325 genomes) composed of Monodnaviria or Duplodnaviria using the large embeddings
326 (**Fig. 3B**). In the Monodnaviria cluster, there are DNA binding proteins that esm-large
327 did not cluster together, reflecting the underlying sequence divergence of these two
328 proteins (35.5% sequence identity over $\sim 71\%$ coverage). However, pst-large clustered
329 these proteins with a replication initiation protein, suggesting a detection of broad
330 functional relationships. Furthermore, the esm-large embeddings clustered various
331 structural proteins together with these DNA-interacting proteins that pst-large notably
332 clustered into distinct clusters. There are additionally numerous proteins unable to be
333 annotated by PHROG (**Fig. 3B**) or VOG (**Extended Data Fig. 5C**) that cluster with

334 annotatable proteins regardless of protein embedding used. Similar visual patterns exist
335 for the Duplodnaviria cluster, which prompted us to consider if these were general
336 phenomena of the vPST protein clustering.

337

338 *The vPST co-clusters related protein functions into function modules*

339

340 Given that the vPST leverages genomic context, we suspected that the vPST would
341 be equipped to identify groups of associated protein functions that reflect the underlying
342 genome organization. For example, the late genes encoding for structural, packaging,
343 and lysis proteins are adjacent and transcribed by a single promoter in the Lambda
344 genome³³. We, therefore, assessed protein function co-clustering patterns. For each
345 protein cluster, we calculated the number of times pairs of proteins belonging to different
346 PHROG functional categories co-clustered against the number of times each pair of
347 categories would be expected to co-cluster based on the underlying distribution of the
348 PHROG database categories. The resulting enrichment networks showed that both
349 vPST models could group proteins based on broader function modules (**Fig. 3C**),
350 regardless of the genome embedding used for genome clustering (**Extended Data Fig.**
351 **7**). For example, tail, head and packaging, connector, and lysis proteins, which are
352 notably late gene proteins, consistently co-clustered above background in vPST protein
353 clusters. Additionally, DNA-interacting (nucleotide metabolism, lysogeny, and gene
354 expression), early gene (host takeover, lysogeny), and lysogeny (lysogeny, lysogenic
355 conversion) function modules were enriched in vPST protein clusters. Interestingly,
356 regardless of how the genomes were clustered, using ESM2 protein embeddings to
357 cluster the proteins did not lead to interpretable functional modules emerging (**Fig. 3C**,
358 **Extended Data Fig. 7**). Additionally, while some functional relationships were detected
359 in the GenSLM ORF clusters, this was not consistent depending on how the genomes
360 were clustered. These results were also consistent with the proportion of protein
361 clusters that we considered as belonging to these function modules such as late genes,
362 DNA-interacting, replication, and packaging (**Fig. 3D**). Notably, protein clusters that
363 belonged to these function modules made up a greater proportion of vPST protein
364 clusters than ESM2 or GenSLM clusters when using VOG annotations, regardless of
365 how the genomes were clustered (**Extended Data Fig. 8A**). The effect is less
366 pronounced with PHROG annotations (**Extended Data Fig. 8B**), but this difference may
367 be attributed to the overall decrease in functional annotation with the PHROG database
368 (**Extended Data Fig. 5B**), which led to excluding a greater number of protein clusters
369 that belonged to each functional module. These data demonstrate that considering
370 genome context better enables the vPST to detect broader functional associations
371 implicitly encoded in viral genome organization.

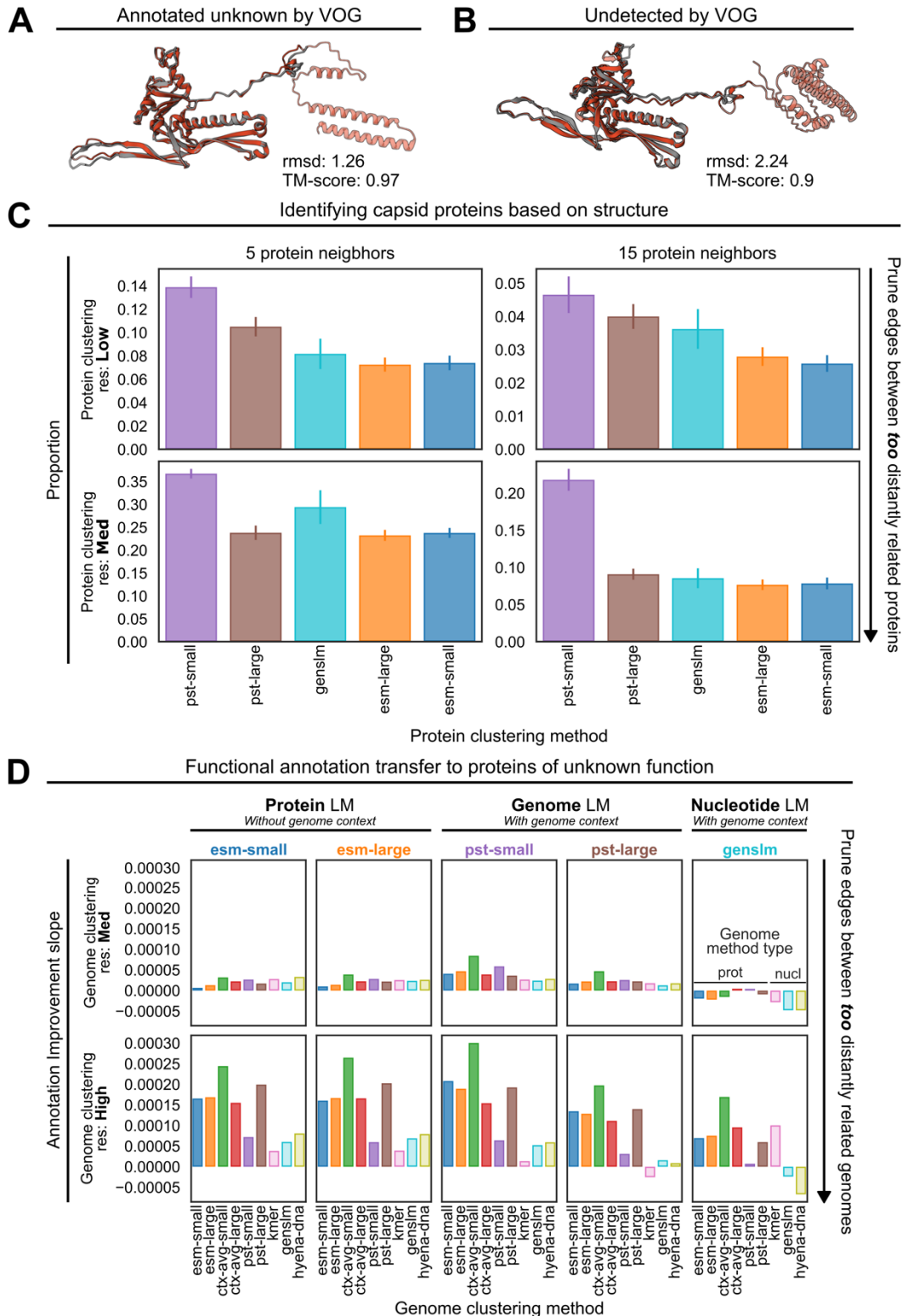
372

373 *The vPST expands our understanding of proteins of unknown function*

374

375 Interestingly, hypothetical proteins unable to be annotated by either the VOG or
376 PHROG databases were considered the most important by pst-large (**Fig. 3A**). One
377 explanation is that since proteins of unknown function make up 70-90% of all proteins in
378 the vPST test dataset (**Extended Data Fig. 5B, Supplementary Table 2**), it is likely that
379 there are true viral hallmark structural and replication proteins that have diverged at the

380 sequence level among the unannotated proteins. To understand if the vPST uses more
381 than sequence-level information to relate proteins, we investigated whether
382 unannotated proteins that cluster with detectable capsid proteins contained conserved
383 capsid-like structural folds as evidence that these unannotated proteins are indeed
384 capsid proteins. We filtered the proteins from the test viruses to maintain proteins
385 belonging to protein clusters that contained only annotated capsid proteins or
386 hypothetical proteins. We then used foldseek³⁴ and ProstT5³⁵ to translate this protein
387 set into a structural alphabet for searching against Protein Data Bank³⁶ structures for
388 structural homology. To validate the structural reasoning of this approach that does not
389 directly infer a protein structure, we independently aligned the structures of the
390 reference HK97 major capsid protein³⁷ with two different AlphaFold 3-predicted³⁸
391 structures using the most structurally similar proteins from our dataset: one that was
392 detected by a VOG profile with unknown function (**Fig. 4A**, pTM=0.66) and one
393 undetected entirely (**Fig. 4B**, pTM=0.6). The strong alignments indicate that our
394 workflow can accurately identify capsid-fold-containing proteins from the protein
395 sequence alone. Using this approach, the vPST models generally showed the greatest
396 average proportion of unannotated proteins with structural homology to known capsid
397 proteins (**Fig. 4C**), regardless of how the proteins or genomes were clustered
398 (**Extended Data Fig. 9A**). GenSLM ORF embeddings were also better than the ESM2
399 protein embeddings for this task, likely due to being pretrained on microbial genomes,
400 which would contain some viral sequences, and finetuned on SARS-CoV-2 genomes.
401



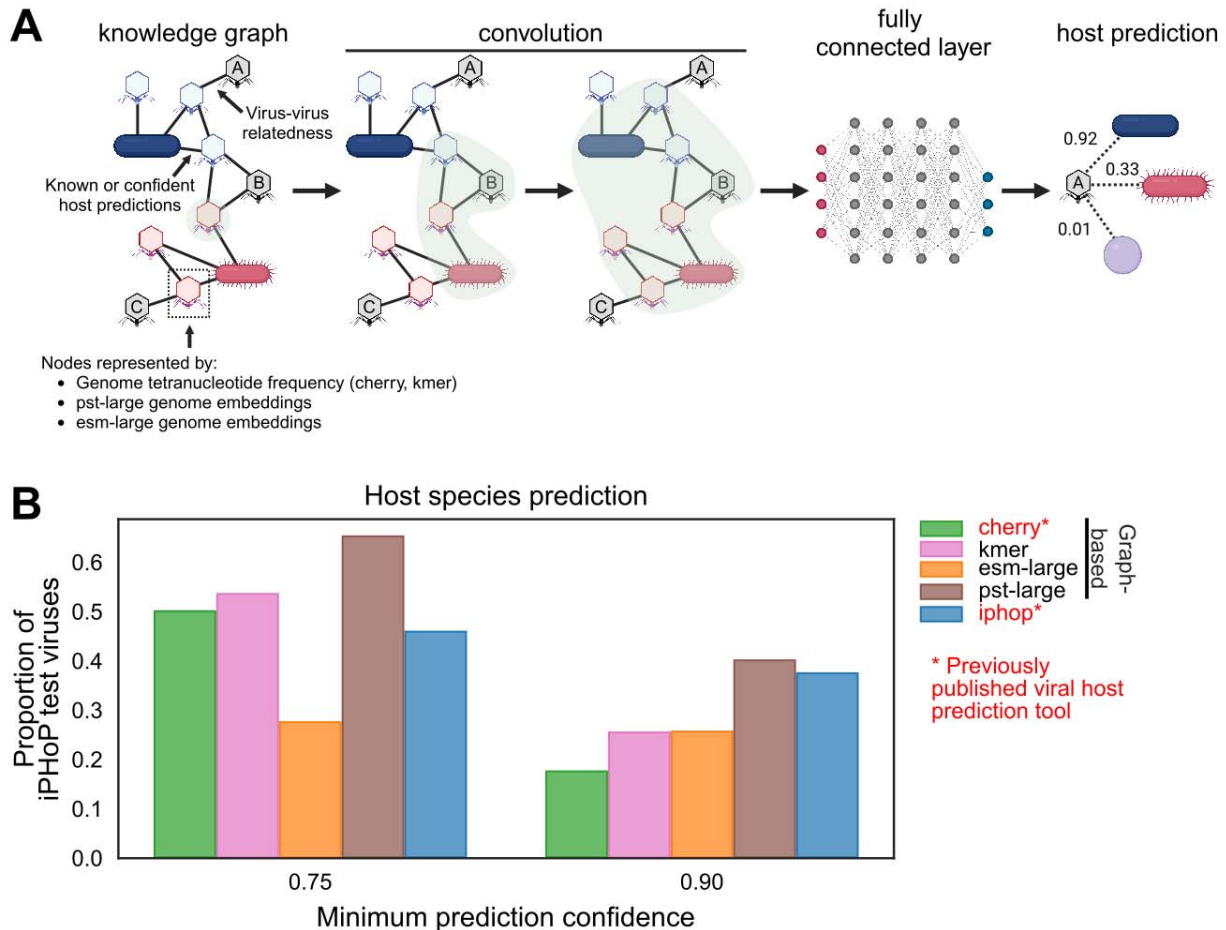
402
403
404
405
406
407

Figure 4. The vPST expands functional annotation of hypothetical proteins. **AB)** Structural alignments with the HK97 major capsid protein (PDB: 2FS3, gray) for a protein annotated by VOG as unknown (**A**, “IMGVR_UViG_2851668853_000002|2851668853|2851668853|1181413-1220308_35”) and another undetected by VOG (**B**, “IMGVR_UViG_3300036770_002539|3300036770|Ga0310126_0001736_19”). The red cartoon diagrams are the query proteins from our dataset and were chosen due to being the

408 most similar to the HK97 capsid protein from each category. **C**) The average proportion of proteins
409 unannotated by VOG clustering with annotated capsid proteins that have structural homology with known
410 capsid folds. Structural homology was detected using foldseek searching against the Protein Data Bank
411 database. Error bars represent the standard deviation over the embedding used for genome clustering.
412 Values are only comparable within each subpanel. **D**) Sensitivity of annotation transfer from annotated to
413 nearby unannotated proteins over the choice of k nearest neighbors for genome clustering. Instances of
414 annotation transfer were detected if the nearest protein (based on cosine distance of the protein/ORF
415 embeddings) to each unannotated protein had a VOG annotation. All analyses performed with genomes
416 and proteins from the test dataset.

417
418 We next considered that similarity in embedding space could be used to propagate
419 functional labels from annotated to unannotated proteins. To evaluate the annotation
420 transfer ability of vPST, we first performed a nearest neighbor sensitivity analysis. For all
421 unannotated proteins in the test set, we identified the nearest protein within each
422 genome cluster using cosine distance of each protein embedding. If the nearest protein
423 was annotated by VOG, we scored that as an improvement in annotation. Protein-based
424 embeddings outperformed the GenSLM ORF embeddings for transferring annotations,
425 regardless of how the genomes were clustered (**Extended Data Fig. 9B**). Additionally,
426 the rate at which this annotation improvement increases as more genome neighbors are
427 considered showed that the vPST was more sensitive (**Fig. 4D**). Specifically, clustering
428 genomes with ctx-avg-small or pst-large genome embeddings led to the greatest rate of
429 improvement as more genome neighbors are allowed. Interestingly, when nucleotide
430 methods were used for the genome clustering or the protein distance searches, the rate
431 decreased, suggesting that adding more genome neighbors impedes annotation
432 transfer. This may be due to the limited range of nucleotide information in capturing
433 distant relationships. This means that as the nucleotide-based genome clusters
434 increase in size, the nearest neighbor in ORF embedding space for an unannotated
435 protein is just another unannotated protein. Further, considering only the single nearest
436 protein is a conservative baseline. It would be possible to improve these results not only
437 by considering more protein neighbors but also by finetuning the vPST with a protein
438 annotation task.

439
440 *The vPST can be applied toward viral host prediction*
441



442
443
444
445
446
447
448

Figure 5. The vPST improves host prediction. **A)** Graph neural network approach for host prediction developed by CHERRY. The node representations are swapped out to the corresponding data type. **B)** Proportion of iPHoP test viruses whose true host species is predicted above the indicated confidence threshold. No test viruses were filtered for similarity to those in the vPST training set. The graph-based models were trained in this study, while “iphop” represents the results of iPHoP on the test set.

449
450
451
452
453
454
455
456
457
458

Since we expect that the vPST can be used as a general-purpose model for downstream viromics tasks, we used the pst-large genome embeddings for viral host prediction as a proof-of-concept (**Extended Data Fig. 10A**). We adopted and modified a graph framework described previously³⁹ that models this scenario as a link prediction task in a virus-host interaction network. Briefly, the objective is to predict for any pair of virus and host whether there should be a link, indicating a prediction for infection of that host by the corresponding virus (**Fig. 5A**). This task can be performed by a Graph Neural Network (GNN), which uses a form of convolutions to aggregate local (more related) parts of the graph to improve link prediction.

459
460
461
462
463

We implemented a variant of the GNN-based CHERRY algorithm³⁹ (**Fig. 5A**), swapping out the node (genome) embeddings of both viruses and hosts with either ESM2, vPST, or the tetranucleotide frequency (kmer) vectors that CHERRY uses. Although this design is likely suboptimal for vPST, which has embeddings specialized for viruses but not hosts, it enables a direct comparison of the choice of genome

464 embedding instead of various virus-host genome embedding combinations. We then
465 trained these models using the training dataset of the host prediction tool iPHoP⁴⁰ to
466 compare with previously published work (**Extended Data Fig. 10A**). Then, each trained
467 model and iPHoP were evaluated using the same iPHoP test dataset. We evaluated
468 whether the true host species for each test virus could be identified with high confidence
469 (**Fig. 5B**). The model using vPST genome embeddings outperformed all other methods
470 at the host species level, although the margin between vPST and iPHoP was small
471 when retaining predictions ≥ 0.9 confidence. Although there are viruses in the iPHoP
472 test set that are similar to those in the vPST training set (**Extended Data Fig. 10B**),
473 excluding these viruses does not change the overall results (**Extended Data Fig. 10D**).
474 Additionally, when evaluated at broader host taxonomic levels, the kmer model
475 performed the best, with CHERRY and vPST close behind (**Extended Data Fig. 10D**).
476 The kmer model notably includes implementation-specific changes to CHERRY that
477 appear to enable greater performance. Further, the lower vPST performance at broader
478 host taxonomic levels could be explained by the fact that the vPST genome
479 embeddings were not tuned for hosts. However, the ESM2-based model, which is more
480 comparable to vPST, does not perform well when evaluated at any confidence threshold
481 or host taxon rank. This directly demonstrates the importance of training on viral
482 datasets for viromics tasks.

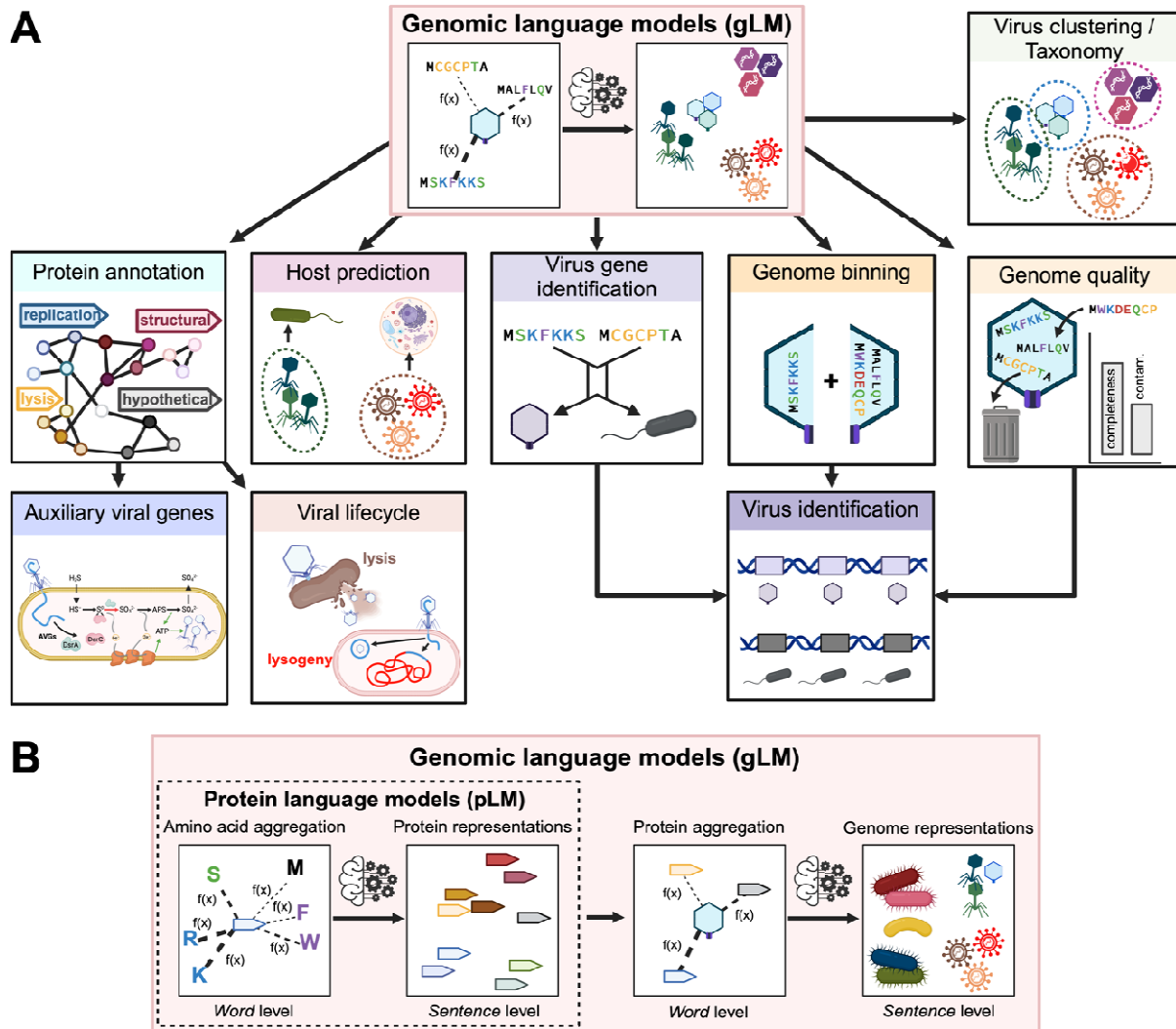
483

484 Discussion

485

486 Here, we presented the PST framework for modeling genomes as sets of proteins,
487 where each protein is initially represented by information-rich ESM2 protein
488 embeddings. The PST contextualizes the input protein embeddings and subsequently
489 yields genome representations as weighted averages of contextualized protein
490 embeddings, which can be targeted toward either protein-level or genome-level
491 downstream tasks. When pretrained on a large, diverse dataset of viral genomes, the
492 vPST demonstrated superior ability in understanding relationships among viral genomes
493 (**Fig. 2E**). At the protein level, the vPST protein embeddings demonstrated patterns of
494 broad function grouping, consistently clustering late gene proteins together (**Fig. 3B**).
495 Additionally, vPST often clustered capsid-fold-containing proteins that could not be
496 annotated by VOG with annotated capsid proteins (**Fig. 4A**), suggesting that vPST uses
497 inferred structural information for relating proteins. The vPST further showed high
498 sensitivity for annotation transfer (**Fig. 4B**). Performance for these protein-level tasks
499 could be further improved by finetuning the ESM2 pLM with viral sequences and by
500 training a vPST with a dual objective that more directly considers protein-protein and
501 genome-genome relationships. Finally, when applied toward a viral host prediction task,
502 the vPST genome embeddings were able to detect the true host species for the greatest
503 number of viruses when compared against two previously published host prediction
504 tools (**Fig. 5B**). We notably refrained from overanalyzing the subtle differences in
505 performance in the proof-of-concept host prediction task since there are numerous
506 training techniques beyond the scope of our work that could have resulted in a superior
507 vPST-based host prediction model. It is, therefore, important to emphasize that the
508 vPST-based host prediction model performed on par with (and sometimes better than)

509 existing host prediction tools without the vPST being initially tasked with host prediction
 510 and without significant training time.
 511



512 **Figure 6.** The PST can be a general-purpose microbial and viral genome language model. **A)** Potential
 513 downstream tasks of the pretrained vPST that represent commonly desired steps of a typical
 514 computational viromics pipeline. **B)** Example workflow of a genome language model based on the PST
 515 that could incorporate both microbial and viral input genome datasets.
 516

517
 518 It is imperative to reiterate that this superior performance in a variety of viromics
 519 tasks emerged despite not training the vPST with these objectives. Taken together, our
 520 results indicate that the vPST is suitable as a foundation model for common viromics
 521 tasks, such as virus identification, taxonomy, host prediction, protein annotation,
 522 genome binning, etc. (**Fig. 6A**). We anticipate that more thorough studies for
 523 downstream viromics problems will benefit from starting from our pretrained vPST.
 524 Additionally, finetuning the vPST can bring even greater performance for these
 525 downstream tasks. For example, finetuning an end-to-end host prediction model with a
 526 virus-host dataset would likely significantly improve predictive power compared to what
 527 we observed. Further, the iPHoP training dataset has limited diversity (**Extended Data**

528 **Fig. 10C**), which could suggest that the results here are not representative of true
529 performance. Nonetheless, our work has provided a guideline for a standalone vPST-
530 based host prediction tool.

531
532 There has been growing caution around biological foundation models due to
533 potential biosecurity threats such as generating novel pathogenic viruses or guiding
534 gain-of-function viral mutations. For example, the AlphaFold 3 web server does not
535 allow predictions for certain viral proteins³⁸, Evo excluded viruses with eukaryotic hosts
536 from its pretraining data⁴¹, and ESM3-open filtered viral sequences and select agents
537 from its training sets⁴². While developing vPST, we have assessed the ethical
538 implications of this viral foundation model and are having independent experts consider
539 these impacts before releasing the vPST code and model weights. We, however,
540 perceive the vPST to have a low biosecurity risk. First, only 0.2% of the vPST training
541 viruses infect humans. Of these, only 4 are on the CDC's list of bioterrorism agents
542 (<https://emergency.cdc.gov/agent/agentlist-category.asp>; Filoviridae viruses: Ebolavirus
543 and Marburgvirus), and 10 more are under surveillance by the National Respiratory and
544 Enteric Virus Surveillance System (<https://www.cdc.gov/nrevss/php/dashboard>). Further,
545 only 1% of the training viruses infect mammals, which would be the most likely viral
546 reservoirs that could spillover into human populations. Since our model was not trained
547 considering host identity, the low abundance of these viruses in the training dataset
548 likely minimizes their influence on the learned vPST embeddings. Second, the lowest
549 resolution of the vPST is at the protein level, meaning that it would be difficult to reverse
550 engineer a *de novo* viral genome using the vPST. While a nucleotide language model
551 reported the ability to generate *de novo* bacterial virus genomes³¹, the similarity of these
552 genomes to the training dataset was not investigated. One pitfall is that the model could
553 have been generating trivially *de novo* genomes that do not differ substantially from the
554 training data. Reverse engineering genomes from our protein-based work is further
555 complicated by the complexities of how human viruses tend to encode and express
556 genes (overlaps, alternative starts, alternative splicing, post-translational processing,
557 etc.). These molecular biology issues likely mean that achieving *in vivo* activity of a
558 generated viral genome would be challenging. We, therefore, perceive that the
559 demonstrated and potential future benefits (**Fig. 6A**) of our work to advance our
560 understanding of viruses outweigh any hypothetical threats that would require significant
561 resources to unleash.

562
563 Finally, our PST architecture, while trained on viral proteins and genomes for this
564 study, is agnostic to the source of the proteins and type of genomes. The only
565 requirements of our framework are the ordered protein sequences and genome strand
566 of each ORF. These requirements are more easily satisfied by microbial genomes,
567 where computational ORF calling is both accurate and common. However, our PST
568 could theoretically work with large enough datasets of experimentally determined ORFs
569 from eukaryotes as well. Nonetheless, we propose that our PST implementation is
570 equally appropriate for developing a microbial foundation model to solve challenges in
571 microbial genomics (**Fig. 6B**), which notably also include poor protein annotation rates
572 and high sequence divergence. In fact, our foundation vPST model was still useful for

573 host genome representations in the virus-host prediction task (**Fig. 5**), despite only
574 being trained on viruses.

575 **Data Availability**

576 Sources for publicly available viral genomes are listed in Supplementary Table 1.
577 Supplementary data specific to this manuscript, including protein FASTA files, protein
578 and genome embeddings, trained vPST model weights, and virus-host interaction
579 graphs, were deposited at DRYAD: (doi: 10.5061/dryad.d7wm37q8w). The repository
580 will be made public after the completion of our biosecurity review.

581
582 **Code Availability**

583 All code for the PST model architecture and analyses specific to this manuscript will be
584 released at: https://github.com/AnantharamanLab/protein_set_transformer. Specifically
585 for manuscript-associated analyses, Jupyter notebooks will be provided for each
586 method section that uses code. We will provide additional repositories for generating the
587 ESM2 protein embeddings, GenSLM ORF and genome embeddings, and HyenaDNA
588 genome embeddings that can be found in the main model repository above. The
589 repositories will be made public after the completion of our biosecurity review.

590
591 **Author contributions**

592 All authors (CM, AG, KA) conceived the project. CM conducted model and software
593 development, all analyses, results visualization, and content organization. AG and KA
594 provided project feedback. CM and KA wrote the manuscript draft. All authors (CM, AG,
595 KA) reviewed the results and edited and approved the manuscript.

596
597 **Acknowledgements**

598 This research was supported by the National Institute of General Medical Sciences of
599 the National Institutes of Health under award (R35GM143024) and by a National
600 Science Foundation award (2226451). CM was supported by a National Science
601 Foundation Graduate Research Fellowship and a University of Wisconsin-Madison
602 SciMed GRS Fellowship. AG acknowledges support from Jeanne M. Rowe. We thank
603 members of the Anantharaman lab for project discussion and feedback on the
604 manuscript. Model training and inference was performed using the resources and
605 assistance of the University of Wisconsin-Madison Center for High Throughput
606 Computing⁴³.

607
608 **Competing interests**

609 The authors declare no competing interests.

610 **Online Methods**

611

612 **Viral genome datasets**

613

614 We acquired viral genomes from 12 different publicly available sources^{1,17–27} as a
615 training dataset. For GTDB (r202), we used PhageBoost⁴⁴ (v0.1.7) with default settings
616 to identify integrated proviruses, filtering predictions that did not encode at least 20
617 proteins. We then filtered genomes that were not considered complete or high-quality as
618 defined by CheckV⁴⁵ (v1.0.1). We then dereplicated this set of genomes using a custom
619 workflow. We first used skani⁴⁶ (v0.1.0 sketch: --fast) to compute pairwise average
620 nucleotide identity (ANI) between all pairs of viruses. We constructed a graph where
621 edges connected viruses with $\geq 95\%$ ANI and $\geq 50\%$ genome coverage of the alignment
622 for both genomes. The edge weights were the product of ANI and coverage. We then
623 used the Markov clustering algorithm⁴⁷ (mcl v14-137 -l 2.0) to cluster this graph, taking
624 one genome from each cluster at random as a representative genome. For the test
625 dataset, we chose the most complete, least contaminated, and longest genome for each
626 viral operational taxonomic unit in IMG/VR v4¹¹, ensuring that each representative was
627 considered high-quality by CheckV. We then dereplicated this putative test dataset with
628 the training dataset using a similar approach as above with skani (--slow, $\geq 95\%$ ANI,
629 $\geq 85\%$ coverage) and mcl. We kept all viruses that did not cluster with training viruses.
630 For both datasets, we filtered out viruses predicted to encode only 1 protein. The final
631 number of viral genomes was 103,589 for the training dataset and 151,255 for the test
632 dataset.

633

634 For all viruses, we predicted protein open reading frames (ORFs) using the Python
635 bindings of prodigal called pyrodigal⁴⁸ (v2.3.0) for single-contig viruses and prodigal-gv
636 (v2.11.0) for viral metagenome-assembled genomes (vMAGs). We did not consider the
637 updates made by prodigal-gv⁴⁹ (include gene models for giant viruses and viruses using
638 alternative genetic codes) to be substantial enough to apply to the entire dataset given
639 the scale and distribution of the data. This led to 6,391,562 proteins for the training
640 dataset and 7,182,220 for the test dataset.

641

642 For the training viruses, viral taxonomy not provided by IMG/VR v3 was assigned
643 using geNomad⁴⁹ (v1.5.0) to get labels that were consistent with the current standards.
644 For the test viruses, we used the provided taxonomic labels since they were consistent
645 with current standards, and most were predicted using geNomad also. We did not
646 perform host prediction on these viruses, so host labels were either predicted by the
647 source database or are known due to integrated provirus prediction. The summary of
648 information for the training and test viruses can be found in **Supplementary Table 1**.

649

650 **ESM2 protein language model embeddings**

651

652 PyTorch (v2.1.0)⁵⁰ and fair-esm³ (v2.0.0) were used to obtain protein embeddings.
653 We refer to the ESM2 models “esm2_t6_8M_UR50D” (6 layers, 8M parameters, 320-
654 dimensional embedding) and “esm2_t30_150M_UR50D” (30 layers, 150M parameters,
655 640-dimensional embedding) as “esm-small” and “esm-large”, respectively. The amino

656 acid embeddings in each protein were averaged for a single d -dimensional vector. For
657 proteins longer than 20,000 amino acids, the sequence was split in half, and the
658 embeddings for each half were then averaged for the final embedding. This only
659 affected 1 bacterial protein in the host prediction analysis.

660

661 **The Protein Set Transformer model architecture**

662

663 The Protein Set Transformer (PST) was built using PyTorch (v2.0.0), PyTorch
664 Geometric⁵¹ (v2.3.1), and PyTorch-Lightning (v2.0.7). The PST draws inspiration from
665 deep learning of set-structured data like the SetTransformer¹² while using modifications
666 that are specific to pointsets¹³, which are sets whose items are d -dimensional vectors.
667 The PST, thus, models genomes as a set of proteins G_i where $g_{ij} \in G_i$ is the j th protein
668 in the i th genome. Each protein is initially represented by its d -dimensional ESM2
669 protein embedding, with $G_i \in R^{n_i \times d}$ where n_i is the number of proteins encoded in
670 genome G_i . We did not finetune the ESM2 models, so the ESM2 embeddings were
671 used as frozen inputs. For each protein g_{ij} , learnable embeddings for both the position
672 in the genome and for the encoding strand were concatenated to the ESM2
673 embeddings. The positional embeddings for proteins were relative to the positions of the
674 proteins in each genome and are used so the model learns relative ordering of proteins.
675 For fragmented genomes such as viral metagenome-assembled genomes (vMAGs), the
676 scaffolds were randomly oriented such that all proteins were numbered continuously
677 from the randomly chosen starting scaffold.

678

679 To account for the large variation in the number of proteins encoded by each
680 genome, we used a memory-efficient graph-based implementation that considers each
681 genome as a graph and each protein as nodes in the genome graph. Notably, each
682 individual genome matrix G_i is stacked for each minibatch, so there is no padding. Then,
683 an indexing pointer keeps track of the offsets (number of rows/proteins) for each
684 genome for efficient access. For memory efficiency and to model real fragmented
685 genomic data, we break each genome graph into subgraphs whose node sets include
686 15-50 mutually exclusive, contiguously located proteins. The size of each subgraph was
687 tuned and is, thus, fixed. These nodes are all fully connected in each subgraph such
688 that all proteins in each genome subgraph attend to each other in the PST encoder. We
689 prevent subgraphs with 1 node by adding possible singleton node cases to the previous
690 subgraph. The subgraph size (“chunk size”) hyperparameter is constant for all
691 genomes. Thus, a minibatch of N genomes is represented by a single graph $\mathcal{G} =$
692 $(X, E) = \text{Stack}(G_i)$ where X is the total number of proteins encoded by the N genomes.
693 E is the total number of protein-protein edges and is a function of the subgraph size and
694 number of proteins per genome.

695

696 The PST uses the encoder-decoder paradigm previously described with the
697 SetTransformer¹². The encoder uses multi-head self-attention to contextualize each
698 protein by the other proteins within the same genome. Then, the decoder uses multi-
699 head attention pooling to summarize the genome as a weighted average of
700 contextualized protein embeddings. To contextualize the proteins in each genome, we

701 used a graph-based implementation of multi-head scaled-dot product self-attention¹⁴ in
 702 each layer of the PST encoder:

$$\alpha_{ij} = \text{GraphSoftmax}\left(\frac{(W^Q x_i)(W^K x_j)}{\sqrt{d}}\right)$$

703

$$\text{MultiHeadAttn}(\mathcal{G}) : x_i^{(l+1)} = W^{Q,(l)} x_i^{(l)} + \sum_{j \in \mathcal{N}(i) \cup \{i\}} \alpha_{ij}^{(l)} W^{V,(l)} x_j^{(l)}$$

704 where $x_i^{(l)}$ is the embedding vector for the i th protein at the l th encoder layer, d is the
 705 protein embedding dimension. Likewise, $W^{\cdot,(l)}$ is the weight matrix for the query, key,
 706 and value at the l th encoder layer. $\mathcal{N}(i)$ is the set of protein neighbors for the i th protein
 707 in the same genome subgraph. α_{ij} is the scaled-dot product attention calculation.
 708 GraphSoftmax is a modified softmax function that only normalizes the attention values
 709 within the set of subgraphs that belong to the same genome. Thus, only proteins in the
 710 same subgraph attend to each other, but the attention values are normalized by all
 711 proteins in the genome. To enable multi-head attention, we split the input protein
 712 embeddings in the same number of chunks as the number of attention heads along the
 713 embedding dimension. After the self-attention calculation, we concatenate the outputs
 714 from each head back together. Further, we followed a pre-normalization strategy in
 715 which we normalized the input protein embeddings before the linear layers. Specifically,
 716 we used GraphNorm (implemented in PyTorch-Geometric) normalization operator that
 717 normalizes the proteins embeddings only within each genome. Additionally, we used the
 718 corresponding skip connections in which untransformed inputs are added to values
 719 post-attention. A full PST encoder layer can thus be mathematically represented as the
 720 following set of equations (1):

$$\begin{aligned} X^1 &= \text{GraphNorm}(X^0) \\ X^2 &= \text{MultiHeadAttn}(\mathcal{G}) \\ X^3 &= X^0 + X^2 \\ X^4 &= \text{GraphNorm}(X^3) \quad \#(1) \\ X^5 &= \text{FF}(X^4) \\ X^6 &= X^3 + X^5 \end{aligned}$$

721 where X represents the intermediate protein representations and X^0 is the input protein
 722 embeddings in the stacked batch. FF represents a 2-layer feedforward network with
 723 GELU (Gaussian error linear unit⁵²) activation and dropout after each layer. After the full
 724 PST encoder, a final GraphNorm operation was applied.

725

726 The PST decoder uses multi-head attention to compute a per-protein attention score
 727 to be used as the weights for a weighted average of protein embeddings over each
 728 genome. As described previously¹², multi-head attention pooling uses a learnable d -
 729 dimensional seed vector S as the query when computing attention. During the attention
 730 calculation, the contextualized protein embeddings X^C output from the PST encoder are
 731 projected onto S :

$$\text{Attn}(X^C, S) = \text{GraphSoftmax}\left(\frac{(W^Q S)(W^K X^C)}{\sqrt{d}}\right) \times (W^V X^C)$$

732 The attention values from this projection are used to weight X^C . After re-weighting, X^C is
733 averaged over each genome to produce the final genome outputs. The full set of PST
734 decoder equations is similar to the encoder (**Equation 1**):

$$\begin{aligned}X^0 &= \text{GELU}(WX^C) \\X^1 &= \text{GraphNorm}(X^0) \\X^2 &= \text{Attn}(X^1, S) \\X^3 &= X^0 + X^2 \\X^4 &= \text{GraphNorm}(X^3) \\X^5 &= \text{FF}(X^4) \\X^6 &= X^3 + X^5 \\X^7 &= \text{GraphPool}(X^6) \\X^G &= \text{FF}(X^7)\end{aligned}$$

735 where W is the weights of a linear layer. GraphPool is a pooling (mean) operator over
736 each genome graph that averages the contextualized weighted protein embeddings for
737 each genome. Each FF is a different 2-layer feedforward network with GELU activation
738 and dropout after each layer. X^G is the final genome embeddings. See **Extended Data**
739 **Fig. 1** for a pictorial representation the PST architecture.

740

741 **Training the viral Protein Set Transformer foundation model with triplet loss**

742

743 The foundation viral Protein Set Transformer (vPST) model was trained using a self-
744 supervised triplet loss objective $\mathcal{L}(\mathcal{G})$ as described previously¹³:

$$\begin{aligned}D(G^a, G^p, G^n) &= [\|f(G^a) - f(G^p)\|_2^2 - \omega_i \|f(G^a) - f(G^n)\|_2^2 + \alpha]_+ \\ \mathcal{L}(\mathcal{G}) &= \frac{1}{2N} \sum_{i=1}^N C_i [D(G_i^a, G_i^p, G_i^n) + D(G_i^a, G_i^{ip}, G_i^{in})]\end{aligned}$$

745 where G_i^a is the i th genome treated as an anchor point, G_i^p is the positive genome for
746 the i th genome, G_i^n is the negative genome for the i th genome, and G_i^i is the
747 augmented genome for the i th genome created using the PointSwap sampling
748 method¹³. $f(\cdot)$ is the function modeled by the full PST neural network, and $\|x - y\|_2^2$ is
749 the L2 (Euclidean) distance between the vectors x and y . C_i is class weight to amplify
750 the contribution to the loss for classes that are less abundant than others. We used the
751 viral realm of each virus as the class and compute C_i as an inverse abundance
752 frequency. Suppose that the i th genome belongs to viral realm k , then the class weight
753 is computed as:

$$\begin{aligned}F_k &= \frac{n_k}{N} \\ C_i &= \frac{1}{F_k}\end{aligned}$$

754 where n_k is the number of genomes in the training dataset belonging to viral realm k out
755 of N total genomes.

756

757 To account for the self-supervised choice of the negative genome, the scale factor
758 ω_i reweights the anchor-negative distance according to the following exponential decay
759 equation:

$$\omega_i = \exp\left(-\frac{CD(G_i^a, G_i^n)}{2(c\sigma)^2}\right) \#(2)$$

760 where $CD(X, Y)$ is the Chamfer distance between the genomes X and Y , c is a scaling
761 factor, and σ is the standard deviation of all Chamfer distances. $[\cdot]_+ = \max(0, \cdot)$, which
762 means that there is no contribution to the loss function for cases where the positive
763 genome is already closer to the anchor genome than the negative by a margin of α .
764 Thus, α is the farthest distance the negative genome needs to be from the anchor
765 compared to the positive genome. This restraint notably prevents representation
766 collapse that could occur in the naïve case of embedding the anchor and positive
767 genomes in the same position.

768
769 For a training minibatch, positive mining occurs in the input ESM2 embedding space
770 using Chamfer distance, before the PST forward pass and before concatenating
771 positional and strand embeddings. The Chamfer distance $CD(X, Y)$ between genomes X
772 and Y always uses the input ESM2 embeddings and is defined as follows:

$$CD(X, Y) = \frac{1}{|X|} \sum_{x \in X} \min_{y \in Y} \|x - y\|_2^2 + \frac{1}{|Y|} \sum_{y \in Y} \min_{x \in X} \|x - y\|_2^2 \#(3)$$

773 where $x \in X$ are the proteins from genome X and $y \in Y$ are the proteins from genome Y .
774 Intuitively, this means that the positive genome is defined as the most similar genome
775 based on cumulative distance of ESM2 protein embeddings, which should choose a
776 positive genome that encodes the most similar proteins.

777
778 Negative mining occurs in the PST embedding space and requires the positive
779 genome for a semi-hard sampling scenario. The only candidates for a negative genome
780 are those that are farther than the positive genome in the PST embedding space using
781 Euclidean distance, and we choose the first genome that is farther than the positive as
782 the negative in the semi-hard case. In cases where there are no genomes farther than
783 the positive genome in the PST embedding space, such as at the beginning of training
784 when the model weights have not been well-optimized, we loosen the semi-hard
785 sampling requirement and choose the genome closest to the positive genome as the
786 negative genome. Since negative mining is self-supervised, we use the exponential
787 decay reweighting factor ω_i to down-weight poor choices of a negative genome that are
788 actually very similar to the anchor genome. Notably, the ω_i reweighting factor depends
789 on the Chamfer distance (**Equation 2**) and, subsequently, the input ESM2 embeddings.
790 Thus, we implicitly consider the ESM2 embeddings as a ground truth for protein
791 representation when mining both the positive and negative genomes.

792 793 PointSwap sampling

794
795 When training the vPST, we used the data augmentation technique PointSwap
796 sampling¹³. During positive mining, we keep track of the most similar protein from the
797 positive genome X^p for each protein in the anchor genome $x \in X$ (**Equation 3**) as the
798 flow $x_i \rightarrow x_j^p$. We create the augmented genome X' as follows:

$$X' = \text{PointSwap}(X, X^p) = \{x'_0, \dots, x'_{n_i}\}, \text{ where}$$

$$x'_i = \begin{cases} x_j^p & \text{if } u_i < p \\ x_i & \text{otherwise} \end{cases}$$

799 where u_i is a set of samples from a standard uniform distribution $[0,1]$ and p is a rate of
800 protein swapping between genomes. This means that the augmented genome X' differs
801 from the anchor genome by swapping related proteins with the most related positive
802 genome, which intuitively mimics genetic variation. To form an augmented triplet with
803 the augmented genome as the positive genome, the negative genome is selected from
804 the set of augmented genomes in a minibatch using the procedure described above.

805

806 **Modified Leave-One-Group-Out cross validation and hyperparameter tuning**

807

808 To optimize the model hyperparameters (**Supplementary Table 3**), we used
809 Optuna⁵³ (v3.3.0) to iteratively sample hyperparameters in a direction that optimizes the
810 objective function using a Bayesian Tree-structured Parzen Estimator method. Model
811 performance was evaluated using a modified version of the Leave-One-Group-Out
812 (LOGO) cross validation (CV) strategy. Here, we considered the viral taxonomic realm
813 to be the group with 5 total groups: Duplodnaviria, Monodnaviria, Riboviria, Varidnaviria,
814 and Unknown / Other. We modified the LOGO strategy to always include Duplodnaviria
815 in each training fold, since this group of viruses accounted for 65.4% of the training
816 dataset. This resulted in training 4 separate models validated on the remaining viral
817 realms. Each of the 4 folds were synchronized during training to enable overall
818 performance monitoring as the average of each fold. This enabled real-time monitoring
819 of each tuning trial's performance. Thus, we were able to stop trials early depending on
820 several criteria using the average validation loss of each fold: (1) if the loss plateaued
821 (std of change less than $1e-6$) after having trained 3 epochs, (2) if the loss did not
822 decrease by 0.05 within 5 epochs, (3) if the current performance was worse than the
823 median performance of previous trials at the same training epoch, (4) if the model was
824 trained for 20 epochs, (5) if 24 hours passed, (6) if the loss was not finite. For number
825 3, this was maintained by the Optuna framework, and we required at least 1 complete
826 trial before this was enabled. In the case of early stopping due to reasons 1, 2, 3, and 6,
827 these trials were marked as pruned and not used by Optuna's median performance
828 calculation.

829

830 In total, we trained 16 complete, 16 failed, and 22 pruned trials using "esm-large"
831 protein embeddings as input and 45 complete, 1 failed, and 29 pruned trials using "esm-
832 small" protein embeddings as input. The only reason for failing was due to out-of-
833 memory errors on A100 80GB vRAM GPUs hosted by the University of Madison-
834 Wisconsin Center for High Throughput Computing⁴³. All trials were tuned using 1 GPU
835 since Optuna has limited support for GPU parallelism.

836

837 The final performance for each training iteration was the average validation loss from
838 each of the 4 models. Once the triplet loss of the best model setup decreased below
839 20.0, we chose the best hyperparameter configuration and trained 2 vPST models
840 corresponding to esm-small and esm-large protein embeddings that we refer to as pst-
841 small and pst-large, respectively. Each vPST model was trained on all genomes in the
842 training dataset without validation for 15 (pst-large, 33.7 hours) or 50 epochs (pst-small,

843 10.2 hours). Training was stopped once the training loss plateaued and did not
844 decrease by 0.05 within 5 epochs. During training of the final models, a learning rate
845 scheduler was used that linearly decreased the learning each epoch, and 50 (pst-large)
846 or 100 (pst-small) minibatches were accumulated before backward passes. We tested
847 batch accumulation sizes of 1, 25, 50, 100, and 250, and the above values led to the
848 best model.

849
850 For both tuning and training the final models, gradients were clipped to keep all
851 values below a magnitude of 1.0, and we used mixed precision training, using bfloat-16
852 data when available. These choices helped stabilize training. Our fold training
853 synchronization strategy and modified LOGO CV approach were implemented in a
854 custom package called “lightning-cv” available from the main model repository. This
855 package heavily relies upon and extends functionality in the lightning-fabric sub-library
856 of PyTorch-Lightning (v2.0.7).

857

858 **GenSLM open reading frame (ORF) and genome embeddings**

859

860 We used the 25M parameter GenSLM²⁸ foundation model (“genslm_25M_patric”,
861 downloaded September 2023) for our analyses since the output embedding dimension
862 (512) was on par with other protein and genome embeddings used. The GenSLM
863 foundation model is pretrained only on bacterial and archaeal nucleotide genes where
864 the gene sequences were broken into codons as input. The authors then finetuned the
865 foundation models on a dataset of SARS-CoV-2 genomes. However, it is not clear if
866 only the open reading frames (ORFs) from SARS-CoV-2 were included or if entire viral
867 genomes were used as input during finetuning. This is further complicated by the fact
868 that the protein-coding density of the SARS-CoV-2 genome is 71.2% (based on the
869 NCBI RefSeq reference sequence NC_045512.2). We chose to mimic the pretraining
870 setup and input the protein-coding ORFs for each virus in our datasets. Notably, we
871 used GenSLM as a nucleotide analog of ESM2, producing ORF embeddings akin to the
872 ESM2 protein embeddings. We used these ORF embeddings for protein/ORF analyses
873 and the average of these over each genome as genome embeddings for genome
874 analyses.

875

876 **HyenaDNA genome embeddings**

877

878 We used the HyenaDNA²⁹ model with the longest context size (1M nucleotides,
879 “large-1m”, downloaded from HuggingFace in November 2023) that has 6.5M
880 parameters. We converted all non-ACGTN nucleotides to Ns. Genomes larger than 1M
881 nucleotides were split into non-overlapping fragments of 1M nucleotides at most. Then
882 each fragment was tokenized and fed to the “large-1m” HyenaDNA model. The
883 embedding of each genomic fragment was averaged to produce the final genome
884 embedding. We also used this same averaging approach for fragmented genomes (ie
885 vMAGs), where the final genome embedding was the average of each fragment.

886

887 **Tetranucleotide frequency vectors as simple genome embeddings**

888

889 For each genome, we computed tetranucleotide frequency vectors
890 ([AAAA ... TTTT]) using the bionumpy⁵⁴ package (v1.0.8). We filtered all nucleotides
891 not in the canonical ACGT alphabet before calculation. For RNA viruses, U nucleotides
892 were represented by T for simplicity. For multi-scaffold viruses, these frequency vectors
893 were computed for each scaffold and then averaged over each scaffold. Throughout the
894 paper, these are referred to as “kmer”.

895

896 **Clustering genome and protein embeddings**

897

898 We constructed a similarity-weighted k -nearest neighbors (k NN) graph. The set of
899 k NN was computed using the faiss⁵⁵ (v1.8.0) Python bindings. For genome
900 embeddings, we used the divide-and-conquer IndexIVFFlat search index that splits the
901 input embeddings into n_{cells} Voronoi cells for faster retrieval. For the training dataset
902 ($n=103,589$ viral genomes), $n_{cells} = 2650$, and for the test dataset ($n=151,255$ viral
903 genomes), $n_{cells} = 3875$. Then, the L2 (Euclidean) distance D was used to identify the
904 closest genome neighbors. The L2 distances were converted to similarity scores S using
905 a Gaussian kernel:

$$S = \exp\left(-\frac{D^2}{\sqrt{d}}\right) \in [0, 1] \#(4)$$

906 where d is the dimensionality of the genome embedding. These similarity scores were
907 used as edge weights in the k NN graph.

908

909 For protein embeddings, we restricted the k NN search to only consider proteins that
910 belong to genomes within the same genome cluster. The protein embeddings were unit
911 normalized such that the L2-norm for each protein embedding equaled 1. Then, we
912 used cosine similarity to select the set of k NN, and the cosine similarity scores were
913 used as the k NN graph edge weights. For genome clusters with fewer than 78 proteins,
914 cosine similarity was brute-force computed for all pairs of proteins. For larger genome
915 clusters, the IndexIVFFlat partitioning method was used where $n_{cells} = \left\lceil \frac{n_{proteins}}{39} \right\rceil$ since
916 39 is the minimum number of data points per Voronoi cell. Genome or protein clusters
917 were then detected in the similarity-weighted k NN graph using the Leiden³² algorithm
918 Python implementation of iGraph (v0.11.3). The resolution values we used were 0.1
919 (“med”) and 1.0 (“high”) for genome clustering and 0.1 (“low”) and 0.5 (“med”) for protein
920 clustering. We do not include singletons as clusters for downstream analyses.

921

922 **Genome and protein clustering evaluation**

923

924 To compare clusters formed using different input embeddings, we computed the
925 cluster-wise average amino acid identity (AAI) between genomes, viral and host
926 taxonomic purity, and protein function purity for each cluster. Each of these cluster-level
927 metrics was weighted by the size of each cluster, specifically including unlabeled
928 genomes or proteins in the size calculation, and then summarized with a weighted
929 average:

$$C_{\text{summary}} = \sum_{i=1}^{n_{\text{clusters}}} C_i w_i$$

930 where

$$w_i = \frac{n_i}{\sum_{i=1}^{n_{\text{clusters}}} n_i} \#(5)$$

931 and C_i is the cluster-level metric.

932 Finally, the summary score C_{summary} was weighted by the proportion of non-singletons
933 for the given dataset P to penalize clustering iterations that did not include all genomes
934 or proteins:

$$C'_{\text{summary}} = C_{\text{summary}} \times P$$

935

936 Protein functional purity was computed using curated functional categories from
937 VOG or PHROG. To compute purity of clustering (viral or host taxonomy, protein
938 function), we used the information gain ratio I as a proxy for purity. For taxonomic purity,
939 we considered the case of clustering all genomes into a single cluster as the
940 background. For functional purity, we used the distribution of functional categories from
941 the annotation databases as the background. In either case, unlabeled proteins and
942 genomes were excluded during the entropy computation but included for the cluster size
943 weighting. Then we computed I as follows:

$$I = \frac{H_{\text{background}} - \sum_{i=1}^{n_{\text{clusters}}} H_i w_i}{H_{\text{background}}} \in [-h, 1]$$

944 where h is the information gain of the background compared to a uniform distribution
945 and H is the information entropy of each cluster with respect to a set of labels related to
946 viral taxonomy, host taxonomy, or protein function. The cluster size weight w_i is
947 computed as previously described (**Equation 5**). Values of I close to 0 indicate
948 clustering patterns with no improvement above background, while values of I close to 1
949 suggest maximal purity since there are few clusters with multiple labels. It is possible to
950 interpret I as a purity score since the backgrounds are not pure with respect to the
951 labels. Thus, a maximum I means that there is only a single label for a given cluster. I is
952 further weighted by the proportion of non-singletons:

$$I' = I \times P$$

953

954 **Average amino acid identity (AAI)**

955

956 We used mmseqs2⁵⁶ (v13.45111) and polars (v0.20.6) to compute the AAI between
957 pairs of viruses at scale. Given the large number of viruses in this study (>250k), we did
958 not exhaustively compute the AAI between all pairs of viruses (~32.5B). Instead, we
959 used heuristics implemented by mmseqs2 to only consider the AAI between viruses that
960 had any protein similarity detectable when using the mmseqs2 search settings: -s 7.5 -c
961 0.3 -e 1e-3. For each pair of viral genomes, we only retained the best hits for each
962 protein from each genome. Then, AAI was computed as the mean of protein-protein
963 sequence similarities computed by mmseqs2.

964

965 **Average amino acid identity (AAI) genome clustering**

966

967 During calculation of AAI for a pair of viral genomes, we tracked the proportion of
968 shared proteins relative to the total number of proteins from each genome u and v as S_u
969 and S_v , respectively. To cluster viral genomes using AAI, we constructed an edge-
970 weighted graph with edge weights corresponding to:

$$E_{uv} = \min(S_u, S_v) \times AAI \times 100 \in [0,1]$$

971 The edge weights, therefore, penalize cases where only a few proteins relative to the
972 total number of proteins in the genome with fewer proteins are used for the AAI
973 calculation. We then applied the Markov clustering algorithm⁴⁷ (mcl v14.137 -l 2.0 for
974 “med” resolution or 4.0 for “high” resolution), which uses edge-weight-guided
975 probabilistic random walks to cluster the AAI graph.

976

977 We only considered two levels of clustering, genus-level and family-level, using
978 thresholds previously described⁵⁷. Genus-level clustering sets the minimum AAI to 0.4
979 and requires either at least 16 shared proteins or $\min(S_u, S_v) \geq 0.2$. Family-level
980 clustering sets the minimum AAI to 0.2 and requires either at least 8 shared proteins or
981 $\min(S_u, S_v) \geq 0.1$.

982

983 **Protein functional annotation**

984

985 We used VOG (r219) and PHROG⁵⁸ (v4) databases for the annotation of viral
986 proteins. For VOG, which supplies profile Hidden Markov models (HMMs), we used
987 pyhmmmer (v0.9.0) with a bit score cutoff of 40. For PHROG, we used mmseqs2
988 (v13.45111) with the recommended search settings
989 (<https://phrogs.lmge.uca.fr/READMORE.php>). In both cases, we kept the best hit for
990 each protein with the max bit score. For each database, we curated the functional
991 categories of each annotation that we describe below.

992

993 For PHROG, which already provides an extensive set of 10 categories (including
994 unknown function), we manually readjusted certain categories. Our manual curation of
995 the PHROG database affected 1,937 out of 38,880 profiles. We renamed the following
996 categories for better intuition of the functional category: “DNA, RNA and nucleotide
997 metabolism” to “nucleotide metabolism”, “integration and excision” to “lysogeny”, and
998 “transcription regulation” to “gene expression”. We then dissolved the “moron, auxiliary
999 metabolic gene and host takeover” category for being too broad and relatively
1000 unrelated. These 461 profiles were split into the already existing “other”; the newly
1001 created “host takeover”, “lysogenic conversion”, “metabolic gene”; and the renamed
1002 “gene expression”, “lysogeny”, and “nucleotide metabolism” categories. Generic
1003 annotations like “membrane associated protein” and “ABC transporter” were put in the
1004 “other” category. We considered proteins involved in host replication and cell division
1005 inhibition, superinfection exclusion, anti-sigma factors, and defense against host
1006 antiviral proteins to be “host takeover”. Proteins that encoded toxins or
1007 antitoxins/resistance proteins were categorized as “lysogenic conversion.” Proteins
1008 directly involved in specific metabolic transformations were put in “metabolic gene,”
1009 while accessory or generic proteins like “nicotinamide mononucleotide transporter” were
1010 considered as “other”. These changes can be found in **Supplementary Table 5**.

1011
1012 VOG provides very broad categories: “Xr” for replication, “Xs” for structural, “Xh” for
1013 host-benefitting, “Xp” for virus-benefitting, and “Xu” for hypothetical proteins. The “Xh”
1014 and “Xp” categories are also ambiguous on what specific function the protein may
1015 perform. We, therefore, used text pattern matching on the specific HMM annotation
1016 descriptions to subdivide all HMMs into 9 categories: anti-host defense, exit, gene
1017 expression, integration, packaging, replication, structural, other, and unknown. Briefly,
1018 we separated terminases, portal proteins, and head packaging proteins from other
1019 structural proteins into a “packaging” category. Lysis, virion export, and budding HMMs
1020 were considered collectively as the “exit” group. “Integration” includes both integrases
1021 and excisionases as well as transposases. We considered all nucleotide metabolism
1022 and genome replication to be part of “replication”. To account for overlap in text
1023 matching, we enforced the following hierarchy: structural > packaging > exit >
1024 integration > gene expression > anti-host defense > replication > unknown > “RNA
1025 polymerases” > other. The final category for each HMM was, therefore, the highest in
1026 the hierarchy. We added RNA polymerases that did not indicate if they were replicative
1027 RNA-directed or transcriptive DNA-directed at the bottom to put these specific RNA
1028 polymerases in the “gene expression” category. Additionally, HMMs without matches
1029 were thus considered in the “other” category. The category for each VOG r219 HMM
1030 can be found in **Supplementary Table 6**, and the regex patterns used to categorize
1031 each HMM can be found in **Supplementary Table 7**.

1032 1033 **Protein attention scaling and analysis**

1034
1035 We computed the attention values as follows: Let $A_{ij} \in A_i$ be the mean attention
1036 score across all attention heads for the j th protein from the i th genome:

$$A_{ij} = \frac{1}{n_{\text{heads}}} \sum_{k=1}^{n_{\text{heads}}} A_{ijk}$$

1037 The sum of per-protein attention values for each genome is 1.0:

$$\sum_{A_{ij} \in A_i} A_{ij} = 1.0$$

1038 Given n_i , the number of proteins in the i th genome, and n_k , the number of proteins in
1039 the k th genome, and $n_i \neq n_k$, it follows that A_i and A_k are not directly comparable since
1040 the number of proteins each genome is not the same. More explicitly stated, consider
1041 $n_i = 2$ and $n_k = 4$, and $A_i = [0.5 \ 0.5]$ and $A_k = [0.5 \ 0.3 \ 0.05 \ 0.15]$. For the
1042 genome i , the model has randomly split attention to both proteins, while for genome k ,
1043 the model clearly has attended to the first protein more than the others, despite the
1044 numerical values being equivalent.

1045
1046 Therefore, to compare the vPST attention values per protein for each genome, we
1047 normalized the attention scores. We considered the background case for the attention
1048 distribution A_i to be a uniform distribution, ie $A_i \sim U(0,1;n_i)$ where $U(0,1;n_i)$ is a
1049 standard uniform distribution with probability $\frac{1}{n_i}$ of attending any of the n_i proteins. We

1050 then computed the distance between A_i and $U(0,1;n_i)$ using the normalized Kullbach-
1051 Leibler (KL) divergence:

$$D_i = \frac{H_i^{U(0,1;n_i)} - H_i^{A_i}}{H_i^{U(0,1;n_i)}} \in [0, 1]$$

1052 H^X is the entropy of the probability distribution X :

$$H(X) = - \sum_{x \in X} p(x) \log_2 p(x)$$

1053 We then rescale all per-protein attention values in A_i by the KL-divergence D_i to down-
1054 weight misleadingly large attention values that are uniformly (randomly) distributed:

$$A'_{ij} = A_{ij} \times D_i$$

1055 Thus, for cross-genome comparisons, we use the normalized attention scores $A'_{ij} =$
1056 $\{A'_{i1} \dots A'_{in_i}\}$. In our above example, $A'_i = [0.0 \ 0.0]$ and
1057 $A'_k = [0.0881 \ 0.0528 \ 0.0088 \ 0.0264]$.

1058
1059 Then, when analyzing the association of vPST attention with protein function, we
1060 first clustered the proteins using sequence identity (mmseqs2 v13.45111 -e 1e-3 -c 0.5 -
1061 s 7.5). We computed the max scaled attention A'_{ij} for all proteins in the same cluster.
1062 For function association analyses, we additionally retained 50 protein clusters with the
1063 largest A'_{ij} values for each functional category curated in the VOG r219 and PHROG
1064 databases.

1065 **Protein annotation improvement**

1066
1067 We considered all proteins unable to be annotated using the VOG r219 or PHROG
1068 databases as hypothetical proteins, where N_H is the number of hypothetical proteins.
1069 We computed the annotation improvement as a function of a genome clustering
1070 assignment and protein embedding. We used the same protein search settings
1071 described in the Methods section "**Clustering genome and protein embeddings**":
1072 cosine similarity on the unit-normalized protein embeddings, restricted to proteins that
1073 belong to genomes in the same genome cluster. For each protein, we searched for the
1074 closest non-self protein, scoring this as an improvement if this neighbor protein was
1075 annotated.
1076

$$T = \sum_{i=1}^{N_H} \begin{cases} \text{if NearestNeighbor}_i \text{ is annotated} & 1 \\ \text{else} & 0 \end{cases}$$

1077 Then, we computed the overall annotation improvement AP as the proportion of
1078 hypothetical proteins whose nearest neighbor was annotated:

$$AP = \frac{T}{N_H}$$

1079
1080 For a given genome embedding, we also computed the rate of change of AP over
1081 the number of nearest genome neighbors used for genome clustering using the
1082 numpy.polyfit (v1.23.5) function.
1083

1084 Protein function co-clustering

1085

1086 We used curated PHROG functional categories (**Supplementary Table 5**) to
1087 compute functional co-clustering, excluding the category for proteins of unknown
1088 function. Given a genome clustering and protein clustering configuration, for each
1089 protein cluster $P_i \in P$, we count the co-occurrence between pairs of functional
1090 categories u and v as the product of the number of proteins belonging to each category
1091 in the cluster:

$$C_i^{uv} = n_i^u \times n_i^v$$

1092 where n_i^u is the number of proteins in the i th protein cluster that belongs to category u .

1093 The observed co-occurrence C^{uv} between the functional categories u and v is defined
1094 as the sum of cluster-level co-occurrences:

$$C^{uv} = \sum_{i=1}^{|P|} C_i^{uv}$$

1095 To account for the distribution of PHROG annotation profiles, we computed an
1096 enrichment score against the background of the distribution of the 38,800 PHROG
1097 profiles:

$$E^{uv} = \frac{C^{uv}}{C_{background}^{uv}} \in [0, \infty]$$

1098 where $C_{background}^{uv}$ is computed analogously to C^{uv} except using relative abundances of
1099 the annotation profiles themselves instead of annotated proteins. To identify functional
1100 categories that co-occur frequently, we constructed a fully-connected graph with all
1101 PHROG functional categories as nodes and the corresponding edge weights E^{uv}
1102 between categories u and v . We then applied the Leiden algorithm with resolution 1.0 to
1103 identify sub-communities of co-occurring functions enriched above background.

1104

1105 Protein functional module detection

1106

1107 We defined the following protein functional modules based on curated functional
1108 categories (**Supplementary Tables 5 and 6**) and annotation text searches. For
1109 replication proteins in the PHROG database, we included proteins that were initially
1110 categorized as “nucleotide metabolism” and had matches to the following regex pattern
1111 “(i)DNA pol|single strand DNA binding|Par[AB]|DNA primase|(DNA)?[
1112]?helicase|repl|primosom|terminal|ribonucleo[st]ide(.*)?reductase|NDP reductase”. For
1113 VOG, we considered all hits in the replication category. For PHROG packaging
1114 modules, we included hits that belong to the “head and packaging” category and
1115 specifically matched the regex pattern “(i)terminase|portal”. For VOG, we only
1116 considered those in the “packaging” category. For PHROG DNA-interacting modules,
1117 we included all hits that belonged to either “nucleotide metabolism”, “lysogeny”, or “gene
1118 expression” categories. For VOG, all hits belonging to “replication”, “integration”,
1119 “packaging”, and “gene expression” were included. For PHROG late genes, annotations
1120 in the categories “tail”, “head and packaging”, “connector”, and “lysis” were retained.
1121 Likewise, for VOG, the categories “structural”, “exit”, and “packaging” were included.

1122

1123 We considered protein clusters to correspond to a specific functional module if they
1124 met the following module-specific criteria: For searches that only considered a single
1125 functional category (replication, packaging), we required at least 2 proteins from that
1126 category with different annotations. Due to the volume of data, we could not ensure that
1127 the 2 different annotations referred to truly different protein functions and not just the
1128 same function worded differently. For multi-category searches (late genes, DNA-
1129 interacting), we required at least 2 categories to be represented. In either case, we
1130 excluded protein clusters that had any annotated proteins outside the indicated
1131 functional categories to focus on protein clusters that most strongly fit our definition of
1132 functional modules.

1133 1134 **Capsid structure searches**

1135
1136 To quantify the frequency at which embedding-based protein clusters co-cluster
1137 VOG-detectable capsid proteins (VOG bit score ≥ 75) with proteins unable to be
1138 assigned function by VOG, we excluded all embedding-based protein clusters that did
1139 not solely consist of annotated capsids and hypothetical proteins. We then filtered this
1140 candidate set of proteins to keep those that fit the previous criteria at least 10 times
1141 among all clustering configurations. We additionally include sequence identity-based
1142 clusters (mmseqs2 v13.45111 cluster -s 7.5 -c 0.5) that also consisted of unannotated
1143 and capsid proteins as positive controls. This led to a total of 100,704 proteins for this
1144 analysis.

1145
1146 We used foldseek³⁴ (v9.427df8a) to convert our protein sequence database into a
1147 3Di-structure database using the ProstT5³⁵ model (downloaded July 2024; foldseek
1148 createdb with “—prostt5-model” option), which uses language tokens to represent
1149 structural features. We searched our 3Di-structure database against 295k structures
1150 from the Protein Data Bank³⁶ (PDB; downloaded using foldseek in July 2024) using
1151 default settings. We excluded all alignments with bit scores less than 100 and manually
1152 annotated the PDB structures as viral capsids using the following query at the PDB web
1153 server (<https://www.rcsb.org>): “capsid, major capsid, coat, minor capsid, virion”. We
1154 validated this approach by aligning AlphaFold 3-modeled³⁸ (<https://alphafoldserver.com>)
1155 monomer structures with the HK97 major capsid protein (2FS3) using TM-align⁵⁹
1156 implemented in the PDB web server. We choose 2 proteins with the highest scoring
1157 structural alignment as determined by foldseek, each from either proteins annotated
1158 with a VOG profile of unknown function or proteins undetected by VOG, for this
1159 analysis.

1160
1161 We then scored the proportion of unannotated proteins in each cluster that had a
1162 structural alignment with a PDB capsid protein. To summarize these proportions for
1163 each combination of clustering hyperparameters, genome embedding, and protein
1164 embedding, we computed a weighted average of these proportions, using the cluster
1165 size as the weight.

1166 1167 **Embedding UMAP visualization**

1168

1169 We used the Python implementation of the UMAP algorithm⁶⁰ (umap-learn v0.5.3)
1170 for embedding visualization only. For genome embeddings, we used 15 nearest
1171 neighbors defined using Euclidean distance. When computing the reduced embeddings,
1172 we jointly embed the genome embeddings of both the training and test datasets for
1173 each type of genome embedding into the same space. For protein embeddings, we first
1174 unit-normalized each protein embedding to have an L2-norm of 1. Then, we used 8
1175 nearest neighbors defined using cosine distance as this value gave the best visual
1176 separation. In both cases, we did not reduce the dimensionality before visualization, so
1177 the embeddings themselves were directly used as inputs to UMAP algorithm.

1178

1179 **Graph-based host prediction framework**

1180

1181 For the virus host prediction proof-of-concept, we modeled our framework off
1182 CHERRY³⁹, which applies graph learning on a virus-host interaction network $G = (X, E)$.
1183 Our implementation uses PyTorch (v2.2.2), PyTorch-Geometric (v2.5.2) and PyTorch-
1184 Lightning (v2.2.4). The interaction network is bipartite, meaning that there are 2 types of
1185 nodes: viral nodes $v_i \in V$ and host nodes $h_i \in H$. The total node set X is thus $X = V \cup$
1186 H . The edges E represent known virus-host pairs and may also include confident virus-
1187 host predictions that come from virus-host genome alignments (see **Host prediction**
1188 **training and test datasets** for more detail). Given G , the objective is a link prediction
1189 task to infer for any virus-host pair (v_i, h_i) the probability of an edge existing in the
1190 interaction network.

1191

1192 For the most comparable analyses, we designed our neural network architecture
1193 based on CHERRY: an encoder consisting of multiple Graph Convolution⁶¹ (GCN)
1194 layers and a decoder that performs the link prediction. The encoder propagates
1195 information in the genome embeddings among local neighborhoods. Specifically, the
1196 GCN encoder layers can mathematically be represented as:

$$e^{(l+1)} = \phi \left(\tilde{D}^{-\frac{1}{2}} \tilde{A} \tilde{D}^{\frac{1}{2}} e^{(l)} W^{(l)} \right) \# (6)$$

1197 where l is the layer index, \tilde{A} is the adjacency matrix with self-connections ($\tilde{A} = A + I$, I
1198 is the identity matrix,), and \tilde{D} is the diagonal matrix where $\tilde{D}_{ii} = \sum_j \tilde{A}_{ij}$. ϕ is the
1199 activation function, and $W^{(l)}$ is the l -layer model weights. $e^{(0)} \in \mathbb{R}^{|X| \times N}$ is the input
1200 genome embedding where N is the size of the embedding. To compare our work to
1201 CHERRY, which uses tetranucleotide frequency genome embeddings, we substitute the
1202 genome embedding with either the vPST genome embeddings or the simple average of
1203 the ESM2 protein embeddings over each genome for both the viruses and hosts.

1204

1205 The decoder is a 2-layer feedforward neural network that takes the outputs from the
1206 encoder as input. We consider all possible virus-host pairs $(v_i, h_j) \in X$ as a query set Q
1207 where each pair is represented by the difference in encoder embedding:

$$q_{ij} = \text{encoder}(v_i) - \text{encoder}(h_j)$$

1208 The decoder, therefore, is mathematically written as:

$$\begin{cases} q_{ij}^{(l+1)} = \phi(q_{ij}^{(l)} \theta^{(l)}) \\ \text{decoder}(q_{ij}) = \text{sigmoid}(q_{ij}^{L-1}) \end{cases}$$

1209 where $q_{ij}^{(l)}$ is the hidden feature in the l th layer out of L total decoder layers, and
1210 $q_{ij}^{(0)} = q_{ij}$. ϕ is the activation function, and $\theta^{(l)}$ represents the weights of the l th fully
1211 connected layer.

1212

1213 **Host prediction training and test datasets**

1214

1215 For the virus host prediction proof-of-concept, we modeled our framework off
1216 CHERRY³⁹, which applies graph learning on a virus-host interaction network. To
1217 construct the network of known virus-host pairs, we used the train and test datasets
1218 from iPHoP⁴⁰. Specifically, the train dataset included 3628 complete bacterial and
1219 archaeal viruses from NCBI RefSeq prior to 2021. The iPHoP test dataset consisted of
1220 1636 complete bacterial and archaeal viruses from NCBI GenBank, distinct from the
1221 training dataset. Although both datasets indicate the taxonomy of the host, they do not
1222 provide specific genome accessions to link the viruses, which are necessary to
1223 construct the interaction network.

1224

1225 For the training dataset, we used the Virus-Host Database⁶² (accessed April 2024) to
1226 determine the full host taxonomy. We then selected either the NCBI RefSeq
1227 representative sequence associated with the host taxonomy, if one existed, or the most
1228 complete (longest and assembly_level == "Complete Genome") genome from NCBI
1229 GenBank (accessed May 2024). We included all hosts in the Virus-Host Database if
1230 there were multiple such as in the case of viruses with a relatively broad host range.
1231 The set of hosts notably includes multiple strains of the same species or species of the
1232 same genus as indicated in the Virus-Host Database. Then, any strain information was
1233 ignored, so the lowest level of evaluation was at the host species.

1234

1235 We performed a similar search for the test dataset using the information provided in
1236 Supplementary Table 2 of iPHoP⁴⁰. We divided the test virus hosts whose species ranks
1237 were unknown (ie "Wolbachia sp.") into 2 different sets. If these hosts were already in
1238 the set of hosts for the training dataset, we did not retrieve any new host genomes.
1239 Instead, we considered all hosts currently in the set of hosts with the same genus as
1240 potential hosts for these viruses. For new hosts, we used the same search criteria as
1241 above to add a single new host for each of these viruses. This resulted in a total of 805
1242 host genomes, corresponding to 594 unique host species.

1243

1244 **Constructing the virus-host interaction network**

1245

1246 To construct the virus-host interaction network, we constructed the heterogeneous
1247 graph G that has 2 node types (virus, host) and 2 edge types (virus-related to-virus,
1248 virus-infects-host). For the virus-host edges, we included all virus-host pairs identified
1249 above, meaning that G includes both training and test viruses. We notably deviated from
1250 the CHERRY implementation by excluding confident host predictions that would have
1251 come from virus-host BLASTn genome alignments (proviruses) or CRISPR spacers.
1252 This deviation is not concerning since we focused on the relative performance of our
1253 vPST genome embeddings compared to other tools and genome embeddings, rather
1254 than absolute predictive ability.

1255
1256
1257
1258
1259
1260
1261
1262
1263
1264
1265

To select virus-virus edges representing pairs of similar viruses, we used a protein sharing network clustering approach when using tetranucleotide frequency genome vectors. We first excluded all singleton proteins since these do not inform about genome-genome relatedness and only serve to account for the proportion of gene sharing relative to the total number of proteins/protein clusters in each genome. After protein clustering using mmseqs2 (v13.45111 -s 7.5 -e 1e-3 -c 0.5) and filtering singleton proteins, we constructed a sparse ($|V| \times n_{pc}$) presence-absence matrix where n_{pc} is the total number of protein clusters in the dataset. Each row represents what protein clusters are encoded in the indicated genome. We then computed the dice similarity S_{ij} for each pair of genomes as:

$$S_{ij} = \frac{2|v_i \cdot v_j|}{|v_i| + |v_j|}$$

1266 where v_i and v_j are the row presence-absence vectors for the i th and j th genomes,
1267 respectively. We then constructed a graph with all viruses where the edges are
1268 $S_{ij} | S_{ij} > 0$. To understand which viruses were considered related, we clustered this
1269 graph with the Leiden algorithm with a resolution of 0.1. Edges were created in the
1270 interaction graph between all viruses in the same gene-sharing clusters. For the other
1271 genome embeddings we tested, we considered pairs of viruses to be related if their
1272 genome embeddings were at least 90% similar based on a Gaussian-kernel of
1273 Euclidean distances (**Equation 4**). We then pruned these embedding-based virus-virus
1274 connections to only maintain the top 15 neighboring viruses for each virus.

1275
1276
1277

Host prediction model training

1278 We trained new graph-based host predictions models using the iPHoP training
1279 dataset, swapping the genome representations for vPST genome embeddings or the
1280 simple average of the ESM2 protein embeddings over each genome. Our
1281 implementation used PyTorch (v2.1.2) and PyTorch-Geometric (v2.4.0). We used a
1282 binary cross entropy loss objective for the link prediction task to classify the edge E_{ij} as
1283 existing (1) or not (0):

$$\mathcal{L} = -\frac{1}{N} \sum_{k=1}^N y_k \log(p(y_k)) + (1 - y_k) \log(1 - p(y_k))$$

1284 where y_k is the discretized final output for the k th virus-host pair from the model
1285 decoder, given a probability threshold for whether an edge E_{ij} is predicted to exist.

1286

1287 During training, we randomly split all virus-host edges $E = \{E^{(T)}, E^{(V)}\}$ into disjoint
1288 training $E^{(T)}$ and validation $E^{(V)}$ sets at an 80:20 ratio. We additionally randomly
1289 sampled negative edges $E' = \{E'^{(T)}, E'^{(V)}\}$ that do not exist in the virus-host interaction
1290 network G to provide the model with negative examples (implemented by PyTorch-
1291 Geometric). The negative edge sets $E'^{(T)}$ and $E'^{(V)}$ are also disjoint, and $|E| = |E'|$ so
1292 that there was not label imbalance. During the message-passing stage of the model
1293 encoder, only the real edges E are used. After message passing updates the node

1294 representations, we used $E \cup E'$ for decoding and inference with both real and negative
1295 edges. Therefore, $N = |E^{(\cdot)} \cup E'^{(\cdot)}|$ for either the training or validation edges. Since the
1296 prediction task does not depend on virus-virus edges, these edges were not split or
1297 negatively sampled. This means that the graph structure and message passing consider
1298 all viruses, not just training viruses. Thus, during training, we masked any virus-host
1299 edges that contain test viruses in the loss computation to prevent data leakage.

1300
1301 Although we strived to implement a nearly 1:1 model with the original CHERRY
1302 implementation, our implementation and training deviates in 3 ways. (1) We allowed
1303 separate learnable weights for each type of edge (virus-virus, virus-host, and host-virus)
1304 in the message-passing encoder layers by updating **Equation 6**:

$$W^{(l)} = \begin{cases} W_{vv}^{(l)} & \text{virus-virus edges} \\ W_{vh}^{(l)} & \text{virus-host edges} \\ W_{hv}^{(l)} & \text{host-virus edges} \end{cases}$$

1305 $W_{hv}^{(l)}$ and subsequently host-virus edges are required to ensure reciprocal message
1306 passing between virus and host nodes despite the intuitive way of representing virus-
1307 host edges as directed. The native CHERRY implementation does not allow for edge
1308 type-specific weights, instead sharing weights for all edge types.

1309
1310 (2) Due to modeling G as a heterogeneous graph, the message passing layer is not
1311 a true Graph Convolution (GCN) layer, which is not implemented for heterogeneous
1312 graphs in PyTorch-Geometric. Specifically, we use a generalization of the GCN layer⁶³
1313 that allows for heterogeneous graph learning with multiple node and edge types. For
1314 this layer, however, the node update equations for this layer and the GCN layer are
1315 identical, but there may be PyTorch-Geometric implementation-specific differences
1316 beyond changing node representations.

1317
1318 (3) We explored a more sophisticated technique for handling the training and
1319 validation splits for link-level tasks that we refer to as “disjoint training”. Specifically, we
1320 divided the real training edges $E^{(T)}$ into 2 disjoint sets $E^{(T)} = \{E^{(MP)}, E^{(D)}\}$ where $E^{(MP)}$
1321 are edges only used for message passing (node updates) and $E^{(D)} \cup E'^{(T)}$ are edges
1322 only used for supervision (decoding and inference). Specifically, $E^{(D)} \cup E'^{(T)}$ are the
1323 edges used for link prediction. We only considered a 70:30 split for $E^{(MP)}$ and $E^{(D)}$ for
1324 this study when this was enabled. This modification is analogous to splitting training
1325 data into separate training and validation sets to prevent data leakage but only for
1326 training edges.

1327
1328 To decouple the effect of these 3 differences from the choice of node embeddings,
1329 we trained a model that is nearly faithful to the CHERRY implementation without these
1330 changes (barring the required change #2), and then we trained a separate model using
1331 tetranucleotide frequency genome embeddings (kmer) that enables our changes. Thus,
1332 the CHERRY and “kmer” model use the same virus-host interaction graph as input, but
1333 the “kmer” models explored the effects of changes #1 and #3.

1334

1335 To lightly optimize hyperparameters, we sampled from sets of intuitive values for the
1336 number of encoder layers, decoder hidden dimensions, learning rate, whether to enable
1337 disjoint training (at a 70:30 split), and whether to allow edge specific-weights in the
1338 encoder or not. We did not dilate the input embedding dimension in the encoder layers.
1339 For the 2-layer feedforward decoder network, we only chose values smaller or equal to
1340 the input embedding dimension for the first layer. The second layer dimensions were
1341 then required to be strictly less than the first layer dimensions. See **Supplementary**
1342 **Table 9** for the values sampled for each hyperparameter. We applied the same random
1343 seed when training each iteration and chose the best model based good overall
1344 performance and lowest validation loss at the end of 150 training epochs. We defined
1345 “good” overall performance as a validation loss curve that was monotonically
1346 decreasing over or constant at the end of training time. We selected a total of 4 models
1347 that were the best: CHERRY without the above changes and 3 that allowed the above
1348 implementation changes and used different genome embeddings. All models were
1349 trained with a dropout of 0.25 after the encoder and after each decoder feedforward
1350 layer. We used the ReLU activation function after each layer.

1351 1352 **Host prediction model evaluation**

1353
1354 iPHoP (v1.3.3) and each of the 4 trained models were evaluated using the iPHoP
1355 test dataset (see “**Host prediction training and test datasets**”). For the 4 graph-based
1356 models, we considered all test virus-host pairs for link prediction and retained only those
1357 $\geq 75\%$ confidence, which is the minimum for iPHoP, or $\geq 90\%$ confidence. All virus-host
1358 pairs were considered to enable resolution at each host taxonomic rank. However, we
1359 only evaluated if the true host taxon was among the predictions above the confidence
1360 threshold, so not all predictions were analyzed. Specifically, we computed the proportion
1361 of the iPHoP test viruses whose true host taxon was confidently predicted.

1362
1363 Since there were notably a nontrivial number of viruses in the iPHoP test dataset
1364 that were similar to those in the vPST training dataset based on AAI (see “**Average**
1365 **amino acid identity (AAI) genome clustering**”), we filtered these viruses out using
1366 several similarity cutoffs to evaluate their effects on our interpretation of the host
1367 prediction results.

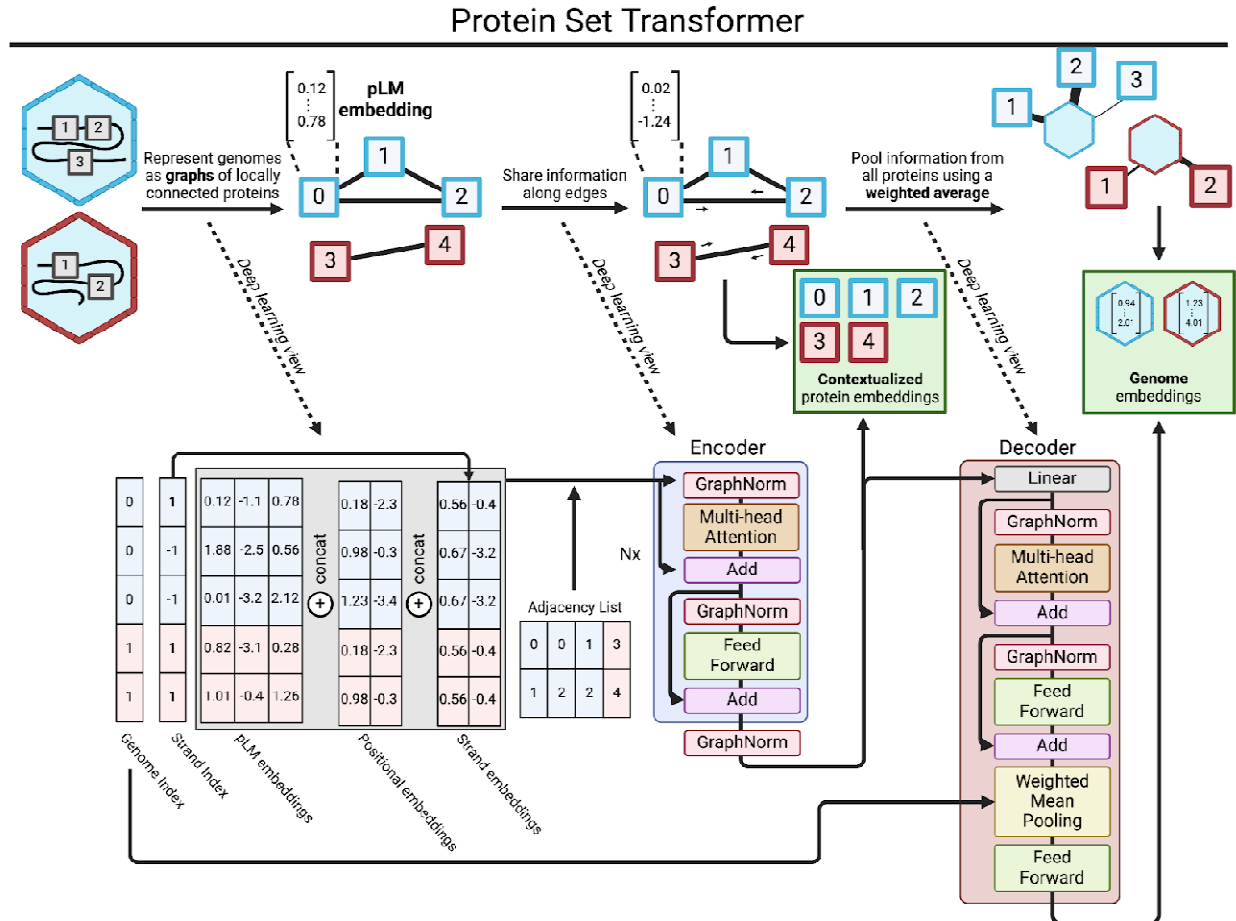
1368 **References**

- 1369 1. Roux, S. *et al.* Ecogenomics and potential biogeochemical impacts of globally
1370 abundant ocean viruses. *Nature* **537**, 689–693 (2016).
- 1371 2. Kieft, K. *et al.* Virus-associated organosulfur metabolism in human and environmental
1372 systems. *Cell Rep* **36**, 109471 (2021).
- 1373 3. Lin, Z. *et al.* Evolutionary-scale prediction of atomic-level protein structure with a
1374 language model. *Science* **379**, 1123–1130 (2023).
- 1375 4. Rives, A. *et al.* Biological structure and function emerge from scaling unsupervised
1376 learning to 250 million protein sequences. *Proceedings of the National Academy of
1377 Sciences* **118**, e2016239118 (2021).
- 1378 5. Chandra, A., Tünnermann, L., Löfstedt, T. & Gratz, R. Transformer-based deep
1379 learning for predicting protein properties in the life sciences. *eLife* **12**, e82819 (2023).
- 1380 6. Flamholz, Z. N., Biller, S. J. & Kelly, L. Large language models improve annotation of
1381 prokaryotic viral proteins. *Nat Microbiol* **9**, 537–549 (2024).
- 1382 7. Li, B. & Liang, G. ESM-PVP: Identification and classification of phage virion proteins
1383 with a large pretrained protein language model and an MLP neural network.
1384 2023.12.29.573676 Preprint at <https://doi.org/10.1101/2023.12.29.573676> (2023).
- 1385 8. Liu, D., Young, F., Robertson, D. L. & Yuan, K. Prediction of virus-host association
1386 using protein language models and multiple instance learning. 2023.04.07.536023
1387 Preprint at <https://doi.org/10.1101/2023.04.07.536023> (2023).
- 1388 9. Andrade-Martínez, J. S. *et al.* Computational Tools for the Analysis of Uncultivated
1389 Phage Genomes. *Microbiol Mol Biol Rev* **86**, e00004-21.
- 1390 10. Hwang, Y., Cornman, A. L., Kellogg, E. H., Ovchinnikov, S. & Girguis, P. R. Genomic
1391 language model predicts protein co-regulation and function. *Nat Commun* **15**, 2880
1392 (2024).
- 1393 11. Camargo, A. P. *et al.* IMG/VR v4: an expanded database of uncultivated virus
1394 genomes within a framework of extensive functional, taxonomic, and ecological
1395 metadata. *Nucleic Acids Res* **51**, D733–D743 (2023).
- 1396 12. Lee, J. *et al.* Set Transformer: A Framework for Attention-based Permutation-
1397 Invariant Neural Networks. in *Proceedings of the 36th International Conference on
1398 Machine Learning* 3744–3753 (PMLR, 2019).
- 1399 13. Arsomngern, P., Long, C., Suwajanakorn, S. & Nutanong, S. Towards Pointsets
1400 Representation Learning via Self-Supervised Learning and Set Augmentation. *IEEE
1401 Transactions on Pattern Analysis and Machine Intelligence* **45**, 1201–1216 (2023).
- 1402 14. Vaswani, A. *et al.* Attention is All you Need. in *Advances in Neural Information
1403 Processing Systems* vol. 30 (Curran Associates, Inc., 2017).
- 1404 15. Devlin, J., Chang, M.-W., Lee, K. & Toutanova, K. BERT: Pre-training of Deep
1405 Bidirectional Transformers for Language Understanding. Preprint at
1406 <http://arxiv.org/abs/1810.04805> (2019).
- 1407 16. Schroff, F., Kalenichenko, D. & Philbin, J. FaceNet: A Unified Embedding for Face
1408 Recognition and Clustering. in *2015 IEEE Conference on Computer Vision and
1409 Pattern Recognition (CVPR)* 815–823 (2015). doi:10.1109/CVPR.2015.7298682.
- 1410 17. Tisza, M. J. & Buck, C. B. A catalog of tens of thousands of viruses from human
1411 metagenomes reveals hidden associations with chronic diseases. *Proceedings of the
1412 National Academy of Sciences* **118**, e2023202118 (2021).

- 1413 18. Cook, R. *et al.* Hybrid assembly of an agricultural slurry virome reveals a diverse and
1414 stable community with the potential to alter the metabolism and virulence of
1415 veterinary pathogens. *Microbiome* **9**, 65 (2021).
- 1416 19. Gazitúa, M. C. *et al.* Potential virus-mediated nitrogen cycling in oxygen-depleted
1417 oceanic waters. *ISME J* **15**, 981–998 (2021).
- 1418 20. Gregory, A. C. *et al.* Marine DNA Viral Macro- and Microdiversity from Pole to Pole.
1419 *Cell* **177**, 1109-1123.e14 (2019).
- 1420 21. Parks, D. H. *et al.* GTDB: an ongoing census of bacterial and archaeal diversity
1421 through a phylogenetically consistent, rank normalized and complete genome-based
1422 taxonomy. *Nucleic Acids Research* **50**, D785–D794 (2022).
- 1423 22. Camarillo-Guerrero, L. F., Almeida, A., Rangel-Pineros, G., Finn, R. D. & Lawley, T.
1424 D. Massive expansion of human gut bacteriophage diversity. *Cell* **184**, 1098-1109.e9
1425 (2021).
- 1426 23. Gregory, A. C. *et al.* The Gut Virome Database Reveals Age-Dependent Patterns of
1427 Virome Diversity in the Human Gut. *Cell Host Microbe* **28**, 724-740.e8 (2020).
- 1428 24. Roux, S. *et al.* IMG/VR v3: an integrated ecological and evolutionary framework for
1429 interrogating genomes of uncultivated viruses. *Nucleic Acids Research* **49**, D764–
1430 D775 (2021).
- 1431 25. Michniewski, S. *et al.* A new family of “megaphages” abundant in the marine
1432 environment. *ISME COMMUN.* **1**, 1–4 (2021).
- 1433 26. ter Horst, A. M. *et al.* Minnesota peat viromes reveal terrestrial and aquatic niche
1434 partitioning for local and global viral populations. *Microbiome* **9**, 233 (2021).
- 1435 27. Brum, J. R. *et al.* Patterns and ecological drivers of ocean viral communities.
1436 *Science* **348**, 1261498 (2015).
- 1437 28. Zvyagin, M. *et al.* GenSLMs: Genome-scale language models reveal SARS-CoV-2
1438 evolutionary dynamics. *The International Journal of High Performance Computing*
1439 *Applications* **37**, 683–705 (2023).
- 1440 29. Nguyen, E. *et al.* HyenaDNA: Long-Range Genomic Sequence Modeling at Single
1441 Nucleotide Resolution. in *Advances in Neural Information Processing Systems* vol. 36
1442 43177–201 (Curran Associates, Inc., 2023).
- 1443 30. Peng, C., Shang, J., Guan, J., Wang, D. & Sun, Y. ViraLM: Empowering Virus
1444 Discovery through the Genome Foundation Model. 2024.01.30.577935 Preprint at
1445 <https://doi.org/10.1101/2024.01.30.577935> (2024).
- 1446 31. Shao, B. A long-context language model for deciphering and generating
1447 bacteriophage genomes. 2023.12.18.572218 Preprint at
1448 <https://doi.org/10.1101/2023.12.18.572218> (2024).
- 1449 32. Traag, V. A., Waltman, L. & van Eck, N. J. From Louvain to Leiden: guaranteeing
1450 well-connected communities. *Sci Rep* **9**, 5233 (2019).
- 1451 33. Liu, X., Jiang, H., Gu, Z. & Roberts, J. W. High-resolution view of bacteriophage
1452 lambda gene expression by ribosome profiling. *Proceedings of the National Academy*
1453 *of Sciences* **110**, 11928–11933 (2013).
- 1454 34. Barrio-Hernandez, I. *et al.* Clustering predicted structures at the scale of the known
1455 protein universe. *Nature* **622**, 637–645 (2023).
- 1456 35. Heinzinger, M. *et al.* Bilingual Language Model for Protein Sequence and Structure.
1457 2023.07.23.550085 Preprint at <https://doi.org/10.1101/2023.07.23.550085> (2024).

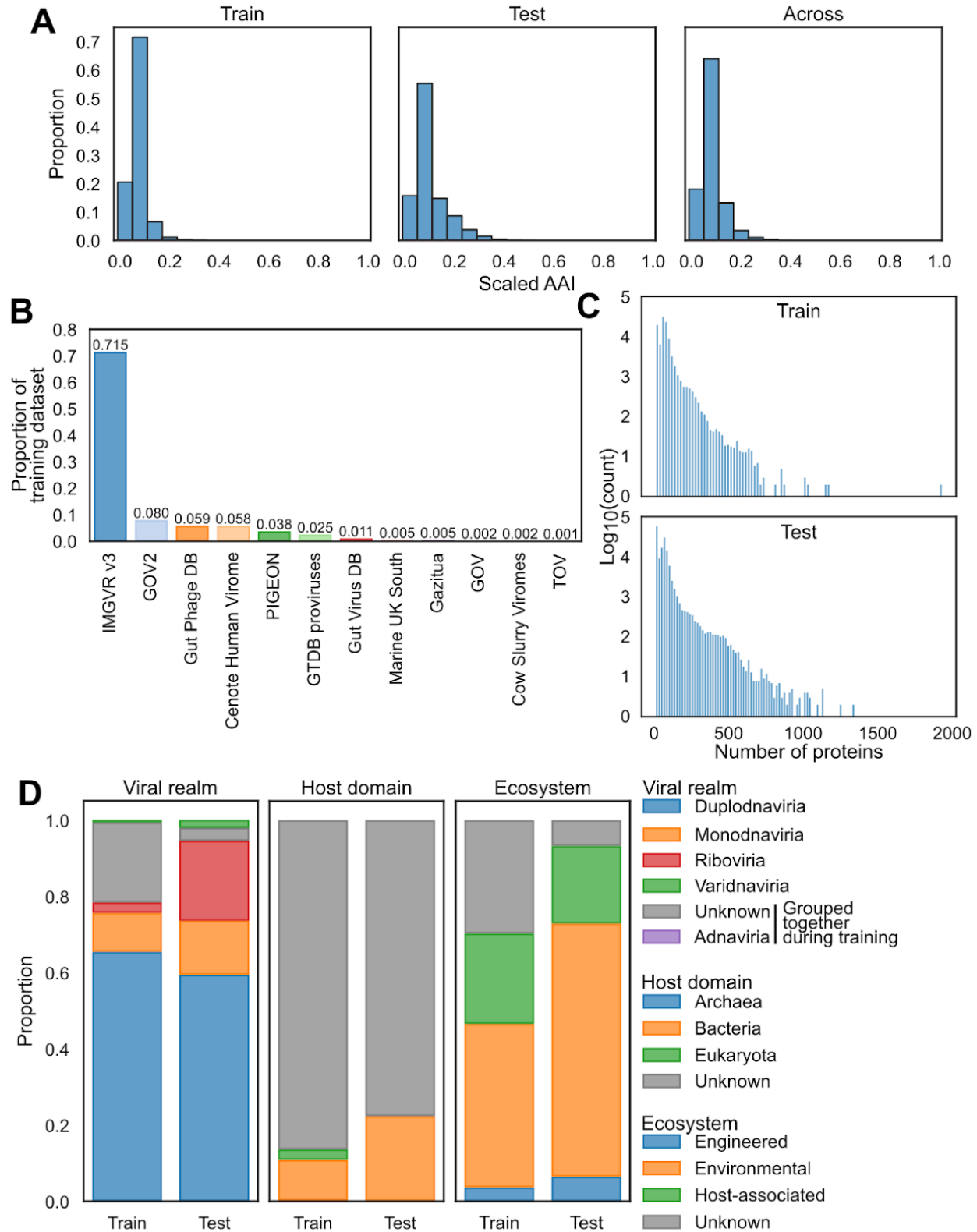
- 1458 36. Varadi, M. *et al.* PDBe and PDBe-KB: Providing high-quality, up-to-date and
1459 integrated resources of macromolecular structures to support basic and applied
1460 research and education. *Protein Science* **31**, e4439 (2022).
- 1461 37. Gan, L. *et al.* Capsid Conformational Sampling in HK97 Maturation Visualized by X-
1462 Ray Crystallography and Cryo-EM. *Structure* **14**, 1655–1665 (2006).
- 1463 38. Abramson, J. *et al.* Accurate structure prediction of biomolecular interactions with
1464 AlphaFold 3. *Nature* **630**, 493–500 (2024).
- 1465 39. Shang, J. & Sun, Y. CHERRY: a Computational methoD for accuratE pRediction of
1466 virus–pRokarYotic interactions using a graph encoder–decoder model. *Briefings in*
1467 *Bioinformatics* **23**, bbac182 (2022).
- 1468 40. Roux, S. *et al.* iPHoP: An integrated machine learning framework to maximize host
1469 prediction for metagenome-derived viruses of archaea and bacteria. *PLOS Biology*
1470 **21**, e3002083 (2023).
- 1471 41. Nguyen, E. *et al.* Sequence modeling and design from molecular to genome scale
1472 with Evo. 2024.02.27.582234 Preprint at <https://doi.org/10.1101/2024.02.27.582234>
1473 (2024).
- 1474 42. Hayes, T. *et al.* Simulating 500 million years of evolution with a language model.
1475 2024.07.01.600583 Preprint at <https://doi.org/10.1101/2024.07.01.600583> (2024).
- 1476 43. Center for High Throughput Computing. Center for High Throughput Computing.
1477 (2006) doi:10.21231/gnt1-hw21.
- 1478 44. Sirén, K. *et al.* Rapid discovery of novel prophages using biological feature
1479 engineering and machine learning. *NAR Genomics and Bioinformatics* **3**, lqaa109
1480 (2021).
- 1481 45. Nayfach, S. *et al.* CheckV assesses the quality and completeness of metagenome-
1482 assembled viral genomes. *Nat Biotechnol* **39**, 578–585 (2021).
- 1483 46. Skani enables accurate and efficient genome comparison for modern metagenomic
1484 datasets. *Nat Methods* **20**, 1633–1634 (2023).
- 1485 47. Van Dongen, S. Graph Clustering Via a Discrete Uncoupling Process. *SIAM J.*
1486 *Matrix Anal. Appl.* **30**, 121–141 (2008).
- 1487 48. Larralde, M. Pyrodigal: Python bindings and interface to Prodigal, an efficient
1488 method for gene prediction in prokaryotes. *Journal of Open Source Software* **7**, 4296
1489 (2022).
- 1490 49. Camargo, A. P. *et al.* Identification of mobile genetic elements with geNomad. *Nat*
1491 *Biotechnol* 1–10 (2023) doi:10.1038/s41587-023-01953-y.
- 1492 50. Paszke, A. *et al.* Automatic differentiation in PyTorch. (2017).
- 1493 51. Fey, M. & Lenssen, J. E. Fast Graph Representation Learning with PyTorch
1494 Geometric. Preprint at <http://arxiv.org/abs/1903.02428> (2019).
- 1495 52. Hendrycks, D. & Gimpel, K. Gaussian Error Linear Units (GELUs). Preprint at
1496 <https://doi.org/10.48550/arXiv.1606.08415> (2023).
- 1497 53. Akiba, T., Sano, S., Yanase, T., Ohta, T. & Koyama, M. Optuna: A Next-generation
1498 Hyperparameter Optimization Framework. in *Proceedings of the 25th ACM SIGKDD*
1499 *International Conference on Knowledge Discovery & Data Mining* 2623–2631
1500 (Association for Computing Machinery, 2019). doi:10.1145/3292500.3330701.
- 1501 54. Rand, K., Grytten, I., Pavlovic, M., Kanduri, C. & Sandve, G. K. BioNumPy: Fast and
1502 easy analysis of biological data with Python. 2022.12.21.521373 Preprint at
1503 <https://doi.org/10.1101/2022.12.21.521373> (2022).

- 1504 55. Douze, M. *et al.* The Faiss library. Preprint at
1505 <https://doi.org/10.48550/arXiv.2401.08281> (2024).
- 1506 56. Steinegger, M. & Söding, J. MMseqs2 enables sensitive protein sequence searching
1507 for the analysis of massive data sets. *Nat Biotechnol* **35**, 1026–1028 (2017).
- 1508 57. Nayfach, S. *et al.* Metagenomic compendium of 189,680 DNA viruses from the
1509 human gut microbiome. *Nat Microbiol* **6**, 960–970 (2021).
- 1510 58. Terzian, P. *et al.* PHROG: families of prokaryotic virus proteins clustered using
1511 remote homology. *NAR Genomics and Bioinformatics* **3**, lqab067 (2021).
- 1512 59. Zhang, Y. & Skolnick, J. TM-align: a protein structure alignment algorithm based on
1513 the TM-score. *Nucleic Acids Res* **33**, 2302–2309 (2005).
- 1514 60. McInnes, L., Healy, J. & Melville, J. UMAP: Uniform Manifold Approximation and
1515 Projection for Dimension Reduction. Preprint at <http://arxiv.org/abs/1802.03426>
1516 (2020).
- 1517 61. Kipf, T. N. & Welling, M. Semi-Supervised Classification with Graph Convolutional
1518 Networks. Preprint at <http://arxiv.org/abs/1609.02907> (2017).
- 1519 62. Mihara, T. *et al.* Linking Virus Genomes with Host Taxonomy. *Viruses* **8**, 66 (2016).
- 1520 63. Morris, C. *et al.* Weisfeiler and Leman Go Neural: Higher-order Graph Neural
1521 Networks. Preprint at <http://arxiv.org/abs/1810.02244> (2021).
- 1522
- 1523



1524
1525
1526
1527
1528
1529
1530
1531
1532
1533
1534
1535
1536
1537
1538
1539

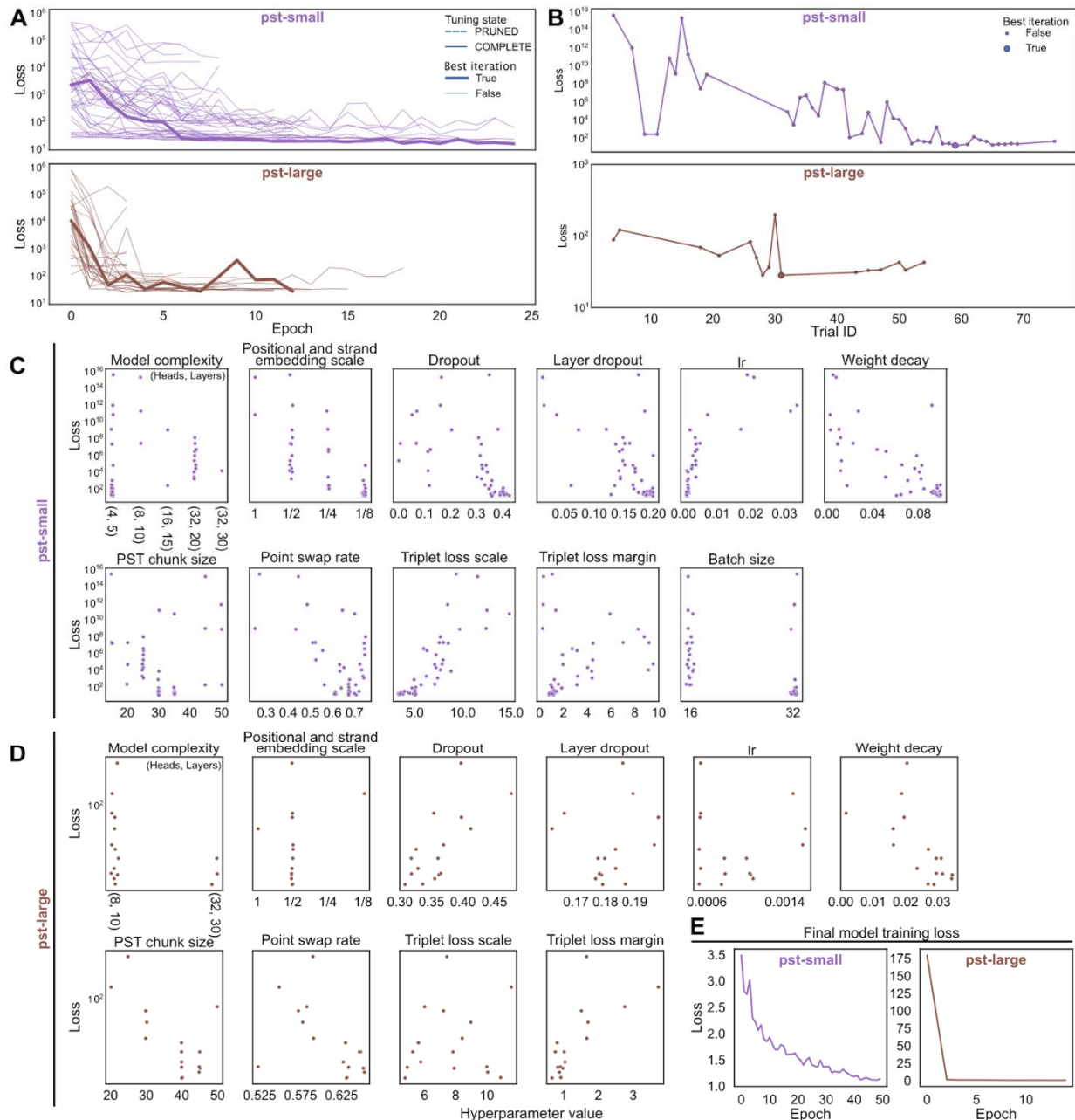
Extended Data Figure 1. A machine-learning-centric view of the Protein Set Transformer (PST) architecture. Each genome is internally represented as a graph, composed of subgraph genome chunks whose size is tunable. A minibatch of genomes is represented in a memory-efficient stacked matrix, where the boundaries for the set of proteins from each genome are tracked using indices and offset pointers for efficient random access. At the beginning of training, the ESM2 protein embeddings are concatenated with learnable positional embeddings based on the relative position in each genome and encoding strand embeddings. This is then input to the PST encoder, which uses multi-head attention for pairs of proteins defined by the initial adjacency matrix. This only allows each protein to attend to its neighbors in the same genome subgraph. The output from the PST encoder are genome-contextualized protein embeddings, which are also the inputs to the PST decoder. The PST decoder uses multi-head attention pooling to project each contextualized protein embedding onto a learnable seed vector. This learns weights for each protein, which are used to pool each protein representation into a final genome representation.



1540
1541
1542
1543

Extended Data Figure 2. The training and test datasets for the vPST. **A)** Distributions of scaled average amino acid identity (AAI) for viruses only within the training (Train), only within test (Test), and between (Across) datasets. The AAI between pairs of

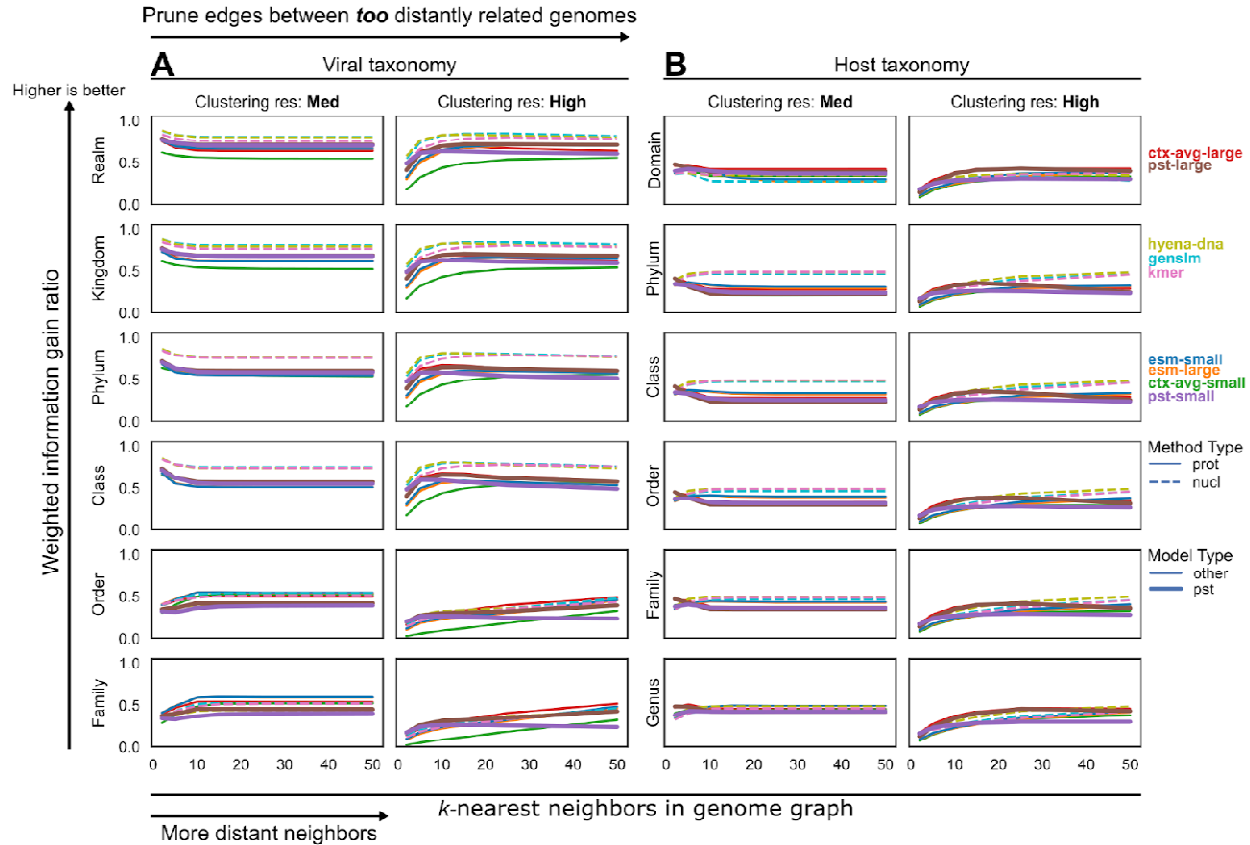
1544 viruses was scaled by multiplying the minimum proportion of proteins from each
1545 genome used in computing the AAI. Only the most similar connections are maintained
1546 for each virus in each dataset. In the case of across train-test boundaries, the most
1547 similar hit for each test virus is reported. **B)** Proportional source of 103,589 training set
1548 viral genomes (see Supplemental Table 1 for the source publication of each virus). **C)**
1549 Genome size distribution of training and test dataset genomes. **D)** The relative
1550 distributions of viral realm, host domain, and broad ecosystem for both the train and test
1551 datasets. Viral realm, if not provided by the source database, was predicted by
1552 geNomad v1.5.0. Host domains not provided by the source database, excluding
1553 predicted proviruses, were considered unknown.



1554
1555
1556
1557
1558
1559
1560
1561
1562
1563
1564
1565

Extended Data Figure 3. Hyperparameter tuning for the vPST. **A**) Triplet loss curves as a function of training step (epoch) for the small vPST (*pst-small*, top, purple) and the large vPST (*pst-large*, bottom, brown). Each line is a different tuning iteration that did not fail. Dashed lines indicate tuning trials that were pruned due to poor performance in early epochs, relative to previous trials. Solid lines indicate trials that completed training. The thicker line highlights the trial that was chosen as the best by cross validation. **B**) The final triplet loss value for each tuning trial. The larger circle indicates the trial chosen as the best by cross validation. **CD**) Hyperparameter values tuned with cross validation for the small vPST (**C**, *pst-small*, 45 complete trials) and the large vPST (**D**, *pst-large*, 16 complete trials). Model complexity refers to the number of attention heads and the number of encoder layers in the PST. Positional and strand embedding scale is

1566 the size of each of these concatenated learned embeddings relative to the input ESM2
1567 protein language model embedding for each protein. Lr is learning rate, and weight
1568 decay is for the AdamW optimizer. The PST chunk size is the number of proteins per
1569 genome chunk. Point swap rate is the proportion of proteins swapped between the
1570 anchor and positive genome during point swap sampling. Point swap scale is the
1571 negative exponential decay scale factor to adjust the weight of the choice of the
1572 augmented negative samples in the triplet loss function. Batch size for pst-small is in
1573 units of genomes and was constant at 8 genomes for pst-large. **E)** Training loss curves
1574 for the final trained models using the optimal hyperparameters selected for each model.
1575



1576

1577

1578 **Extended Data Figure 4.** Viral and host taxonomic purity of viral genome clusters.

1579 Weighted information gain ratio computed as cluster-size-weighted average over all

1580 genome clusters, penalized for the proportion of genome singletons, using **A)** viral

1581 taxonomy or **B)** host taxonomy as the genome labels. Viral taxonomy that was either

1582 not provided by source databases or outdated was predicted using geNomad v1.5.0.

1583 Unlabeled taxa were excluded during calculation of the relative proportion of taxa

1584 comprising a cluster. However, the cluster size used for weighted average included all

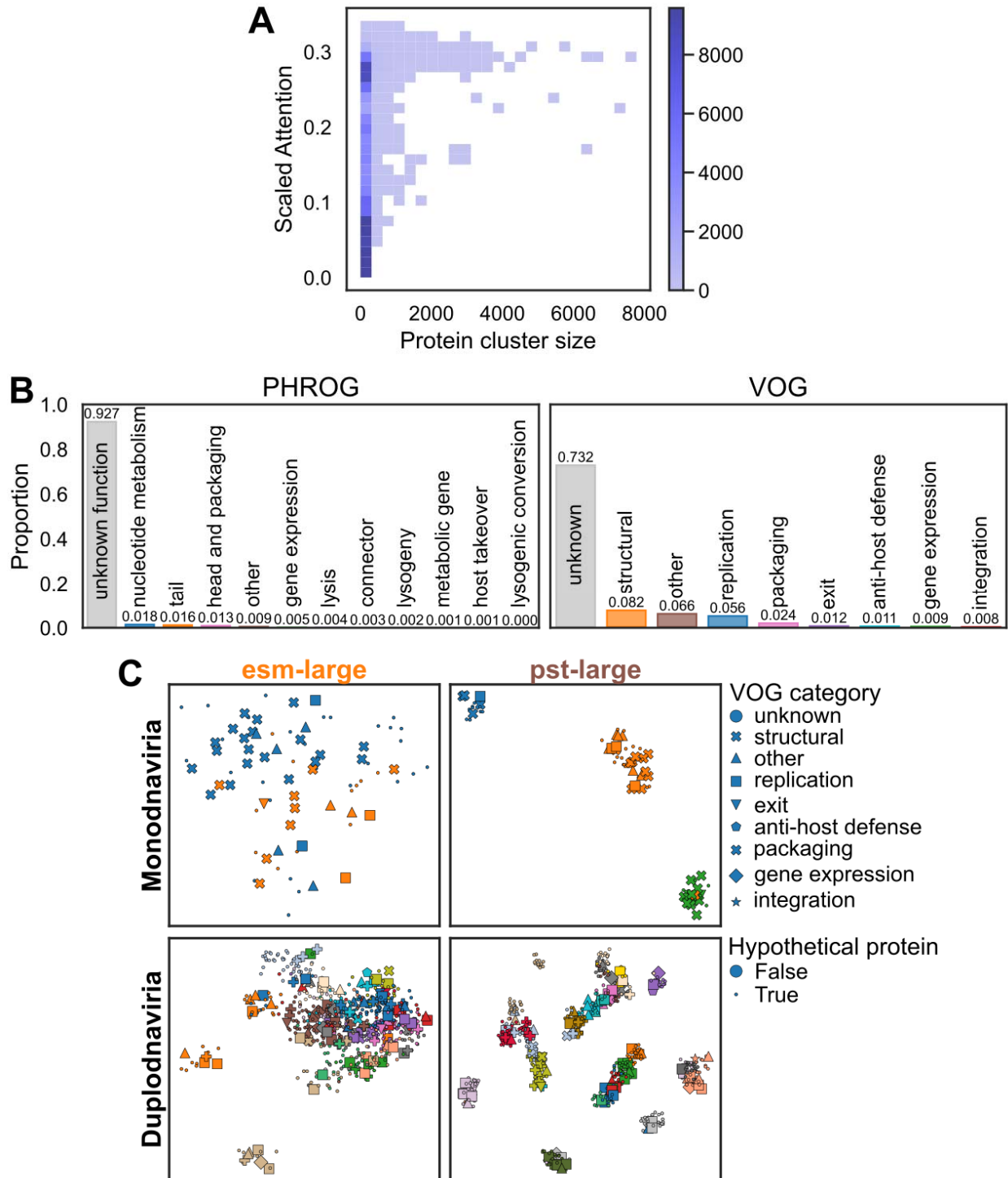
1585 viruses in the cluster. The Leiden algorithm on a similarity-weighted k -nearest neighbor

1586 graph was used to cluster genomes. The edges and edge weights (similarity values)

1587 were computed using a similarity-transformed Euclidean distance on the indicated

1588 genome embeddings (represented by the color of the lines). All analyses were

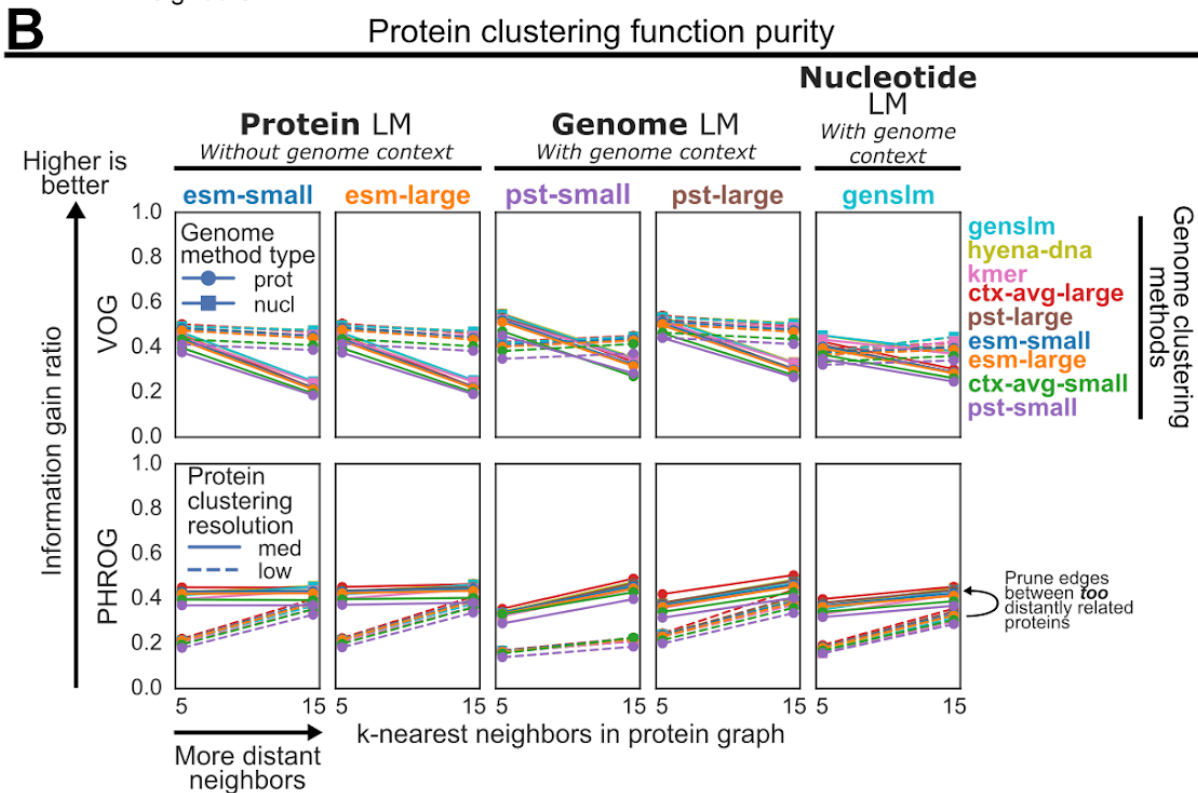
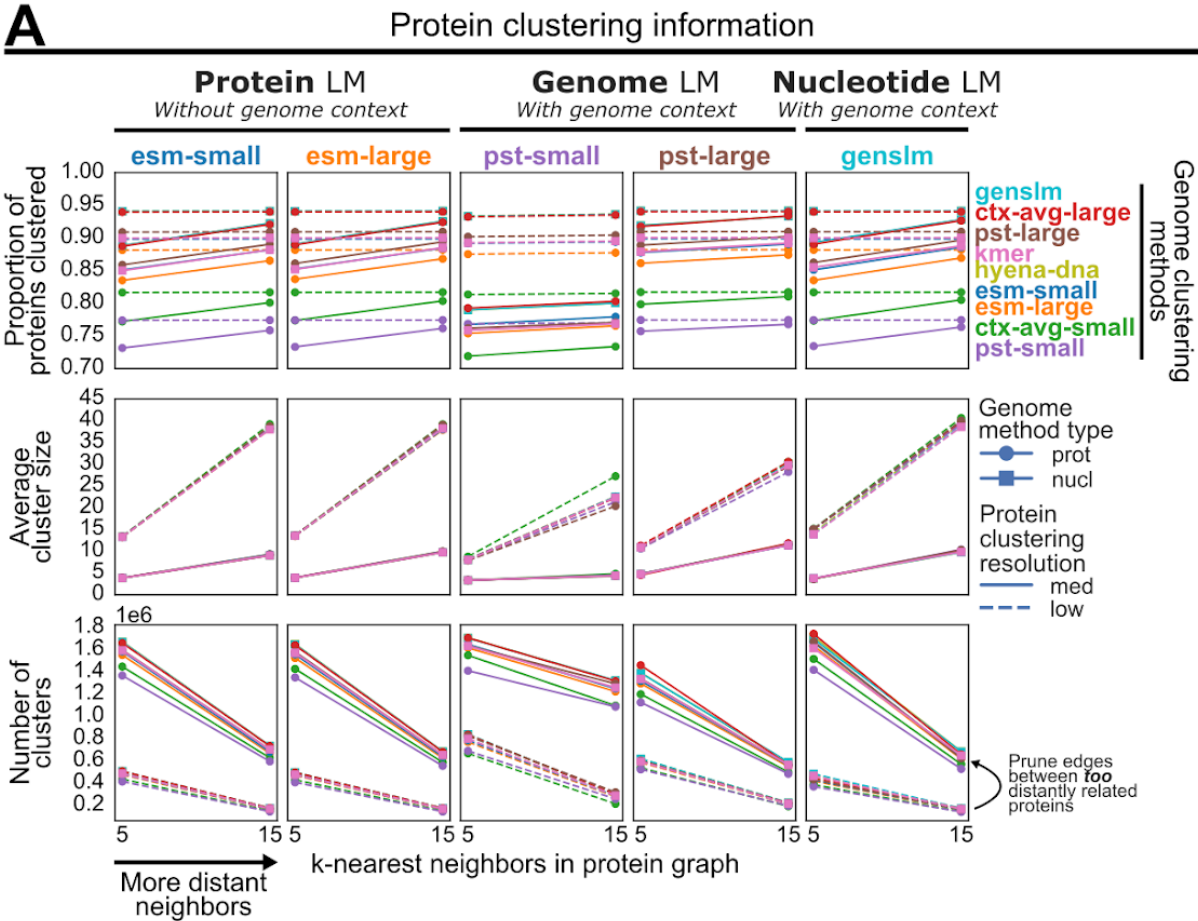
performed on the vPST test dataset.



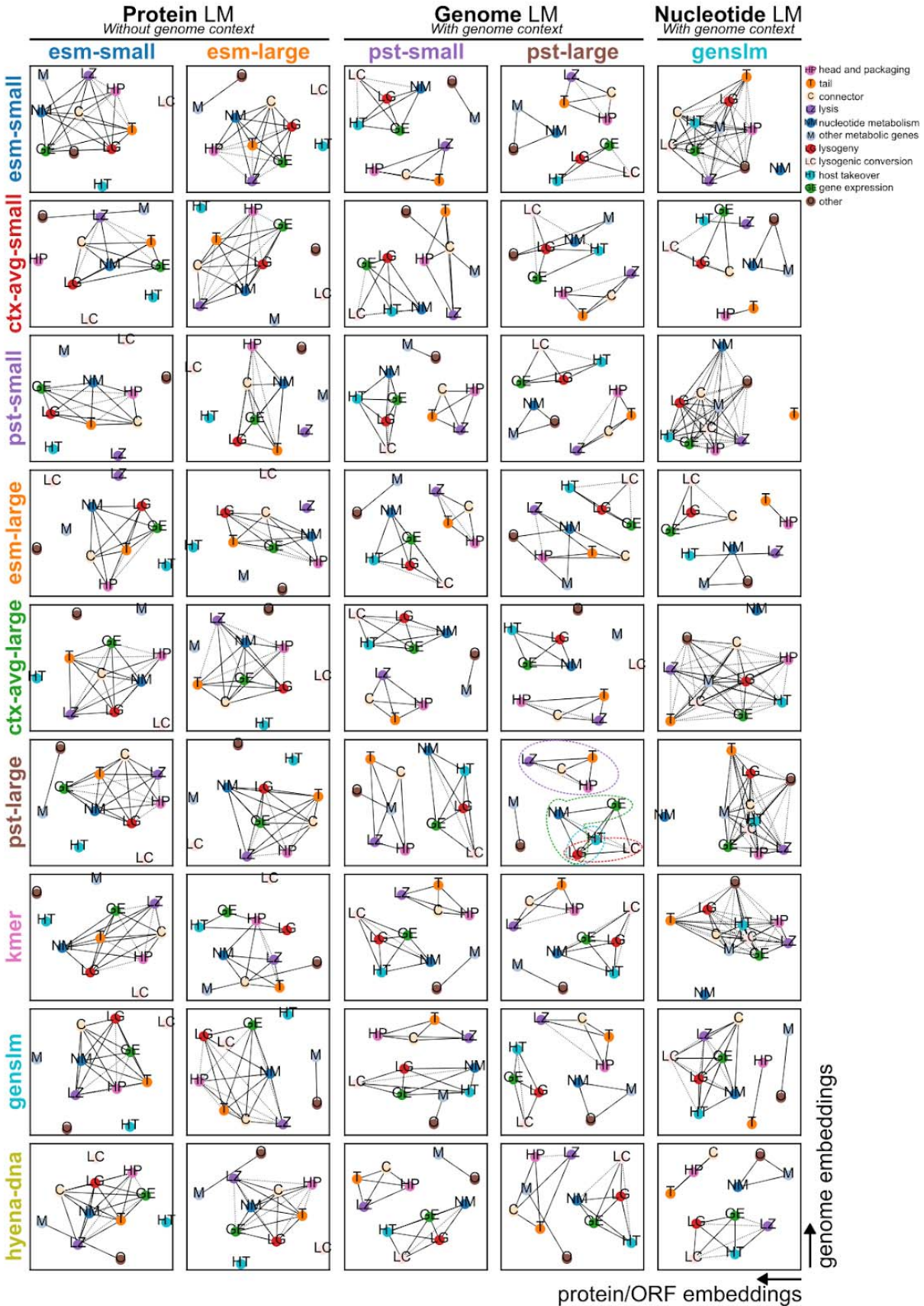
1589
1590
1591
1592
1593
1594
1595

Extended Data Figure 5. Protein function associations detected by vPST. **A)** Scaled attention for the top 50 sequence identity-based protein clusters (mmseqs2) compared to the number of proteins in the cluster in a 2D histogram. Counts greater than 10,000 were clipped for ease of visualization. **B)** Distribution of protein functional annotations using curated categories from the PHROG or VOG databases. **C)** UMAP dimensionality reduction plots for 2 genome clusters primarily composed of ($\geq 85\%$) Monodnaviria (top)

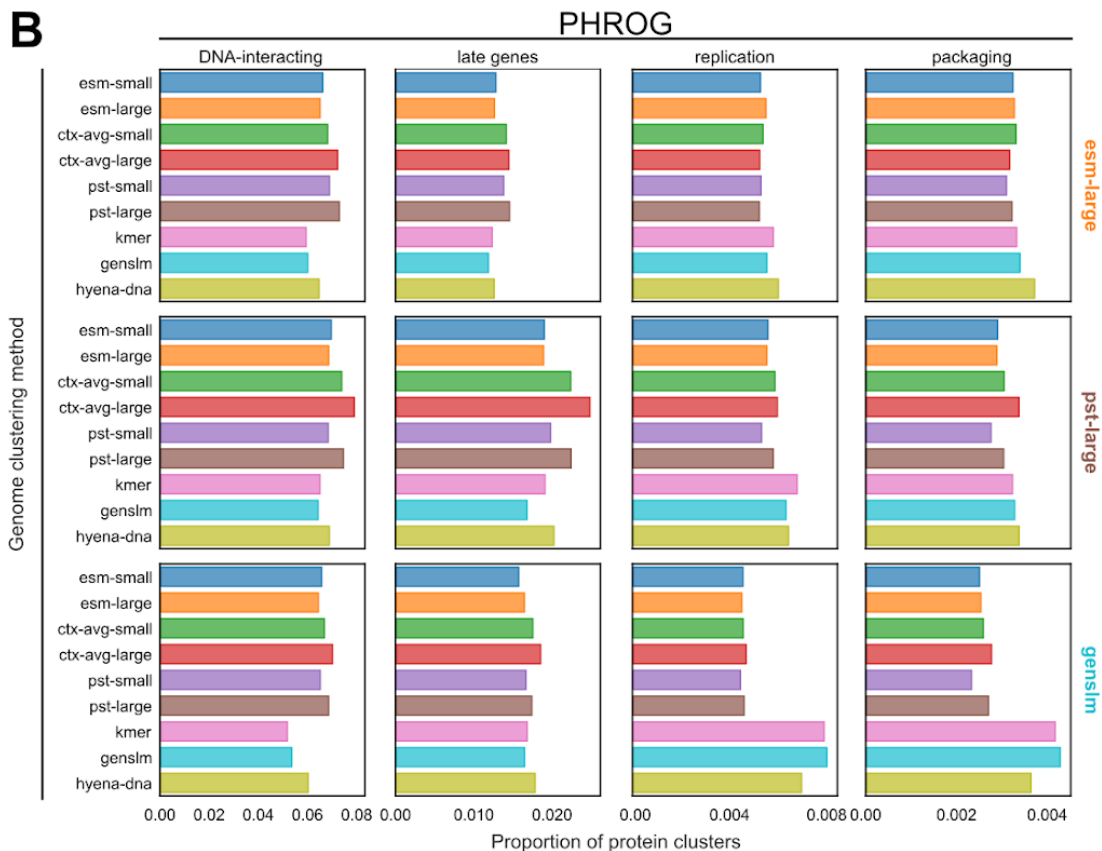
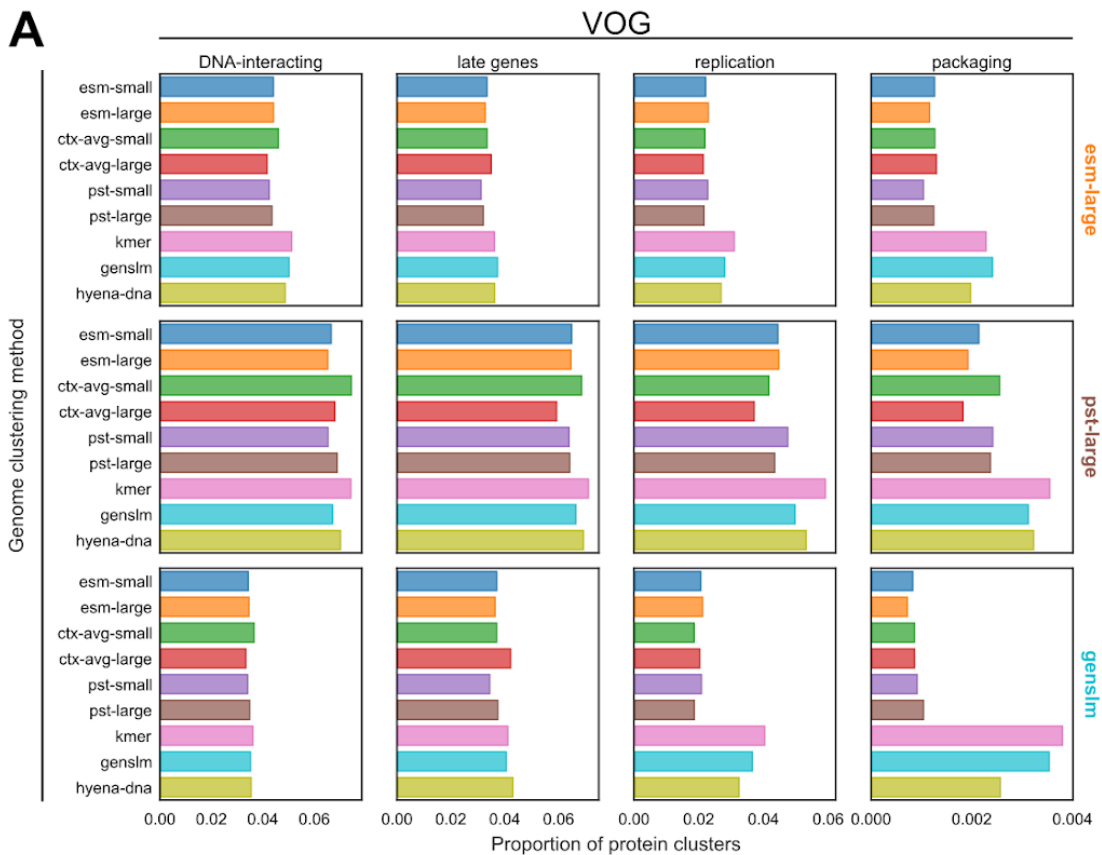
1596 or Duplodnaviria (bottom). Shapes indicate the VOG r219 curated functional category,
1597 and colors represent the protein cluster membership using the Leiden algorithm on the
1598 kNN graph. These are the same genome clusters as in Fig. 3C. All analyses were done
1599 using the proteins from the vPST test dataset.



1601 **Extended Data Figure 6.** Protein clustering information and evaluation. **A)** Statistics of
1602 embedding-based protein clusters detected by the Leiden algorithm on the similarity-
1603 weighted k -nearest neighbors graph. The edges and edge weights (similarity values)
1604 were computed using cosine similarity after length normalizing the indicated protein
1605 embedding (columns). Protein neighbors were only considered if the source genomes
1606 belonged to the same genome cluster, which were clustered using the embedding
1607 indicated by the color of the line. The genomes were clustered using the parameters
1608 that maximized intra-cluster AAI (“High” Leiden resolution, $k=15$). **B)** Weighted
1609 information gain ratio computed as cluster-size-weighted average over all protein
1610 clusters, penalized for the proportion of protein singletons, using curated function
1611 categories from VOG and PHROG as protein labels. Unannotated proteins were
1612 excluded during calculation of the relative proportion of each function category
1613 comprising a cluster. However, the cluster size used for weighted average included all
1614 proteins in the cluster. All analyses were done using the proteins from the vPST test
1615 dataset.



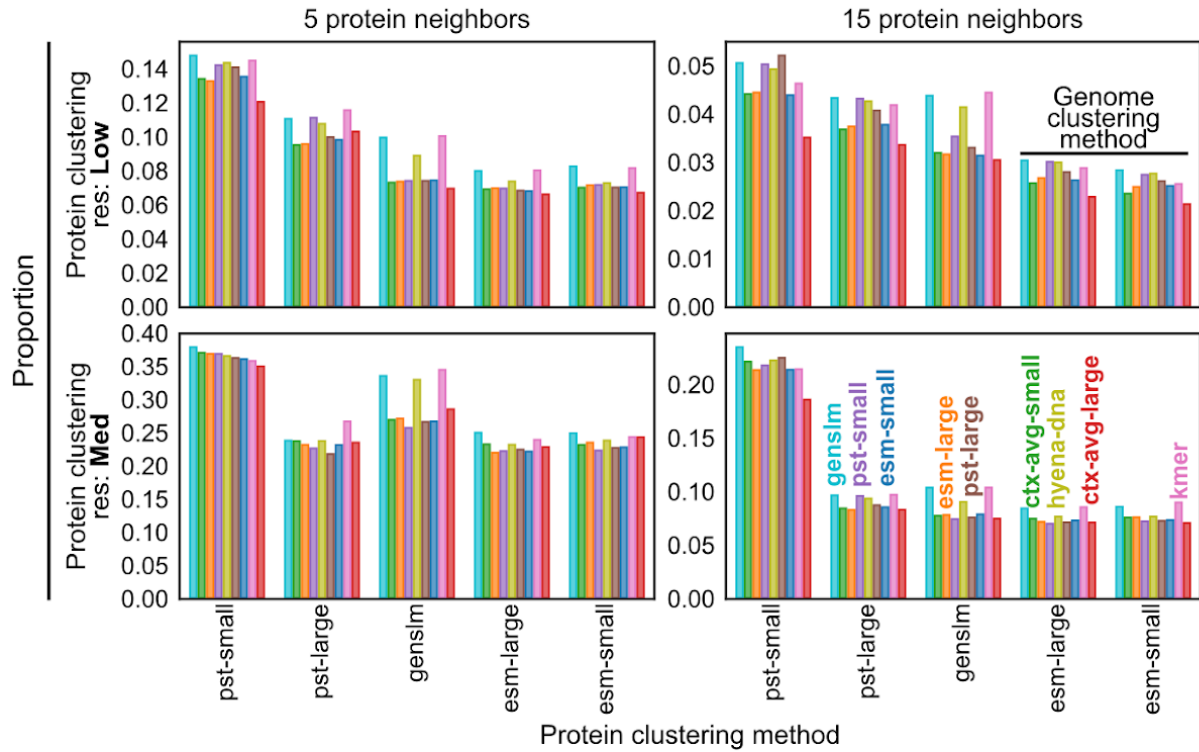
1617 **Extended Data Figure 7.** Summary of protein function co-clustering. Rows indicate the
1618 genome embedding used for genome clustering (k=15, “High” Leiden clustering
1619 resolution), and the columns indicate the protein embedding used for protein clustering
1620 (k=15, “Med” Leiden clustering resolution). The “pst-large” column refers to the
1621 intermediate contextualized protein embeddings, while “ctx-avg-large” row refers to the
1622 average of the previous protein embeddings over each genome. The functional
1623 categories are curated from the PHROG database. Each connected component was
1624 clustered in a co-occurrence graph using the Ledian algorithm with resolution of 1.0.
1625 Edges indicate functional categories that were more enriched in the protein clusters
1626 compared to the joint occurrence of these categories in the PHROG database. The
1627 length of the edges reflects the degree of enrichment. Dotted lines indicate connections
1628 that were less enriched than expected. This is a blow up of **Fig. 3C**. All analyses were
1629 done using the proteins from the vPST test dataset.



1631 **Extended Data Figure 8.** The viral Protein Set Transformer (vPST) protein embeddings
1632 detects functional modules more frequently than other protein clustering methods. The
1633 proportion of protein clusters that correspond to 1 of 4 functional modules (columns)
1634 using either **A)** VOG or **B)** PHROG annotations. Proteins were clustered within each
1635 genome cluster with k=15 and clustering resolution="med" using the protein embedding
1636 indicated by the row. Genomes were clustered with k=15 and clustering
1637 resolution="high" using the genome embedding indicated by the y-axis and color. The
1638 data used here were averaged over the genome clustering methods for **Fig. 3D**. All
1639 analyses were done using the proteins from the vPST test dataset.

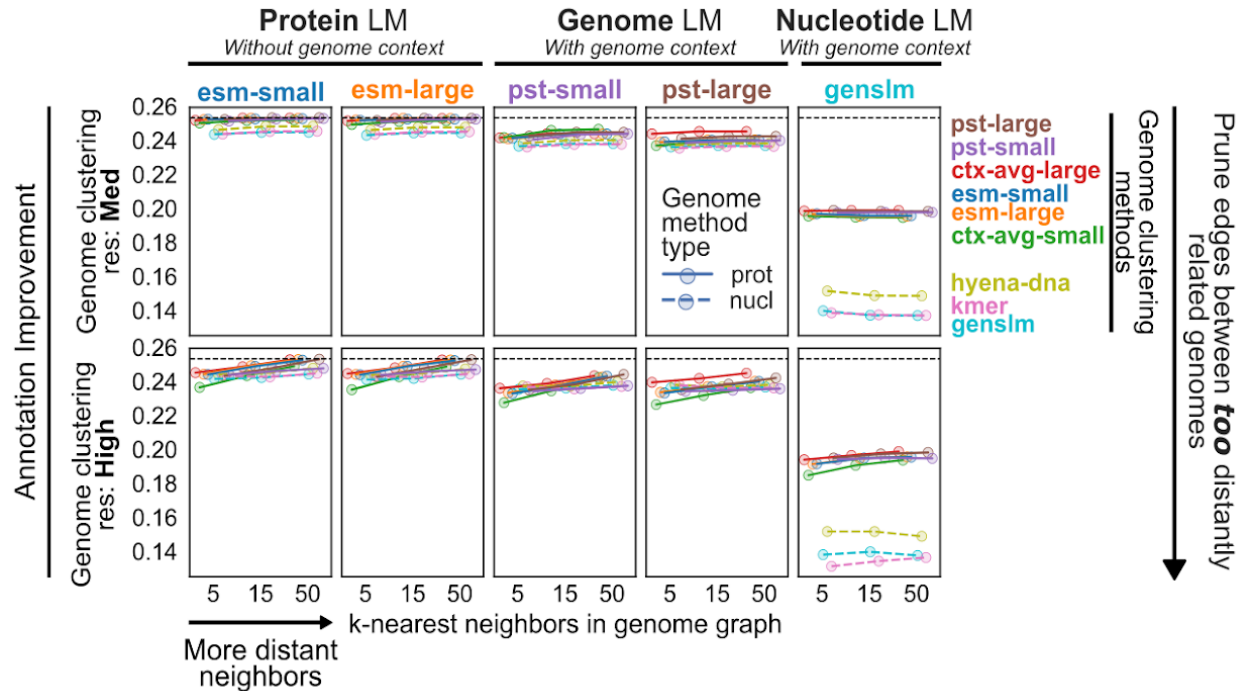
A

Identifying capsid proteins based on structure



B

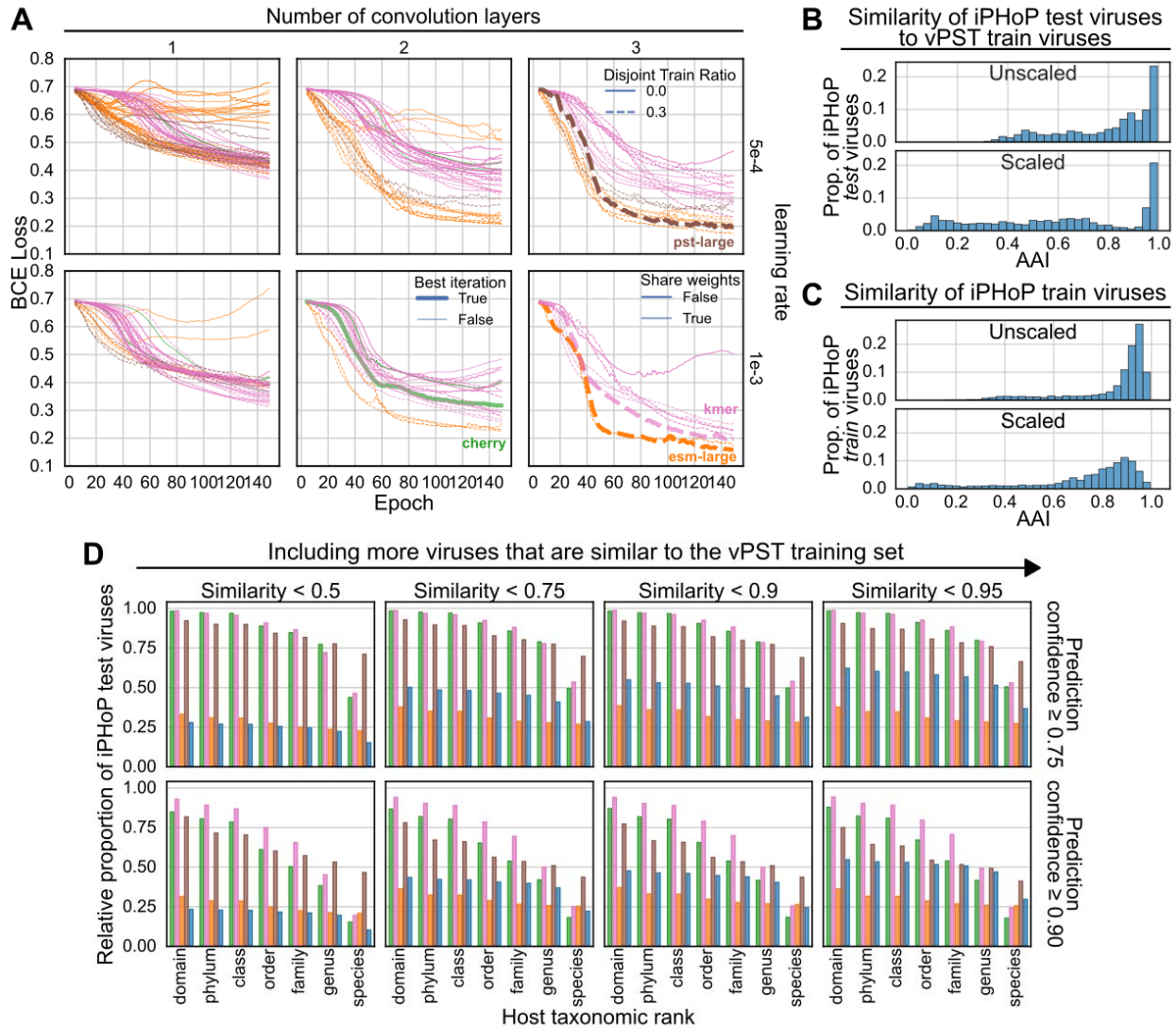
Functional annotation transfer to proteins of unknown function



1640
1641
1642

Extended Data Figure 9. Expanding our view of unannotated hypothetical proteins. **A)** The average proportion of proteins unannotated by VOG clustering with annotated

1643 capsid proteins that have structural homology to known capsid folds. Structural
1644 homology was detected using foldseek searching against the Protein Data Bank
1645 database. Values are only comparable within each subpanel. **B)** Annotation
1646 improvements for proteins unable to be assigned a function by remote homology to the
1647 VOG database. Annotation improvement is defined as the fraction of unannotated
1648 proteins whose closest neighbor (cosine distance) has an assigned function. For both
1649 panels, protein clustering or similarity searching was only allowed for proteins belonging
1650 to genomes in the same genome cluster (“High” Leiden resolution, k=15). The color of
1651 the bars (**A**) and lines (**B**) indicates the embedding used for genome clustering.



1652
1653
1654
1655
1656
1657
1658
1659
1660
1661
1662
1663
1664
1665
1666
1667
1668

Extended Data Figure 10. Comparing host prediction tools on the iPHoP test dataset. **A)** Training curves for graph-based host prediction models. Performance of different iterations of host prediction models evaluated using binary cross entropy (BCE) loss on the same validation dataset. Disjoint train ratio refers to the proportion of edges used for inference only, with the remaining edges used only for message passing. In the case of 0.0, all training edges were used for both message passing and inference. The opacity of the lines indicate if the model weights were shared between the virus-host and virus-virus edges during message passing in the same layer. The best iterations for each input node (genome) embedding type are indicated by the thick lines and were chosen as the global minimum loss at the final epoch. The additional lines are for different choices of decoder fully-connected layer hidden dimensions. Each model was trained using the same training dataset of 3,639 viruses as iPhOP. **B)** The maximum average amino acid identity (AAI) for each iPHoP test virus when searching against the vPST training dataset. “Unscaled” is the raw AAI, while “Scaled” weights the AAI by the minimum proportion of proteins used to compute AAI for each pair of genomes. **C)** Similar to **B**, except only computing AAI internally among the iPHoP training dataset. **D)**

1669 The proportion of iPHoP test viruses whose true host taxonomic rank is predicted with
1670 confidence \geq the indicated threshold (rows). Based on **B**, we excluded viruses from this
1671 set if they were more similar based on scaled AAI to the vPST training dataset
1672 (columns).