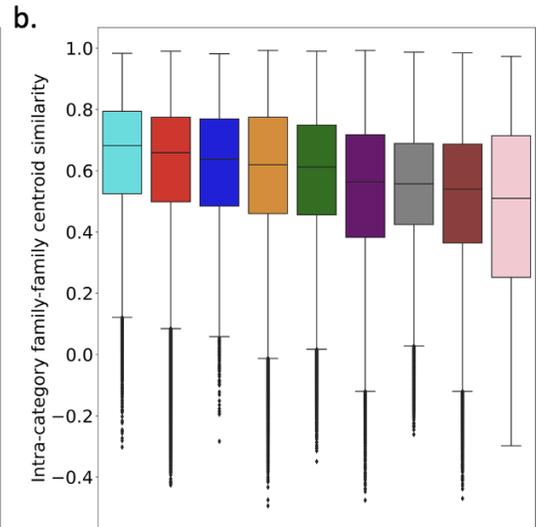
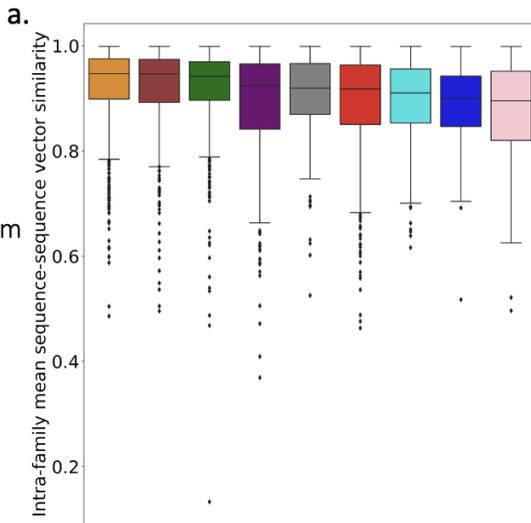


Supplemental Figure 1: Performance of four different PLM-based representations for viral VPF functional classification. Embedded proteins were used to train and evaluate PHROGs functional annotation classification. Performance is across five-fold training-testing splits of PHROGs VPFs. Study is described by the model architecture, protein source, and whether the PLM is trained with a multi-task training objective (MT). Performance is measured as F1-score.

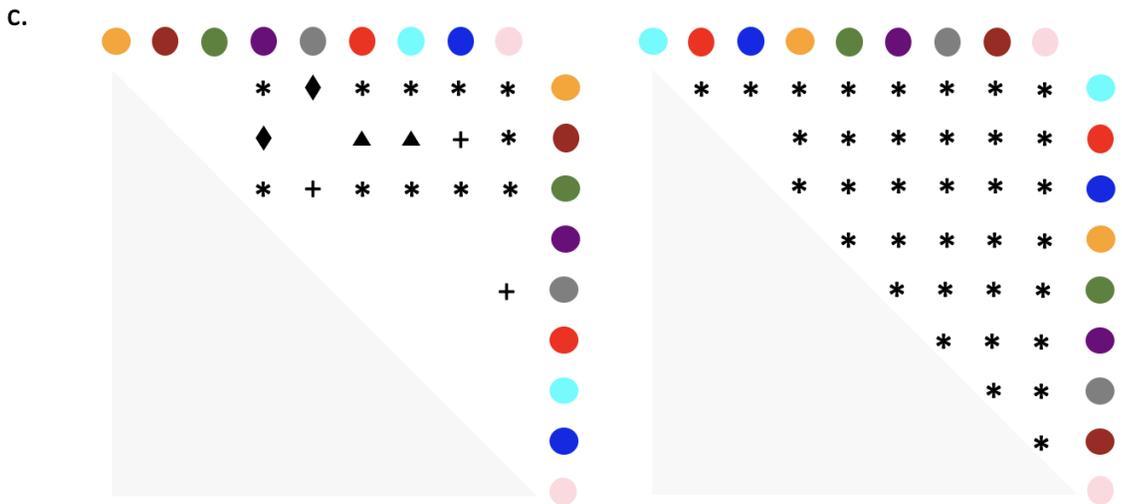
Category

- tail
- head and packaging
- connector
- lysis
- DNA, RNA, and nucleotide metabolism
- transcription regulation
- moron, auxiliary metabolic gene and host takeover
- integration and excision
- other

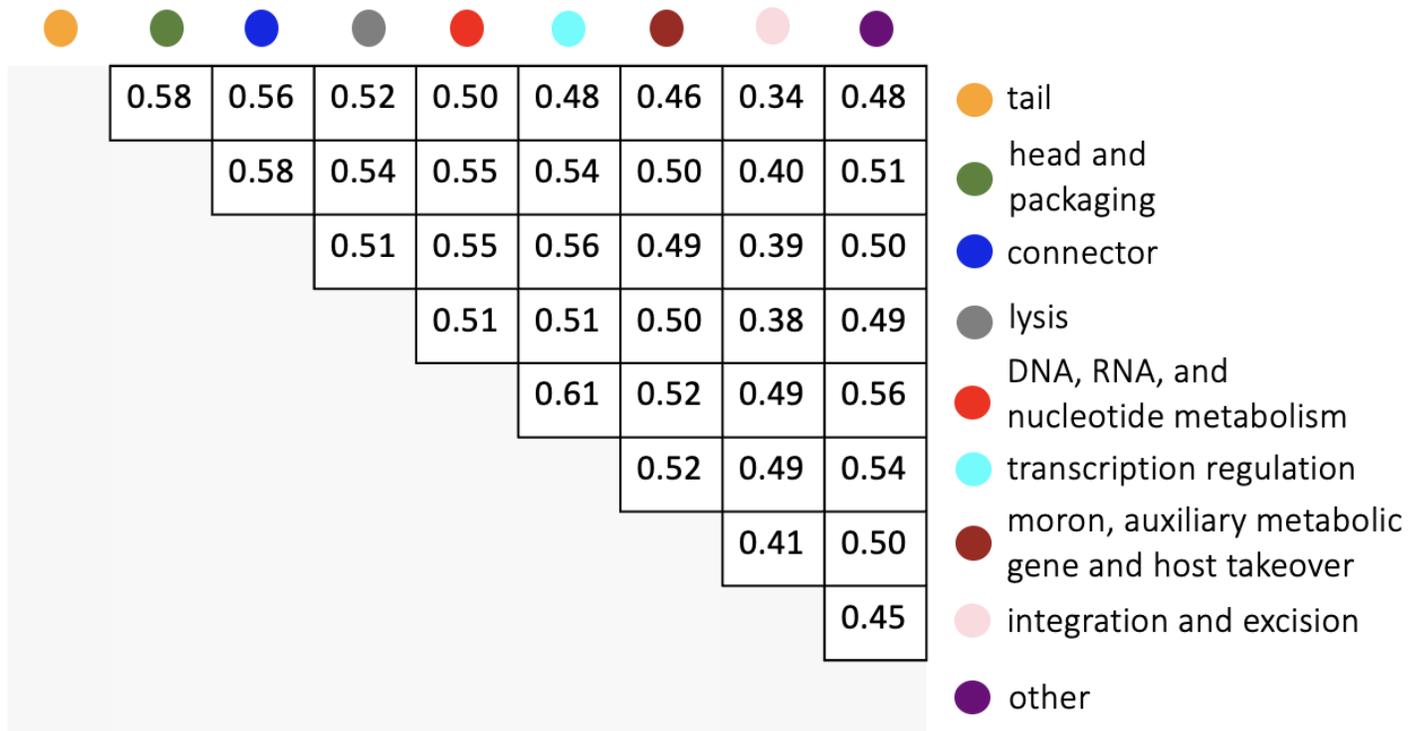


Significance

- + p < 0.01
- ▲ p < 0.001
- ◆ p < 0.0001
- * p < 0.00001



Supplemental Figure 2: Evaluation of functional in embedding similarities of constituent families between functional categories. (a) Distribution of family average sequence-sequence similarity. (b) Distribution of family-family centroid similarity. (c) Significance of pairwise category distribution comparison using independent t-test with Bonferroni correction.



Supplemental Figure 3: Inter-category similarity for PHROGs functional categories. Pairwise family centroid similarities were calculated for every combination of families between the two categories. Score is the average over all comparisons.

a.

```

Model_01:F ELSGGAAKESLTEKQYQNLEFKSHNRENQR 264
Model_01:G ELSGGAAKESLTEKQYQNLEFKSHNRENQR 264
Model_01:M ELSGGAAKESLTEKQYQNLEFKSHNRENQR 264
Model_01:N ELSGGAAKESLTEKQYQNLEFKSHNRENQR 264
1z1b.1.F SSSP----- 196

Model_01:F NWDLDYDALYLEWFFFLRSENFPHLKIEWF 294
Model_01:G NWDLDYDALYLEWFFFLRSENFPHLKIEWF 294
Model_01:M NWDLDYDALYLEWFFFLRSENFPHLKIEWF 294
Model_01:N NWDLDYDALYLEWFFFLRSENFPHLKIEWF 294
1z1b.1.F --(CWLRLAMELAVVT)GQR(VGDD)CEM(KNS)E 224

Model_01:F KDIGDETIVCRLKTKGNRDIQEVN---- 320
Model_01:G KDIGDETIVCRLKTKGNRDIQEVN---- 320
Model_01:M KDIGDETIVCRLKTKGNRDIQEVN---- 320
Model_01:N KDIGDETIVCRLKTKGNRDIQEVN---- 320
1z1b.1.F D--DGLYLV--DQSKTKGVKIAI)TALH)D 249

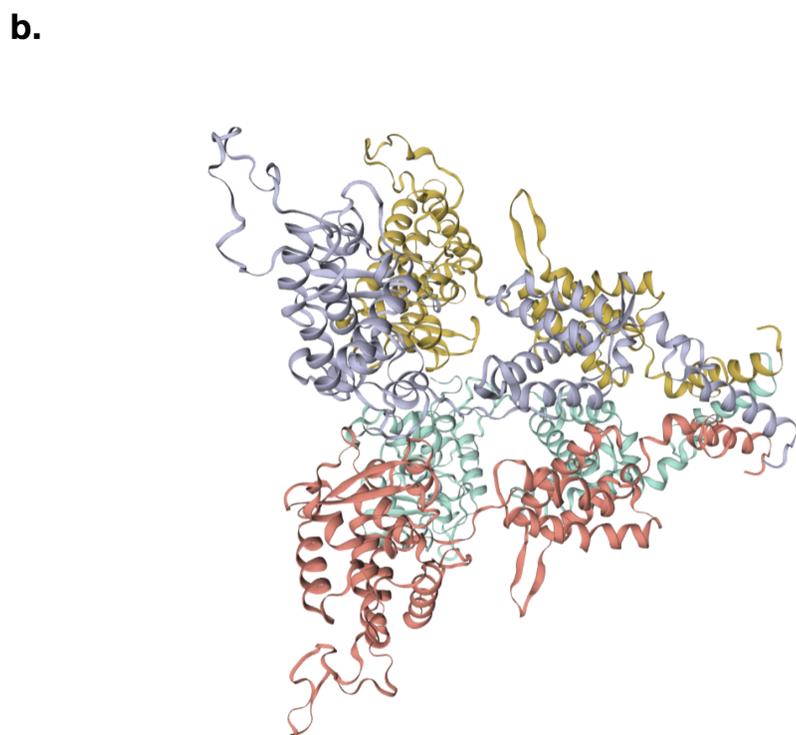
Model_01:F YRPDASKFWKRLSKRRGKEGYLVLPHIKR 349
Model_01:G YRPDASKFWKRLSKRRGKEGYLVLPHIKR 349
Model_01:M YRPDASKFWKRLSKRRGKEGYLVLPHIKR 349
Model_01:N YRPDASKFWKRLSKRRGKEGYLVLPHIKR 349
1z1b.1.F ALG(L)KKE(LDKCKE-I)GG(E)I IAST-RR 277

Model_01:F QEEGGPEAKVVRTLNFLLKSALEKSFNPVD 379
Model_01:G QEEGGPEAKVVRTLNFLLKSALEKSFNPVD 379
Model_01:M QEEGGPEAKVVRTLNFLLKSALEKSFNPVD 379
Model_01:N QEEGGPEAKVVRTLNFLLKSALEKSFNPVD 379
1z1b.1.F EP--LSS(GT)VSRYFMRRARKASGLSFE--GD 303

Model_01:F SELTWTNLRHTALRLTLEEMPELGIPPLIN 409
Model_01:G SELTWTNLRHTALRLTLEEMPELGIPPLIN 409
Model_01:M SELTWTNLRHTALRLTLEEMPELGIPPLIN 409
Model_01:N SELTWTNLRHTALRLTLEEMPELGIPPLIN 409
1z1b.1.F -PPTFH(ELR)SL(SAR-LYEK---)IS(DK-- 325

Model_01:F DFAAN-AGTSPKMLHETYLKYIQQGVTSKR 438
Model_01:G DFAAN-AGTSPKMLHETYLKYIQQGVTSKR 438
Model_01:M DFAAN-AGTSPKMLHETYLKYIQQGVTSKR 438
Model_01:N DFAAN-AGTSPKMLHETYLKYIQQGVTSKR 438
1z1b.1.F FAQHL(D)GHKSDTM(A-SQY)RDDRGRE----- 349

```



Supplementary Figure 4: Comparative protein structure modelling of a novel integrase sequence supports annotation as a tyrosine recombinase. (a) Target/template alignment between the *Prochlorococcus* PAC1 sequence (indicated as Model_01), and the template sequence 1Z1B, the phage lambda integrase. Active site residues Arg 212, Lys 235, His 308, Arg 311, His 333, and Tyr 342, are indicated with red arrows. (b) Homology model of *Prochlorococcus* PAC1 sequence based on template 1Z1B. Colors indicate individual monomers of the homo-tetramer template protein structure in both a and b.

Category	Number of families	Total proteins
connector	133	34,054
DNA, RNA, and nucleotide metabolism	1,065	98,104
head and packaging	946	90,584
integration and excision	105	16,688
lysis	299	30,997
moron, auxiliary metabolic gene and host takeover	458	22,566
other	560	37,349
tail	1,219	102,955
transcription regulation	303	39,386
unknown function	33,792	395,697

Supplemental Table 1: Viral protein families (VPFs) and total number of proteins manually annotated to 10 function categories in the Prokaryotic virus Remote Homologous Groups (PHROGs) database.

Project	Model name	Model architecture	Protein source	Mutl-task training objective	Number of proteins	Training objective
ProSE ³⁰	MT-LSTM	LSTM	Unifref90	yes	76,215,872	Masked token + predicting contact between residues in protein structure + structural similarity of proteins by SCOP hierarchy
ProSE ³⁰	DLM-LSTM	LSTM	Unifref90	no	76,215,872	Masked token
ProteinBERT ²⁹	ProteinBERT	Transformer	Uniref90	yes	~106M	Masked token + GO term annotation
ProtTrans ²⁸	protbert_bfd	Transformer	BFD	no	2,122M	Masked token

Supplemental Table 2: PLMs evaluated for viral protein VPF functional classification. Four PLMs were used to embed PHROGs VPF proteins for training and testing the viral function classifier. Model name is name from original study. PLMs were chosen to vary in the model architecture, protein source, and use of multiple training objectives.

Category	Support	Precision	Recall	F1-score
connector	2,145	0.61	0.62	0.62
DNA, RNA, and nucleotide metabolism	18,972	0.83	0.90	0.86
head and packaging	14,520	0.75	0.77	0.76
integration and excision	402	0.94	0.89	0.92
lysis	2,046	0.60	0.62	0.61
moron, auxiliary metabolic gene and host takeover	4,097	0.63	0.68	0.65
other	7,087	0.65	0.60	0.63
tail	11,060	0.61	0.90	0.72
transcriptional regulation	535	0.58	0.71	0.64
unknown function	20,078	0.71	0.45	0.55
weighted average	80,939	0.71	0.71	0.70

Supplemental Table 4: Functional classifier performance for EFAM VPFs labeled with PHROGs annotation. PHROGs functional annotation was assigned to EFAM VPFs using HMM matching. Support indicates how many EFAM VPFs matched PHROGs VPFs per category.