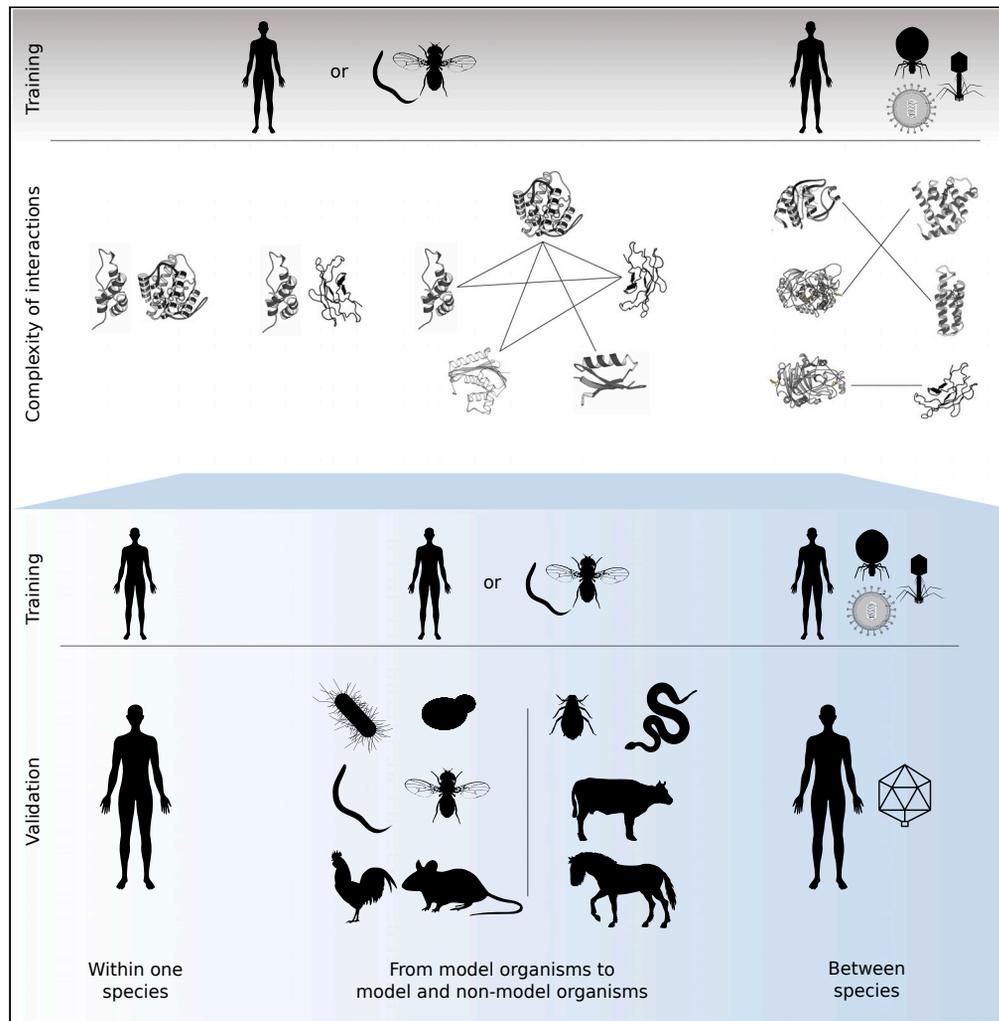


Article

SENSE-PPI reconstructs interactomes within, across, and between species at the genome scale



Konstantin Volzhenin, Lucie Bittner, Alessandra Carbone

alessandra.carbone@lip6.fr

Highlights

SENSE-PPI enables PPI reconstruction at the genome level from sequences

10,000 proteins can be screened against themselves in a matter of hours

It finds interactions in non-model organisms, disordered proteins, human-virus pairs

Specific functional modules are sharply detected

Volzhenin et al., iScience 27, 110371
July 19, 2024 © 2024 The Author(s). Published by Elsevier Inc.
<https://doi.org/10.1016/j.isci.2024.110371>



Article

SENSE-PPI reconstructs interactomes within, across, and between species at the genome scale

Konstantin Volzhenin,¹ Lucie Bittner,^{2,3} and Alessandra Carbone^{1,3,4,*}

SUMMARY

***Ab initio* computational reconstructions of protein-protein interaction (PPI) networks will provide invaluable insights into cellular systems, enabling the discovery of novel molecular interactions and elucidating biological mechanisms within and between organisms. Leveraging the latest generation protein language models and recurrent neural networks, we present SENSE-PPI, a sequence-based deep learning model that efficiently reconstructs *ab initio* PPIs, distinguishing partners among tens of thousands of proteins and identifying specific interactions within functionally similar proteins. SENSE-PPI demonstrates high accuracy, limited training requirements, and versatility in cross-species predictions, even with non-model organisms and human-virus interactions. Its performance decreases for phylogenetically more distant model and non-model organisms, but signal alteration is very slow. In this regard, it demonstrates the important role of parameters in protein language models. SENSE-PPI is very fast and can test 10,000 proteins against themselves in a matter of hours, enabling the reconstruction of genome-wide proteomes.**

INTRODUCTION

Protein-protein interaction (PPI) networks play a key role in biology and medicine in the interpretation of protein functions in cellular processes. In the past two decades, working with networks has significantly advanced our understanding of the relationships between molecules.^{1–6} This was possible thanks to many computational attempts,^{7–20} among which are recent deep learning architectures^{21–24} and high-throughput experimental methods such as yeast two-hybrid^{25,26} or tandem purification^{27,28} that have been extensively developed. A particular concern is their level of noise and incompleteness.^{29–31} In addition, various studies give insights into the precise spatial organization and dynamic temporal remodeling of local protein interaction networks within the cell,^{32,33} highlighting that PPI networks are only “projections” of particular spatiotemporal PPI realizations. For instance, within each individual, genomic alterations contribute to different PPI realizations. Accordingly, understanding any biological process demands defining three parameters: the composition of the “underlying protein network,” its organization in space, and its evolution over time. Future technological developments will bring an overwhelming amount of precise information on these genomics-spatio-temporal dimensions, and sophisticated computational tools for extracting information from them are mandatory. Ultimately, PPI networks should be understood within biological frameworks including transcriptomic and epigenomic data. Here, we address the first step of this development, which is the construction of the “underlying protein network,” which will serve as the basis for more sophisticated and realistic reconstructions. Because of the difficulties explained above, which are intrinsic to experimental data, *ab initio* computational reconstruction of PPIs is supposed to provide invaluable information, leading to the identification of protein partners.

Today, the problem of PPI network reconstruction is becoming particularly important due to the wealth of sequence data available from different species and the need to understand proteomes within and between species, which most of the time are non-model organisms, that is species that cannot grow in the laboratory and have a long life cycle, low fecundity, or poor genetics for instance. At the same time, the identification of protein partners in PPIs puts new deep learning approaches to the test, for discovering sophisticated correlations within pairs of interacting protein sequences and for estimating the absence of such correlations when the interaction does not exist in the cell. We take up this challenge and propose SENSE-PPI, an *ab initio* deep learning approach for protein partner identification, coupling layers of gated recurrent units (GRU)³⁴ with the ESM2 protein language model (PLM) encoding sequences,³⁵ to achieve optimal identification for protein partners spanning a large spectrum of biological functions and leading to the reconstruction of PPI networks. Indeed, the complexity of the problem may vary depending on the features and the origins of proteins. SENSE-PPI has undergone extensive training and testing across a variety of datasets to assess its performance and generalizability: (1) the human dataset STRING11.0,²³ used for initial training and validation.

¹Sorbonne Université, CNRS, IBPS, UMR 7238, Laboratoire de Biologie Computationnelle et Quantitative (LCQB), 75005 Paris, France

²Institut de Systématique, Evolution, Biodiversité (ISYEB), Muséum national d'Histoire naturelle, CNRS, Sorbonne Université, EPHE, Université des Antilles, Paris, France

³Institut Universitaire de France, Paris, France

⁴Lead contact

*Correspondence: alessandra.carbone@lip6.fr

<https://doi.org/10.1016/j.isci.2024.110371>



(2) A second human dataset derived from the STRING database, which incorporates the “neighboring-exclusion” condition (see [STAR Methods](#) section), to simulate more challenging and realistic interaction scenarios. (3) Datasets from diverse species such as fly, worm, yeast, mouse, chicken, and bacteria. These were employed to test SENSE-PPI when trained on the human proteome alone or combined with other species like fly, worm, or chicken. This approach helps explore how evolutionary distances between the training data and test organisms impact the model’s performance. (4) Datasets from non-model organisms including horse, cow, snake, and aphid, used to further evaluate the model’s robustness and applicability to a wider range of biological contexts. (5) The human-virus dataset, specifically designed to study PPIs between species, providing insights into cross-species interaction dynamics. (6) The IDPpi dataset,³⁶ which focuses on interactions involving human intrinsically disordered proteins, offering a specialized perspective on this challenging category of proteins. (7) A reference yeast dataset that allows for comparative analysis against other existing models (can be found in the [supplemental information](#)). Each dataset serves a specific purpose: the first two are used to examine the effects of different dataset construction methods on the results, with the second also serving as the primary training dataset for testing the model under various conditions. This comprehensive testing framework enables us not only to evaluate the generalizability of SENSE-PPI but also to delineate its limitations across a broad spectrum of biological conditions.

RESULTS

We leverage the power of deep learning to identify correlations in interacting sequences and to estimate their absence when there is no interaction. Our architecture, SENSE-PPI, exploits the latest generation of PLMs, ESM2, and recurrent neural networks, GRU, to predict the interaction of protein pairs with different features on a large scale. Validation is carried out using a wide range of interaction datasets of different complexities, from model and non-model organisms, inter-species interactions, human-virus interactions, and interactions with intrinsically disordered human proteins. A large new dataset of interactions in *Homo sapiens*, comprising over 1 million interactions, is designed based on a restrictive “neighboring-exclusion” condition that guarantees more accurate predictions, a property that is particularly required when laboratory tests are planned. SENSE-PPI is compared with state-of-the-art deep learning approaches, PIPR,²¹ D-SCRIPT,²³ Topsy-Turvy,²² and STEP³⁷ (see [supplemental information](#)), in various scenarios.

Training on the human proteome and testing on model and non-model organisms

Recent global efforts to sequence the biodiversity of species^{38–42} make PPI predictions in yet unexplored organisms a major challenge. Here, we test SENSE-PPI on non-model species and assess its generalization capabilities for predicting interactions in model organisms that are phylogenetically distant, to varying degrees, from the species used for training. The SENSE-PPI model was trained on *Homo sapiens* and tested on PPI data from several different species: *Homo sapiens*, *Mus musculus*, *Drosophila melanogaster*, *Caenorhabditis elegans*, *Saccharomyces cerevisiae*, *Escherichia coli*, *Gallus gallus*, *Equus caballus*, *Bos taurus*, *Notechis scutatus*, and *Acyrtosiphon pisum*. The first seven organisms are model organisms but the last four are not.

Benchmark testing for model organisms

To compare the performance of SENSE-PPI with state-of-the-art approaches on model organisms, we repeated the training and testing procedures described in ref.²³ SENSE-PPI was trained on the large STRING11.0 human dataset, specially designed to avoid protein redundancies in training and validation sets. [Table 1](#) presents the comparison of SENSE-PPI performance with those of PIPR, D-SCRIPT, and Topsy-Turvy trained on the same human data. The overall performance remains high for all four DL models when testing is carried out within *H. sapiens*. However, scores decline progressively when tests are carried out on the other four model species. As the species’ evolutionary distance from *H. sapiens* increases, performance decreases progressively: the most accurate predictions are obtained for *M. musculus* and the worst for *S. cerevisiae*. SENSE-PPI outperforms all other DL methods on all model species tested and, more importantly, offers greater generalization ability in several respects. First, the ability of SENSE-PPI to distinguish between positive and negative interactions is measured with the AUROC score (see [STAR Methods](#)), which remains above 0.9 for all test sets, while D-SCRIPT goes from an AUROC score of 0.833 for *H. sapiens* down to 0.790 for *S. cerevisiae*. Second, SENSE-PPI’s performance results in an F1-score (see [STAR Methods](#)) of 0.848 within *H. sapiens* that remains above 0.737 for *M. musculus*, *D. melanogaster*, and *C. elegans*. This result shows that PPI networks of model species whose common ancestor is estimated at 708 million years from *H. sapiens*⁴³ remain well predicted by SENSE-PPI trained on *H. sapiens*. The SENSE-PPI F1-score decreases to 0.555 for *S. cerevisiae*, sharing a common ancestor with *H. sapiens* at 1,300 million years. It should be noted, however, that the performance of SENSE-PPI on this very distant species is superior in terms of F1-score to that of PIPR, with 0.555 versus 0.456 on *M. musculus*, and D-SCRIPT, with 0.402 on *H. sapiens*. Overall, SENSE-PPI significantly outperforms previously developed methods on both single-species and cross-species tests.

New human training dataset improves accuracy

Motivated by the design of a fair training dataset on which to evaluate SENSE-PPI, we constructed the STRING11.5_neighbor_exclusion human dataset from interactions in the STRING database v11.5. The set of negative protein pairs was defined by satisfying a stringent “neighboring-exclusion” condition, ensuring that A-B is a negative interaction in the training set if (1) A and B are not known to interact in STRING, (2) no homolog of A is known to interact with a homolog of B at more than 40% sequence identity in STRING, and (3) no homolog, at more than 40% sequence identity, of a known interactor of B is known to interact with a homolog of A. Note that the first two conditions are also used in

Table 1. Evaluation of SENSE-PPI trained on the STRING11.0 human dataset

Species	Model	AUPRC	AUROC	F1-Score
<i>M. musculus</i>	PIPR	0.526	0.839	0.456
	D-SCRIPT	0.663	0.901	–
	Topsy-Turvy	0.735	0.934	–
	SENSE-PPI	0.859	0.973	0.782
<i>D. melanogaster</i>	PIPR	0.278	0.728	0.196
	D-SCRIPT	0.605	0.890	–
	Topsy-Turvy	0.713	0.921	–
	SENSE-PPI	0.847	0.969	0.742
<i>C. elegans</i>	PIPR	0.346	0.757	0.235
	D-SCRIPT	0.550	0.853	–
	Topsy-Turvy	0.700	0.906	–
	SENSE-PPI	0.821	0.963	0.737
<i>S. cerevisiae</i>	PIPR	0.230	0.718	0.140
	D-SCRIPT	0.399	0.790	–
	Topsy-Turvy	0.534	0.850	–
	SENSE-PPI	0.657	0.914	0.555
<i>H. sapiens</i>	PIPR	0.835	0.960	0.763
	D-SCRIPT	0.516	0.833	0.402
	Topsy-Turvy	0.703	0.895	0.711
	SENSE-PPI	0.917	0.984	0.848

Testing was carried out on species not present in the training data: mouse, fly, worm, and yeast (top). Testing on a reserved fraction of the human dataset is also shown (bottom). The best values for every specific metric and dataset are highlighted in bold. Scores for PIPR, D-SCRIPT, and Topsy-Turvy in the first four test sets were taken from ref.²² The scores for PIPR and D-SCRIPT in the human test set were taken from ref.²³ The evaluation of Topsy-Turvy on human data was recomputed.

the construction of the STRING11.0 human dataset. This dataset is double the size of the STRING11.0 human dataset and contains more than 86,000 positive and 860,000 negative pairs. We have trained SENSE-PPI on this new dataset and have compared SENSE-PPI performance on the four STRING11.0 test sets from four different species.²³ Table 2 shows the F1-score, precision, and recall. Even though the STRING11.5_neighbor_exclusion human dataset leads to better performance (F1-score) in three cases out of four, the increase remains on par. For practical usage, it is important to note that the trained models provide different precision/recall ratios. In particular, the SENSE-PPI human dataset obtains higher values of precision (trading off some recall instead) and might be more valuable when lab experiments should be conducted and target partners identified. It will produce fewer false positives, thus increasing the proportion of relevant interactions among all positive predictions.

Testing across the evolutionary tree on model and non-model species

Furthermore, we asked whether SENSE-PPI behavior is preserved across the phylogenetic tree and, in particular, for non-model species. For this, we constructed STRING11.5_neighbor_exclusion datasets for the four non-model and six model species (see STAR Methods). A coherent behavior of SENSE-PPI across species and time evolution would support the use of SENSE-PPI on species at a fixed evolutionary distance from the training one for which little is known about protein interactions. Consistent with our expectations, the behavior of SENSE-PPI is illustrated in Figure 1 (see also Table S1).

The evaluation of SENSE-PPI on 10 testing datasets across species provides valuable insights, although it does not offer a comprehensive assessment of performance. Hence, to tackle the redundancy problem inherent in the input data seen in training,^{10,44} we systematically filtered the proteins in the test sets according to different levels of sequence identity from proteins in the training set and evaluated SENSE-PPI precision accordingly. Previous attempts to analyze the effects of the proximity of training and testing sets mostly used the concepts of C1/C2/C3 classes.^{10,44} These classes divide the testing set into three subgroups based on how many proteins in a pair were already present in training (two, one, or none, respectively). In this work, however, we focus on how performance varies with increasing evolutionary distance between the species tested and the species in the training data, without relying on this discrete division.

For this, different extents of sequence identity were measured with MMseqs2⁴⁵ by searching for the closest matches (i.e., consecutive k-mer matches, based on sequence identity without gaps) between testing and training sequences. For each protein pair in the testing set, we computed a so-called “mean pair sequence identity”: first, we computed the maximum sequence identities for both proteins in a pair with respect to all proteins in the training set, and then, we considered the mean of these two values (see STAR Methods). Figure 1A

Table 2. Evaluation of SENSE-PPI trained on two human interaction datasets

Species	Training dataset	F1-score	Precision	Recall
<i>M. musculus</i>	D-SCRIPT	0.782	0.755	0.811
	SENSE-PPI	0.802	0.896	0.727
<i>D. melanogaster</i>	D-SCRIPT	0.742	0.658	0.850
	SENSE-PPI	0.763	0.814	0.718
<i>C. elegans</i>	D-SCRIPT	0.737	0.692	0.788
	SENSE-PPI	0.701	0.815	0.614
<i>S. cerevisiae</i>	D-SCRIPT	0.555	0.437	0.760
	SENSE-PPI	0.603	0.623	0.585

Training was carried out on the STRING11.5_neighbor_exclusion and STRING11.0 datasets. Testing was carried out on the four model species in ref.²³ For each species and each metric, the best performances are highlighted in bold.

shows the distribution of MCC scores with respect to the mean pair sequence identity for 10 different species. We observe a clear positive correlation: the larger the mean pair sequence identity is, the better the performance. This agrees well with the evolutionary distance between species (Figure 1B): species that are phylogenetically closer to the training data (*H. sapiens*) tend to have higher scores. However, one can observe a tendency of non-model organisms to have slightly lower scores: Figures 1C and 1D show the MCC vs. pair sequence identity plots for both model and non-model organisms. Here, instead of simply taking the mean pair sequence identity for all test entries, the data were divided into bins of size 0.1. The two plots illustrate that the primary difference in evaluating model versus non-model organisms lies in the lower values of mean pair sequence identity. Specifically, the average MCC remains between 0.7 and 0.8 for identity values greater than 0.8 in both scenarios. However, when the identity drops below 0.1, the difference doubles: 0.3 for model organisms compared to 0.15 for non-model organisms. Except for *E. coli*—which is evolutionarily the most distantly related organism to the training data—the MCC scores for model organisms do not fall below 0.3. Conversely, non-model organisms only begin to achieve MCC scores above 0.3 when the average identity of sequence pairs exceeds 40%.

Evolutionary close training data improve performance

To demonstrate the sensitivity of SENSE-PPI to training data including species closer to the target organism, we focused on the non-model organism *A. pisum* and trained the system using a dataset enhanced with species evolutionarily closer to aphids than to humans. This enhanced dataset combines the original STRING11.5_neighbor_exclusion human dataset with additional datasets for *D. melanogaster* and *C. elegans*. These species share a common ancestor with the aphid at approximately 362 million years and 572 million years ago, respectively, in contrast to their common ancestor with *H. sapiens* at around 708 million years ago. As expected, the MCC test scores for *A. pisum* increased from 0.660 to 0.722 (Table S2). Notably, even though the overall quality of classification has increased, the most substantial improvements were observed in species that were initially quite distant from the original training set, such as *A. pisum*, *S. cerevisiae*, and *E. coli*.

Further testing involved augmenting the training dataset with PPIs from an additional model organism, *G. gallus*, to explore whether incorporating data from more distantly related species would also enhance results. This expansion of the training set confirmed that the most significant improvements in classification accuracy occur in species more closely related to the newly added data. For example, the most notable score improvement was observed in *E. caballus*.

These findings suggest a general guideline for making predictions with SENSE-PPI: incorporating data from evolutionary closer species tends to yield more accurate results. Additionally, the observed enhancement in prediction precision is not paralleled by a substantial increase in the mean pair sequence identity of the test data, as demonstrated in Figure S1. This, coupled with the observed upward shift in the interpolation curves for the new datasets, indicates that the model's improvements may be attributed less to the transfer of homology and more to the expansion and diversification of the training dataset.

Larger protein language models improve performance

PLMs significantly contribute to enhanced performance. To see this, we consider two sizes of ESM2 embedding modules and analyze how MCC scores change at different values of the mean pair sequence identity for the four STRING11.0 datasets of *M. musculus*, *D. melanogaster*, *C. elegans*, and *S. cerevisiae* (Figures 1E and 1F). Figure 1E shows the SENSE-PPI performance obtained with the smaller version of ESM2, based on 35M parameters (with 1M more parameters used in the rest of the SENSE-PPI architecture), and Figure 1F demonstrates the performance of the SENSE-PPI regular version, based on 3B parameters (with 2M more parameters for the rest of the SENSE-PPI architecture). The ESM2 model with 3B parameters highly improves the predictions on pairs of proteins with low values of mean pair sequence identity compared to the smaller version of ESM2. For instance, for a pair sequence identity <0.2, predictions on the fly display an MCC of 0.2 (Figure 1E) versus 0.4 (Figure 1F). Moreover, for entries with high values of mean pair sequence identity, the performance based on both embeddings remains comparable. This remains true for all four datasets and shows that ESM2 based on a higher number of parameters has better generalization capability.

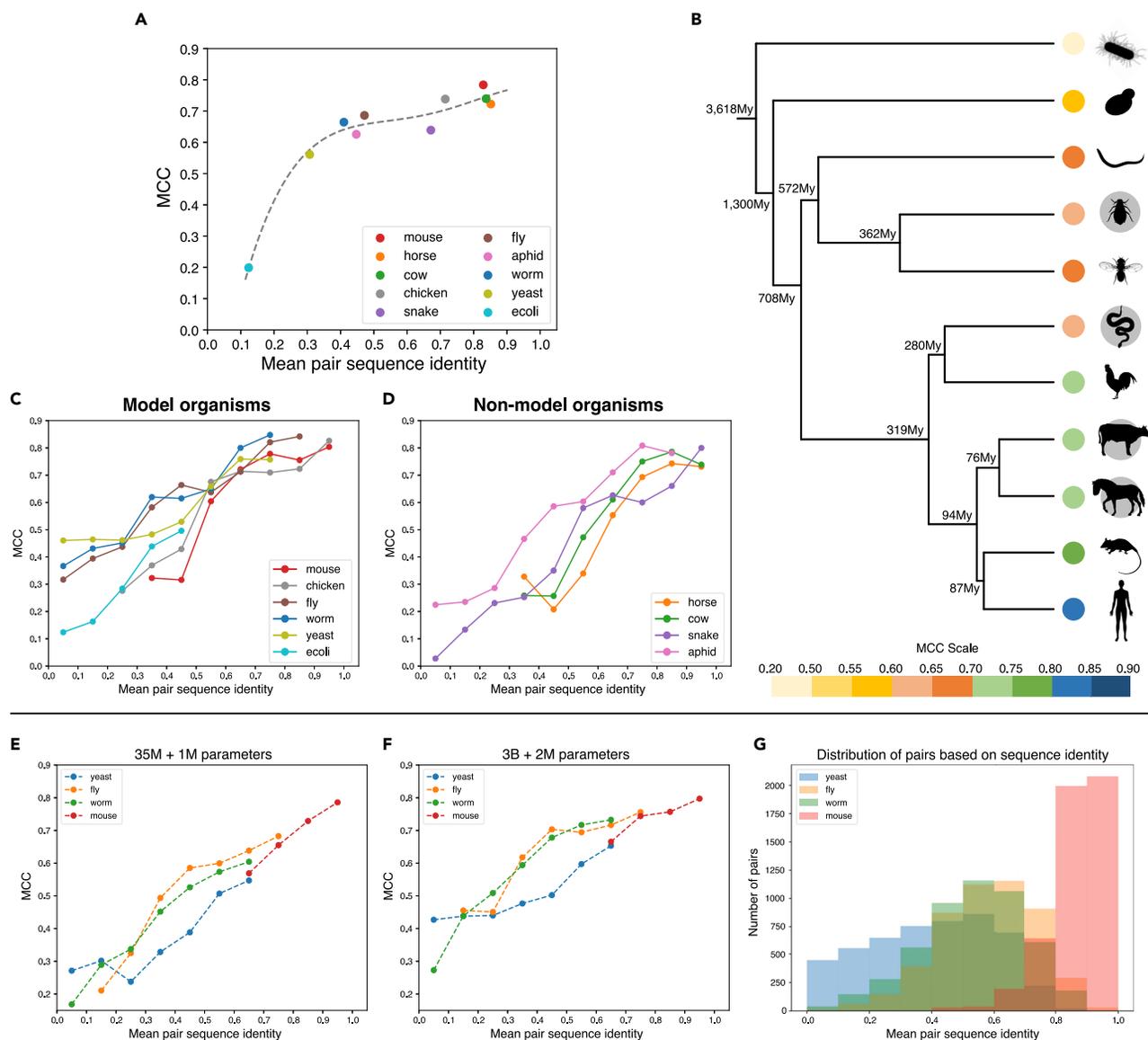


Figure 1. Performance of SENSE-PPI on six model and four non-model organisms

SENSE-PPI is trained on the STRING11.5_neighbor_exclusion human dataset. For each experiment, the testing dataset comprises 5,000 positive and 50,000 negative interactions.

(A) Matthews correlation coefficient (MCC) versus mean pair sequence identity (see STAR Methods) for 10 experiments. Dotted line: least squares polynomial fit with degree 4.

(B) Phylogenetic tree of species used for SENSE-PPI testing and performance evaluation. The color scale, labeling the species in the tree, corresponds to the MCC value obtained in testing (as reported in A); colors go from blue (high MCC) to yellow (low). Non-model species are shown with a gray disk behind them. The tree was reconstructed using phyloT v2 (<http://phyloT.biobyte.de/>). Evolutionary information on species divergence times was obtained with TimeTree.⁴³ F1 and MCC scores are reported in Table S1. The MCC score for *H. sapiens* trained on *H. sapiens* is 0.836.

(C and D) Plots of MCC scores versus mean pair sequence identity for model (C) and non-model (D) organisms.

(E and F) Comparison of two versions of SENSE-PPI differing by the size of the ESM2 embedding module. The trend of SENSE-PPI MCC scores versus the mean pair sequence identity is shown for four different testing datasets coming from²³ for the “light” version (E) and the regular version (F).

(G) Distribution of pairs in four test sets for (E)–(F) based on the mean pair sequence identity. The data on (C)–(G) are presented in bins of 10% across the values of mean pair sequence identity. The MCC on (C)–(F) is computed for a given bin only if it contains at least 40 positive and 400 negative samples.

Evaluation on the human intrinsically disordered protein dataset

Approximately 33% of eukaryotic proteins have significant disordered regions, with an increasing occurrence of disorder in higher organisms,⁴⁶ and predicting interactions between proteins that are possibly disordered becomes an urgent demand.^{47–50} Deep learning

approaches based on sequences are expected to open the way to this challenge. SENSE-PPI was additionally compared to the IDPpi model³⁶ designed to predict the interaction between proteins with stable structure and intrinsically disordered proteins (IDPs). Tables S3 and S4 show the results for both models. The models were tested on the same five test sets used in³⁶ but were trained differently. Whereas the IDPpi training set included the IDPs in the test data, SENSE-PPI was trained on the STRING11.0 human dataset, which contained no information on IDPs. The same performances obtained by the two approaches suggest that by providing more extensive training (note that the STRING11.0 dataset is 20 times larger than the IDPpi training dataset), SENSE-PPI can obtain relatively accurate predictions on IDP interactions. This test once again highlights the generalization capabilities of SENSE-PPI, enabling us to extend its use to process IDPs with fairly good accuracy.

Human-virus PPIs

Predicting and understanding virus-host PPIs is important for developing new therapeutic interventions, but knowledge of virus-host interactions covered by current databases, such as VirHostNet,⁵¹ is limited. We therefore tested the ability of SENSE-PPI to predict interactions between a known human protein and an unknown viral protein. We trained SENSE-PPI on a human-virus PPI dataset²⁴ that contains interactions between human and viral proteins, i.e., only cross-species interactions are included. Data were split between training and testing, so that test data only included interactions for Epstein-Barr and influenza viruses. The model was trained on the set with a large variety of viral species interacting with human proteins (for dataset details see STAR Methods). We required that all human proteins involved in the test have been seen during training (this dataset corresponds to class "C2" in ref.¹⁰), a condition that is reasonable in terms of potential applications. The results of this test are presented in Tables S5 and S6. The overall high scores may be explained by the fact that the model is already familiar with all the human proteins present in the data.¹⁰ However, it should be noted that the viral proteins from training and testing are drastically different: by searching the closest matches for every test viral protein in the training set (using MMseqs2, see STAR Methods), we computed a mean sequence identity of 0.21 for the Epstein-Barr virus and of 0 for the Influenza viruses (no matches were found even with sensitivity parameter values in the range of 7–10⁴⁵). Nonetheless, the protocol can be potentially useful for further exploration of interactions between different species, where the model is familiar with all the target proteins of one of the organisms of interest. The problem of predicting cross-species interactions between two unknown proteins remains wide open.

PPI network reconstruction

In order to verify the performance of SENSE-PPI on the *ab initio* reconstruction of PPI networks, we carried out several tests to reconstruct PPI networks for sets of a few dozen proteins from the STRING database (version 11.5).⁵² For each test, we set one or two proteins as "seeds" and collected known partners on the basis of high confidence for physical interaction according to STRING. We then calculated SENSE-PPI predictions for a comparative analysis with STRING data. Figures 2 and 3 show four examples of these tests chosen to illustrate the general characteristics of SENSE-PPI behavior: (1) false positives are frequently related to indirect interactions, that is non-physical interactions with proteins sharing a common partner in the network (Figure 2A); (2) SENSE-PPI often misses interactions for weakly connected proteins (Figure 2B); (3) it successfully identifies interactions involved in functionally distinguished subnetworks (Figures 2A and 3); (4) it successfully predicts PPIs in other species (Figure 2C). These four sets group interacting human proteins (Figures 2A and 2B), *C. elegans* proteins (Figure 2C), and *M. musculus* proteins (Figure 3). The tests were carried out on the model trained on the STRING11.0 human dataset.

Three sets of proteins are analyzed in Figure 2 where heatmaps display confidence scores for STRING (lower triangle) and prediction scores for SENSE-PPI (upper triangle). Although we used the same color spectrum for both sources, there is no direct correlation between the SENSE-PPI scale and the STRING scale, since our model was trained on binary labels, and STRING provides confidence estimates that were not taken into account during SENSE-PPI training (except for the fact that we only used STRING high-scoring interactions in the training set). In fact, interactions used for training or validation were not used for testing (see black cells in heatmaps). An intuitive representation of the heatmap data is illustrated by three associated graphs.

The first protein set (Figure 2A) is seeded on C1R and RFC5 proteins from *H. sapiens* and encompasses 12 other human proteins meant to interact with at least one seed protein. Proteins for which the majority of positive interactions in a given subset were already present during training have been omitted. C1R and RFC5 were chosen based on their subcellular localization: C1R is secreted in extracellular regions, whereas RFC5 is expressed in the nucleus. Moreover, these proteins possess different functional characteristics: while C1R is the first component of the classical pathway of the complement system within the immune system, RFC5 takes part in DNA replication in the cell cycle. SENSE-PPI clearly distinguishes the two functionally independent networks without any uncertainty, assigning very low scores to all protein pairs across the subnetworks. Some false-positive interactions are present and strictly concentrated in the RFC5 complex. They concern the PCNA protein, which interacts, at different strengths in STRING, with RFC5 and many, but not all, of its partners. For this set, SENSE-PPI inference of false positives is often observed for STRING highly connected proteins. In such cases, reconstructions are not completely false, as a considerable proportion of errors are made on protein pairs that may not interact physically, but are either functionally linked or homologous to actual physical partners.

The second set is seeded on the human TUFM protein and 12 other proteins (Figure 2B), known to belong to the TUFM complex as described in STRING with the highest scores. TUFM is an elongation factor localized in the mitochondrial outer membrane. SENSE-PPI shows an excellent reconstruction ability for this complex, with very high scores provided for the interactions even though multiple proteins have not been seen in training. However, a weakness is observed for proteins with a low degree of interaction in the network, such as TSFM, ATG5, and

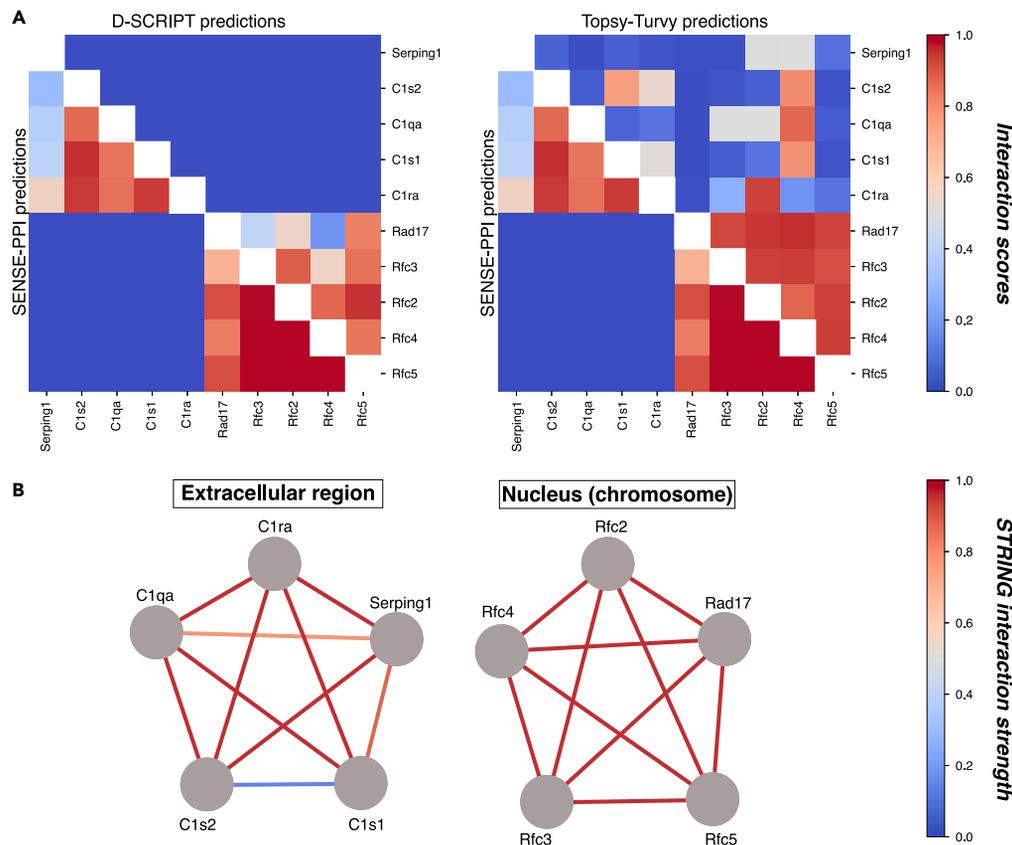


Figure 3. PPI network reconstruction seeded on the protein complexes C1-Serping1 and RFC-Rad17 in *M. musculus*

(A) SENSE-PPI trained on the STRING11.0 human dataset was used to reconstruct two networks with independent functionality in the cell. None of the proteins in the sets were observed during training, hence the full set of predicted interactions was considered in the evaluation. The heatmaps show SENSE-PPI (lower triangle), D-SERIP (upper triangle, left), and Topsy-Turvy (upper triangle, right) predicted interaction scores, for all possible protein pairs in the set. According to STRING v11.5, all RFC-Rad17 partners and C1-Serping1 partners physically interact with each other while no interaction between these two groups of proteins is known (see B).

(B) Two graphs representing STRING data on the two complexes. No mouse protein (gray) has been seen in training, which is done on *H. sapiens*. Edge colors represent the STRING confidence score for the interactions. Scale colors as in Figure 2.

To further check the ability of SENSE-PPI to correctly infer interactions in other organisms, we reconsidered as seeds the C1 complex (C1r, C1s, C1q), with its plasma protease C1 inhibitor Serping1 playing a role in the immune response, and the replication factor RFC complex with Rad17, forming a DNA damage checkpoint complex in the cell cycle, both in *M. musculus*. In the mouse, as in human, these two subnetworks are not supposed to share common partners since the two complexes involved assure drastic differences in functions and subcellular localizations. We compared SENSE-PPI to the D-SERIP and Topsy-Turvy models²² (see also Table 1), where training was performed for all models on the human STRING11.0 dataset. For each seed protein, we chose four partners with the highest confidence score in STRING. Figure 3 shows two heatmaps where network reconstruction shows two separate clusters in SENSE-PPI predictions (lower triangle). D-SERIP (upper triangle, Figure 3A) correctly predicts the RFC-Rad17 complex network but fails on the C1-Serping1 complex. Topsy-Turvy improves D-SERIP predictions on the scoring values of interaction pairs in the RFC-Rad17 complex by being more confident in the predictions. On the C1-Serping1 complex, it identifies some interactions with variable confidence compared to D-SERIP but much less sharply than SENSE-PPI. With an increased number of true positives, the number of false positives increases as well though, and the model incorrectly finds nonexistent interactions between proteins belonging to different complexes. In contrast, with a clear separation between C1-Serping1 and RFC-Rad17 complexes, SENSE-PPI shows the quality of its predictions.

Additionally, predictions at a large proteome scale were conducted to assess the PPI reconstruction capabilities. To substantiate the distinctions between networks from different subcellular localizations, as initially suggested by Figure 2A, we expanded our analysis. A new test dataset was assembled from the UniProtKB/Swiss-Prot database, comprising two distinct groups of proteins: the first group included proteins with experimentally verified subcellular locations in the nucleus (term SL-0191), and the second group consisted of secreted proteins (term SL-0243). Both groups were filtered to include only human proteins with lengths between 50 and 800 amino acids, resulting in 2,939 nuclear and 1,024 secreted proteins. Predictions were made for all possible pairs within these groups, which were subsequently categorized into three types: pairs of nuclear proteins, pairs of secreted proteins, and mixed pairs (nuclear paired with secreted). Figure 2D illustrates the proportion

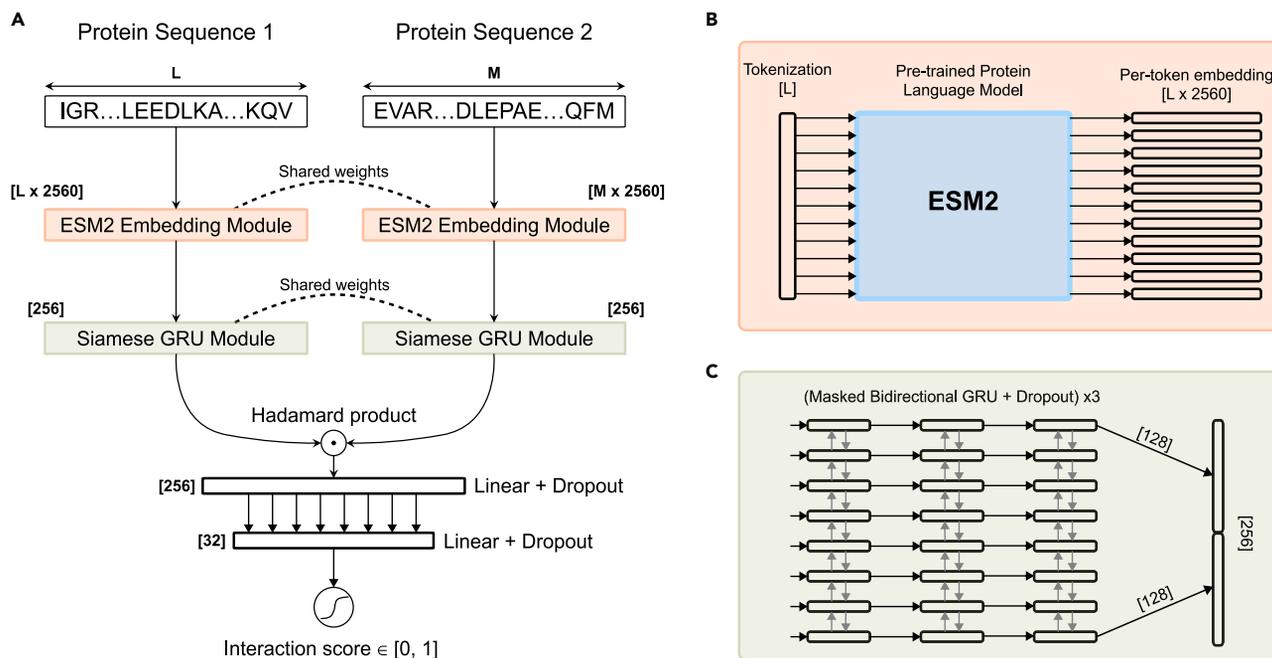


Figure 4. Schematic representation of SENSE-PPI

(A) SENSE-PPI takes two input sequences of lengths L and M through the ESM2 Embedding Module (orange) and transforms them into two large tensors of size $L \times 2560$ and $M \times 2560$, respectively. The two matrices are then independently processed by the Siamese GRU Module (green). The final output vectors are combined together using the Hadamard product and further processed with two linear layers in order to calculate the final interaction score in the interval $[0, 1]$. (B) Details of the ESM2 Embedding Module mapping each amino acid in the input sequence into a vector of length 2,560. The module is essentially the ESM2 model, which comprises 36 layers (3B parameters version). The output of the module is a per-token representation of the input. (C) Details of the Siamese GRU module composed of three bidirectional GRU layers. The last GRU layer follows a “many-to-one” output scheme that produces a vector of a unified length of 256 (128 for one direction and 128 for the opposite one) for proteins of any length. Since all the sequences are padded to the same fixed size during training, every GRU layer was implemented to be dynamical. The layers were configured to accept a sequence mask as a secondary input: this was done in order to process positions in a sequence that belong exclusively to the input data and skip all the padding.

of interactions identified within these categories. Of the 4,317,391 pairs consisting of two nuclear proteins, 149,126 interactions were predicted, representing 3.45% of all possible combinations. For pairs of two secreted proteins, the proportion of positive interactions was 3.02% (15,833 out of 523,776 pairs). In contrast, the mixed group showed a significant decrease to 0.57% (17,110 out of 3,009,536 pairs), indicating a substantial depletion of interactions compared to the first two groups. However, the presence of some interactions in the mixed group can be explained by the overlap of 62 entries in the “nucleus” and “secreted” subcellular location terms. Forty-two percent of the mixed group interactions contain at least one of these 62 shared proteins. At the same time, the percentage of interactions containing shared proteins in the “nuclear+nuclear” and “secreted+secreted” groups is 3% and 16%, respectively.

Lastly, we evaluated the approximate number of interactions within the human proteome. Our test dataset included all human sequences from UniProtKB/Swiss-Prot that are between 50 and 800 amino acids in length, as of February 2024. We analyzed 16,675 proteins, representing 139,019,475 potential protein pairs. The number of interactions identified varied according to the chosen threshold (Figure 2E): we observed 951,286 interactions at a threshold of 0.5; 279,348 interactions at 0.9; and 82,243 interactions at 0.99.

DISCUSSION

Current PPI experimental datasets are still far away from having reached completion, even on well-studied model organisms such as yeast or humans,^{53,54} and their curation is an active area of research.⁵⁵ Today, physical interaction networks obtained by high-throughput techniques are found to include numerous non-functional PPI²⁹ and at the same time, many missing true interactions. This is the main reason that an *ab initio* computational reconstruction of PPIs would provide invaluable information.

We introduced SENSE-PPI, a deep learning model for predicting protein-protein interactions based entirely on sequence information (Figure 4). Trained in a single species, SENSE-PPI predicts the proteome within the same species and generalizes to other species. Trained on physical interactions between proteins with stable structures, it generalizes to proteins interacting with IDPs. Trained on the interaction between proteins of human and viral species, it generalizes to the interaction between human proteins and proteins of new viruses.

The high performance of SENSE-PPI in single-model species predictions has been demonstrated here to open up new possibilities in PPI network reconstruction for non-model organisms, which are often characterized by very sparse or absent biological information. Our analysis shows that training on a model species such as *H. sapiens* and testing on a phylogenetically distant species such as *C. elegans* yet provides

excellent inference. Moreover, we show that if we construct a training dataset out of model organisms that are evolutionary close to a non-model species, we are able to improve the performance: we succeeded in improving the test scores on *A. pisum* by switching the training data from *H. sapiens* to a combination of *C. elegans* and *D. melanogaster*. As a general PPI reconstruction strategy for a non-model organism, the user must first identify the best model organism(s) close to the non-model species, train on it, and infer PPI interactions in the non-model species.

SENSE-PPI is fast, exceeding the time limits of other DL approaches and outperforming the particularly time-consuming docking approaches.¹³ The major advantage of ESM2, used in SENSE-PPI, over other PLM models is that it is much faster (several minutes versus several hours compared to ProtT5-XL-UniRef50,⁵⁶ for example) while providing highly informative matrices. Indeed, ESMFold, which works with a similar embedding, was able to compete with AlphaFold in terms of performance, being 60 times faster and enabling the reconstruction of over 600 million protein structural models.³⁵ The advantage of deep learning architectures over other computational attempts in PPI reconstruction is that training concentrates all the computational weight, while the combinatorics of the problem (i.e., the quadratic number of potential pairwise interactions among the proteins) is handled efficiently by the trained architecture. As expected, SENSE-PPI shows that the computational issue can be overcome and that PPIs of tens of thousands of proteins can now be handled. Indeed, the identification of protein partners in a 100-protein dataset takes around 2 min with a single NVIDIA A100 80Go PCIe. Given that the computational time for ESM2 embeddings scales linearly in the number of proteins (the tensors are computed only once for each protein and then reused) and the prediction time scales quadratically, we can estimate that, using the same hardware, it takes around 2 h to predict the partners of a set of 1,000 proteins (including the computation of ESM2 embeddings for 25 min) and around a week for 10,000 (with 4.5 h for ESM2 embeddings).

Another major advantage of SENSE-PPI over other structure-based approaches to protein interactions, including docking, is that SENSE-PPI achieves far better results based solely on sequences, opening up new directions for development. The difficulties associated with intrinsic disorder, protein instability, or transient interactions, which are inherent in protein structures, are avoided by sequence-based methods, offering a truly revolutionary opportunity to accelerate and create innovative applications in the field of proteomics.

Even though the *ab initio* reconstruction of PPI networks for model species impressively improved with the years, augmenting the accuracy of the inference within a species remains a problem, especially if we wish to tackle questions on specialized interactions of a protein in different tissues, or differentiating the interaction of paralogous proteins within the organism. For this, one expects to improve in this direction through the development of deep learning architectures that will help us to interpret and disentangle interaction signals better than now. Also, searching for protein interactions in a species without sufficient training information and training in a model species that is phylogenetically very far away from the species used in testing are still difficult tasks. Testing on multiple species, such as human-bacteria interactions, for example, remains a wide-open problem. This is clearly to state that the PPI network reconstruction problem, even when we consider proteins that do not undergo post-transcriptional changes, remains far from being solved but that current achievements can, nonetheless, be applied to the screening of large sets of proteins in search of their partners for many biological questions. In complex predictive scenarios, combining SENSE-PPI with deep learning tools like AlphaFold-Multimer⁵⁷ may offer significant advantages. While both sequence- and structure-based methods encounter challenges, particularly evident in proteins with sparse multiple sequence alignments (Figure S2), each method faces unique difficulties. Structure-based PPI reconstructions struggle with intrinsically disordered proteins (IDPs), where the lack of fixed structures confounds predictions. Conversely, sequence-based methods often fail with orphan proteins that significantly diverge from known sequence data. By integrating these models, we can leverage the strengths of one approach to compensate for the weaknesses of the other, potentially enhancing both the accuracy and reliability of predictions in these challenging situations.

In conclusion, SENSE-PPI is a deep learning model that provides results that were not conceivable a few years ago, with an accuracy that could never be reached by state-of-the-art molecular docking approaches applied to proteins with known structures nor deep learning methods before it.

Limitations of the study

Identifying PPIs in a species with insufficient training data, or in a species with a model trained on phylogenetically distant species, presents significant challenges. The task of predicting cross-species interactions, such as those between human and bacterial proteins, particularly for unknown protein pairs, remains an unresolved issue. Overall, the reconstruction of PPI networks, even when focusing on proteins that do not undergo post-transcriptional modifications, is far from being fully achieved. Despite these limitations, current methodologies can still be effectively applied to screen large sets of proteins to identify potential interaction partners, thereby addressing many biological questions.

STAR★METHODS

Detailed methods are provided in the online version of this paper and include the following:

- [KEY RESOURCES TABLE](#)
- [RESOURCE AVAILABILITY](#)
 - Lead contact
 - Materials availability
 - Data and code availability
- [METHOD DETAILS](#)

- Datasets
- Design of the SENSE-PPI architecture
- Deep learning architectures used for comparison
- Guo's yeast dataset: Comparative results not reported in the main text
- Implementation details
- SENSE-PPI framework
- **QUANTIFICATION AND STATISTICAL ANALYSIS**

SUPPLEMENTAL INFORMATION

Supplemental information can be found online at <https://doi.org/10.1016/j.isci.2024.110371>.

ACKNOWLEDGMENTS

i-Bio initiative at Sorbonne University for the PhD fellowship and travel grants to K.V.; Sorbonne Center for Artificial Intelligence (SCAI) and LIP6 laboratory (UMR 7606, Sorbonne University-CNRS) for access to their clusters; Agence Nationale de Recherches sur le Sida et les Hépatites Virales [ANRS–AAP-2021-CSS-12] (A.C.).

AUTHOR CONTRIBUTIONS

K.V. and A.C. conceived and designed the experiments. K.V. performed the experiments. K.V. and A.C. discussed the biological significance of the datasets. L.B. provided additional insights during these discussions. K.V. and A.C. analyzed the data. A.C. and K.V. wrote the paper. All authors read and approved the final manuscript.

DECLARATION OF INTERESTS

The authors declare no competing interests.

Received: December 12, 2023

Revised: May 4, 2024

Accepted: June 21, 2024

Published: June 25, 2024

REFERENCES

1. Laddach, A., Chung, S.S., and Fraternali, F. (2018a). Prediction of protein-protein interactions: Looking through the kaleidoscope. In *Encyclopedia of Bioinfo and Comput Biol: ABC of Bioinformatics* (Elsevier), pp. 834–848.
2. Laddach, A., Ng, J.C.-F., Chung, S.S., and Fraternali, F. (2018b). Genetic variants and protein-protein interactions: a multidimensional network-centric view. *Curr. Opin. Struct. Biol.* **50**, 82–90.
3. Mosca, R., Céol, A., and Aloy, P. (2013). Interactome3d: adding structural details to protein networks. *Nat. Methods* **10**, 47–53.
4. Orchard, S., Ammari, M., Aranda, B., Breuza, L., Briganti, L., Broackes-Carter, F., Campbell, N.H., Chavali, G., Chen, C., Del-Toro, N., et al. (2014). The mintact project—intact as a common curation platform for 11 molecular interaction databases. *Nucleic Acids Res.* **42**, D358–D363.
5. Rolland, T., Taşan, M., Charlotheaux, B., Pevzner, S.J., Zhong, Q., Sahni, N., Yi, S., Lemmens, I., Fontanillo, C., Mosca, R., et al. (2014). A proteome-scale map of the human interactome network. *Cell* **159**, 1212–1226.
6. Szklarczyk, D., Morris, J.H., Cook, H., Kuhn, M., Wyder, S., Simonovic, M., Santos, A., Doncheva, N.T., Roth, A., Bork, P., et al. (2017). The string database in 2017: quality-controlled protein-protein association networks, made broadly accessible. *Nucleic Acids Res.* **45**, D362–D368.
7. Fodor, E.L., Hassan, S.S., Lemke, N., Barh, D., Silva, A., Ferreira, R.S., and Azevedo, V. (2014). An improved interolog mapping-based computational prediction of protein-protein interactions with increased network coverage. *Integr. Biol.* **6**, 1080–1087.
8. Garcia-Garcia, J., Schleker, S., Klein-Seetharaman, J., and Oliva, B. (2012). Bips: Biana interolog prediction server. a tool for protein-protein interaction inference. *Nucleic Acids Res.* **40**, W147–W151.
9. Guo, Y., Yu, L., Wen, Z., and Li, M. (2008). Using support vector machine combined with auto covariance to predict protein-protein interactions from protein sequences. *Nucleic Acids Res.* **36**, 3025–3030.
10. Hamp, T., and Rost, B. (2015). Evolutionary profiles improve protein-protein interaction prediction from sequence. *Bioinformatics* **31**, 1945–1950.
11. Huynen, M., Snel, B., Lathe, W., 3rd, and Bork, P. (2000). Predicting protein function by genomic context: quantitative evaluation and qualitative inferences. *Genome Res.* **10**, 1204–1210.
12. Laine, E., and Carbone, A. (2017). Protein social behavior makes a stronger signal for partner identification than surface geometry. *Proteins* **85**, 137–154.
13. Lopes, A., Sacquin-Mora, S., Dimitrova, V., Laine, E., Ponty, Y., and Carbone, A. (2013). Protein-protein interactions in a crowded environment: an analysis via cross-docking simulations and evolutionary information. *PLoS Comput. Biol.* **9**, e1003369.
14. Matthews, L.R., Vaglio, P., Reboul, J., Ge, H., Davis, B.P., Garrels, J., Vincent, S., and Vidal, M. (2001). Identification of potential interaction networks using sequence-based searches for conserved protein-protein interactions or “interologs”. *Genome Res.* **11**, 2120–2126.
15. Morilla, I., Lees, J.G., Reid, A.J., Orengo, C., and Ranea, J.A.G. (2010). Assessment of protein domain fusions in human protein interaction networks prediction: Application to the human kinetochore model. *N. Biotech.* **27**, 755–765.
16. Pellegrini, M., Marcotte, E.M., Thompson, M.J., Eisenberg, D., and Yeates, T.O. (1999). Assigning protein functions by comparative genome analysis: Protein phylogenetic profiles. *Proc. Natl. Acad. Sci. USA* **96**, 4285–4288.
17. Sacquin-Mora, S., Carbone, A., and Lavery, R. (2008). Identification of protein interaction partners and protein-protein interaction sites. *J. Mol. Biol.* **382**, 1276–1289.
18. Scott, M.S., and Barton, G.J. (2007). Probabilistic prediction and ranking of human protein-protein interactions. *BMC Bioinf.* **8**, 239.
19. Wass, M.N., Fuentes, G., Pons, C., Pazos, F., and Valencia, A. (2011). Towards the prediction of protein interaction partners using physical docking. *Mol. Syst. Biol.* **7**, 469.

20. Yu, H., Luscombe, N.M., Lu, H.X., Zhu, X., Xia, Y., Han, J.-D.J., Bertin, N., Chung, S., Vidal, M., Gerstein, M., et al. (2004). Annotation transfer between genomes: protein-protein interologs and protein-dna regulogs. *Genome Res.* 14, 1107–1118.
21. Chen, M., Ju, C.J.T., Zhou, G., Chen, X., Zhang, T., Chang, K.-W., Zaniolo, C., and Wang, W. (2019). Multifaceted protein-protein interaction prediction based on Siamese residual RCNN. *Bioinformatics* 35, i305–i314.
22. Singh, R., Devkota, K., Sledzieski, S., Berger, B., and Cowen, L. (2022). Topsy-Turvy: integrating a global view into sequence-based PPI prediction. *Bioinformatics* 38, i264–i272.
23. Sledzieski, S., Singh, R., Cowen, L., and Berger, B. (2021). D-script translates genome to phenome with sequence-based, structure-aware, genome-scale predictions of protein-protein interactions. *Cell Syst.* 12, 969–982.e6.
24. Tsukiyama, S., Hasan, M.M., Fujii, S., and Kurata, H. (2021). Lstm-phv: prediction of human-virus protein-protein interactions by lstm with word2vec. *Briefings Bioinf.* 22, bbab228.
25. Ito, T., Chiba, T., Ozawa, R., Yoshida, M., Hattori, M., and Sakaki, Y. (2001). A comprehensive two-hybrid analysis to explore the yeast protein interactome. *Proc. Natl. Acad. Sci. USA* 98, 4569–4574.
26. Uetz, P., Giot, L., Cagney, G., Mansfield, T.A., Judson, R.S., Knight, J.R., Lockshon, D., Narayan, V., Srinivasan, M., Pochart, P., et al. (2000). A comprehensive analysis of protein-protein interactions in *Saccharomyces cerevisiae*. *Nature* 403, 623–627.
27. Bensimon, A., Heck, A.J.R., and Aebersold, R. (2012). Mass spectrometry-based proteomics and network biology. *Annu. Rev. Biochem.* 81, 379–405.
28. Yu, H., Braun, P., Yildirim, M.A., Lemmens, I., Venkatesan, K., Sahalie, J., Hirozane-Kishikawa, T., Gebreab, F., Li, N., Simonis, N., et al. (2008). High-quality binary protein interaction map of the yeast interactome network. *Science* 322, 104–110.
29. Levy, E.D., Landry, C.R., and Michnick, S.W. (2009). How perfect can protein interactomes be? *Sci. Signal.* 2, pe11.
30. Tang, T., Zhang, X., Liu, Y., Peng, H., Zheng, B., Yin, Y., and Zeng, X. (2023). Machine learning on protein-protein interaction prediction: models, challenges and trends. *Briefings Bioinf.* 24, bbad076.
31. Wang, X.-W., Madeddu, L., Spirohn, K., Martini, L., Fazzone, A., Becchetti, L., Wytock, T.P., Kovács, I.A., Balogh, O.M., Benczik, B., et al. (2023). Assessment of community efforts to advance network-based prediction of protein-protein interactions. *Nat. Commun.* 14, 1582.
32. Lobingier, B.T., Hüttenhain, R., Eichel, K., Miller, K.B., Ting, A.Y., von Zastrow, M., and Krogan, N.J. (2017). An approach to spatiotemporally resolve protein interaction networks in living cells. *Cell* 169, 350–360.e12.
33. Scott, J.D., and Pawson, T. (2009). Cell signaling in space and time: where proteins come together and when they're apart. *Science* 326, 1220–1224.
34. Cho, K., van Merriënboer, B., Bahdanau, D., and Bengio, Y. (2014). On the properties of neural machine translation: Encoder-decoder approaches. In *Proc SSSST-8 (Assoc Comput Linguistics)*, pp. 103–111.
35. Lin, Z., Akin, H., Rao, R., Hie, B., Zhu, Z., Lu, W., Smetanin, N., Verkuil, R., Kabeli, O., Shmueli, Y., et al. (2022). Evolutionary-scale prediction of atomic level protein structure with a language model. Preprint at *bioRxiv*. <https://doi.org/10.1101/2022.07.20.500902>.
36. Perovic, V., Sumonja, N., Marsh, L.A., Radovanovic, S., Vukicevic, M., Roberts, S.G.E., and Veljkovic, N. (2018). Idppi: Protein-protein interaction analyses of human intrinsically disordered proteins. *Sci. Rep.* 8, 10563.
37. Madan, S., Demina, V., Stapf, M., Ernst, O., and Fröhlich, H. (2022). Accurate prediction of virus-host protein-protein interactions via a siamese neural network using deep protein sequence embeddings. *Patterns* 3, 100551.
38. Carroll, D., Daszak, P., Wolfe, N.D., Gao, G.F., Morel, C.M., Morzaria, S., Pablos-Méndez, A., Tomori, O., and Mazet, J.A.K. (2018). The global virome project. *Science* 359, 872–874.
39. Lewin, H.A., Robinson, G.E., Kress, W.J., Baker, W.J., Coddington, J., Crandall, K.A., Durbin, R., Edwards, S.V., Forest, F., Gilbert, M.T.P., et al. (2018). Earth biogenome project: Sequencing life for the future of life. *Proc. Natl. Acad. Sci. USA* 115, 4325–4333.
40. Li, F., Zhao, X., Li, M., He, K., Huang, C., Zhou, Y., Li, Z., and Walters, J.R. (2019). Insect genomes: progress and challenges. *Insect Mol. Biol.* 28, 739–758.
41. Rhie, A., McCarthy, S.A., Fedrigo, O., Damas, J., Formenti, G., Koren, S., Uliano-Silva, M., Chow, W., Fungtammasan, A., Kim, J., et al. (2021). Towards complete and error-free genome assemblies of all vertebrate species. *Nature* 592, 737–746.
42. Thompson, L.R., Sanders, J.G., McDonald, D., Amir, A., Ladau, J., Locey, K.J., Prill, R.J., Tripathi, A., Gibbons, S.M., Ackermann, G., et al. (2017). A communal catalogue reveals earth's multiscale microbial diversity. *Nature* 551, 457–463.
43. Kumar, S., Stecher, G., Suleski, M., and Hedges, S.B. (2017). Timetree: a resource for timelines, timetrees, and divergence times. *Mol. Biol. Evol.* 34, 1812–1819.
44. Park, Y., and Marcotte, E.M. (2012). Flaws in evaluation schemes for pair-input computational predictions. *Nat. Methods* 9, 1134–1136.
45. Steinegger, M., and Söding, J. (2017). Mmseqs2 enables sensitive protein sequence searching for the analysis of massive data sets. *Nat. Biotechnol.* 35, 1026–1028.
46. Pentony, M.M., and Jones, D.T. (2010). Modularity of intrinsic disorder in the human proteome. *Proteins* 78, 212–221.
47. Oates, M.E., Romero, P., Ishida, T., Ghalwash, M., Mizianty, M.J., Xue, B., Dosztányi, Z., Uversky, V.N., Obradovic, Z., Kurgan, L., et al. (2012). D2p2: database of disordered protein predictions. *Nucleic Acids Res.* 41, D508–D516.
48. Seoane, B., and Carbone, A. (2021). The complexity of protein interactions unravelled from structural disorder. *PLoS Comput. Biol.* 17, e1008546.
49. Seoane, B., and Carbone, A. (2022). Soft disorder modulates the assembly path of protein complexes. *PLoS Comput. Biol.* 18, e1010713.
50. Van Der Lee, R., Buljan, M., Lang, B., Weatheritt, R.J., Daughdrill, G.W., Dunker, A.K., Fuxreiter, M., Gough, J., Gsponer, J., Jones, D.T., et al. (2014). Classification of intrinsically disordered regions and proteins. *Chem. Rev.* 114, 6589–6631.
51. Guirimand, T., Delmotte, S., and Navratil, V. (2015). Virhostnet 2.0: surfing on the web of virus/host molecular interactions data. *Nucleic Acids Res.* 43, D583–D587.
52. Szklarczyk, D., Kirsch, R., Koutrouli, M., Nastou, K., Mehryary, F., Hachilif, R., Gable, A.L., Fang, T., Doncheva, N.T., Pyysalo, S., et al. (2023). The STRING database in 2023: protein-protein association networks and functional enrichment analyses for any sequenced genome of interest. *Nucleic Acids Res.* 51, D638–D646.
53. Hart, G.T., Ramani, A.K., and Marcotte, E.M. (2006). How complete are current yeast and human protein-interaction networks? *Genome Biol.* 7, 120.
54. Stumpf, M.P.H., Thorne, T., de Silva, E., Stewart, R., An, H.J., Lappe, M., and Wiuf, C. (2008). Estimating the size of the human interactome. *Proc. Natl. Acad. Sci. USA* 105, 6959–6964.
55. Schaefer, M.H., Fontaine, J.-F., Vinayagam, A., Porras, P., Wanker, E.E., and Andrade-Navarro, M.A. (2012). Hippie: Integrating protein interaction networks with experiment based quality scores. *PLoS One* 7, e31826.
56. Elnaggar, A., Heinzinger, M., Dallago, C., Rehawi, G., Wang, Y., Jones, L., Gibbs, T., Feher, T., Angerer, C., Steinegger, M., et al. (2020). Prottrans: Towards cracking the language of life's code through self-supervised deep learning and high performance computing. Preprint at *bioRxiv*. <https://doi.org/10.1101/2020.07.12.199554>.
57. Evans, R., O'Neill, M., Pritzel, A., Antropova, N., Senior, A., Green, T., Židek, A., Bates, R., Blackwell, S., Yim, J., et al. (2021). Protein complex prediction with alphafold-multimer. Preprint at *bioRxiv*. <https://doi.org/10.1101/2021.10.04.463034>.
58. UniProt Consortium (2023). UniProt: the Universal Protein Knowledgebase in 2023. *Nucleic Acids Res.* 51, D523–D531.
59. Szklarczyk, D., Gable, A.L., Lyon, D., Junge, A., Wyder, S., Huerta-Cepas, J., Simonovic, M., Doncheva, N.T., Morris, J.H., Bork, P., et al. (2019). STRING v11: protein-protein association networks with increased coverage, supporting functional discovery in genome-wide experimental datasets. *Nucleic Acids Res.* 47, D607–D613.
60. Szklarczyk, D., Gable, A.L., Nastou, K.C., Lyon, D., Kirsch, R., Pyysalo, S., Doncheva, N.T., Legeay, M., Fang, T., Bork, P., et al. (2021). The STRING database in 2021: customizable protein-protein networks, and functional characterization of user-uploaded gene/measurement sets. *Nucleic Acids Res.* 49, D605–D612.
61. Yang, X., Lian, X., Fu, C., Wuchty, S., Yang, S., and Zhang, Z. (2021). Hvidb: a comprehensive database for human-virus protein-protein interactions. *Briefings Bioinf.* 22, 832–844.
62. Ofer, D., Brandes, N., and Linal, M. (2021). The language of proteins: Nlp, machine learning & protein sequences. *Comput. Struct. Biotechnol. J.* 19, 1750–1758.
63. Alley, E.C., Khimulya, G., Biswas, S., AlQuraishi, M., and Church, G.M. (2019). Unified rational protein engineering with sequence-based deep representation learning. *Nat. Methods* 16, 1315–1322.
64. Heinzinger, M., Elnaggar, A., Wang, Y., Dallago, C., Nechaev, D., Matthes, F., and Rost, B. (2019). Modeling aspects of the

- language of life through transfer-learning protein sequences. *BMC Bioinf.* 20, 723–817.
65. Rao, R., Bhattacharya, N., Thomas, N., Duan, Y., Chen, X., Canny, J., Abbeel, P., and Song, Y.S. (2019). Evaluating protein transfer learning with tape. *Adv. NeurIPS* 32, 9689–9701.
 66. Rives, A., Meier, J., Sercu, T., Goyal, S., Lin, Z., Liu, J., Guo, D., Ott, M., Zitnick, C.L., Ma, J., and Fergus, R. (2021). Biological structure and function emerge from scaling unsupervised learning to 250 million protein sequences. *Proc. Natl. Acad. Sci. USA* 118, e2016239118.
 67. Devkota, K., Murphy, J.M., and Cowen, L.J. (2020). GLIDE: combining local methods and diffusion state embeddings to predict missing interactions in biological networks. *Bioinformatics* 36, i464–i473.
 68. Salwinski, L., Miller, C.S., Smith, A.J., Pettit, F.K., Bowie, J.U., and Eisenberg, D. (2004). The database of interacting proteins: 2004 update. *Nucleic Acids Res.* 32, D449–D451.
 69. Paszke, A., Gross, S., Massa, F., Lerer, A., Bradbury, J., Chanan, G., Killeen, T., Lin, Z., Gimelshein, N., Antiga, L., et al. (2019). Pytorch: An Imperative Style, High-Performance Deep Learning Library (Curran Associates, Inc), pp. 8024–8035.

STAR★METHODS

KEY RESOURCES TABLE

REAGENT or RESOURCE	SOURCE	IDENTIFIER
Deposited data		
STRING	Szklarczyk, D. et al. (2022) ⁵²	https://string-db.org/
UniProt	Consortium, T. U. (2022) ⁵⁸	https://www.uniprot.org/
Datasets used in tests	This paper	http://gitlab.lcqb.upmc.fr/Konstvw/SENSE-PPI
Software and algorithms		
SENSE-PPI	This paper	http://gitlab.lcqb.upmc.fr/Konstvw/SENSE-PPI
D-SCRIPT	Sledzieski, S. et al. (2021) ²³	https://dscrip.csail.mit.edu/
Topsy-Turvy	Singh, R. et al. (2022) ²²	https://cb.csail.mit.edu/cb/topsyturvy/
PIPR	Chen, M. et al. (2019) ²¹	https://github.com/muhaochen/seq_ppi
STEP	Madan, S. et al. (2022) ³⁷	https://github.com/SCAI-BIO/STEP

RESOURCE AVAILABILITY

Lead contact

Further information and requests for materials and code should be directed to and will be fulfilled by the Lead Contact, Alessandra Carbone (alessandra.carbone@lip6.fr).

Materials availability

This study did not generate new materials.

Data and code availability

- All datasets utilized in this study are available in the SENSE-PPI repository and are publicly accessible as of the publication date. The DOI is provided in the [key resources table](#).
- The original code is available in the SENSE-PPI repository and is publicly accessible as of the publication date. The DOI is provided in the [key resources table](#).
- Any additional information needed to reanalyze the data reported in this paper can be obtained from the [lead contact](#) upon request.

METHOD DETAILS

Datasets

The STRING11.0 datasets from human and other species

A large human dataset was constructed to test the D-SCRIPT deep learning model for the prediction of PPIs.²³ The positive interactions were extracted from STRING v11,⁵⁹ a database encompassing diverse PPI networks and consolidating comprehensive information from various primary sources. They are high-confidence interactions supported by a positive experimental evidence score in STRING. Protein sequences are greater than 50 and smaller than 800 amino acids in length, and protein pairs exhibit pairwise sequence identities <40%. By construction, any two protein pairs A-B and C-D in the dataset are *non-redundant*, that is both pairs of sequences A, C, and B, D have <40% sequence identity respectively. Eliminating redundancy is crucial as it prevents the model from relying solely on sequence similarity when predicting interactions. To establish an appropriate balance between positive and negative PPIs, negative interactions have been generated by randomly picking two proteins from distinguished pairs in the non-redundant set of positive interactions. These negative samples were created in a 1:10 positive-negative ratio with the aim of better approximating the frequency of positive interactions observed in nature. These construction criteria lead to a primary dataset consisting of 47,932 positive protein interactions and 479,320 negative protein interactions for *H. sapiens*. These interactions were divided into training and testing sets using an 80% - 20% ratio, respectively. Additionally, we have used 4 smaller datasets of *M. musculus*, *D. melanogaster*, *C. elegans* and *S. cerevisiae* from²³ made using the procedure described above and containing 5,000 positive and 50,000 negative interactions each.

The STRING11.5_neighbor_exclusion human dataset

We constructed a second human PPI dataset directly from the STRING database v11.5⁵² by focusing exclusively on physical interactions from *H. sapiens*. Similar to the STRING11.0 dataset, we only considered proteins in a length range from 50 to 800 amino acids in order to optimize

computational resources and to avoid potential biases introduced by very short or very long sequences. A STRING combined confidence score threshold of 0.5 was applied and interactions obtained by homology transfer were excluded. Note that the STRING combined confidence score is derived from a variety of sources, including gene fusion, text mining, experiments, annotated pathways, and other methods (see <https://string-db.org/> for a comprehensive list of these sources). Additionally, we implemented a clustering strategy to eliminate redundant interactions. This involved using MMseqs2 to cluster proteins at 40% sequence identity, followed by the removal of redundant pairs, a process similar to that used for the STRING11.0, dataset.

The key distinction between the human STRING11.0, dataset and the human STRING11.5_neighbor_exclusion, dataset rests in the use of a stricter criterion for constructing negative training pairs, specifically the “neighboring-exclusion” condition applied to the latter. Specifically, this involves a secondary clustering phase where proteins in the non-redundant dataset are clustered using a 40% sequence identity threshold. Subsequently, a protein pair, A-B, is selected for the negative set only if there is no known interaction between A and B in the STRING database, and if B is not in the same cluster as any of A’s known interactors. This method ensures that the selection of negative pairs is based not only on the absence of known interactions but also on reducing potential homology-related biases within shared clusters.

As above, in order to reflect the natural frequency of interactions observed in biological systems, the positive-to-negative ratio in the number of protein pairs was set to 1:10. Therefore, the full dataset contains 86,304 positive and 863,040 negative interactions. This ratio aims to approximate the relative scarcity of true positive interactions, thus ensuring a more representative distribution of positive and negative examples in the dataset. We tested SENSE-PPI using other positive-to-negative ratios: 1:1 and 1:100, following the same procedure. We did not consider larger ratios, such as 1:1000, due to the “neighboring-exclusion” condition, which restricts the availability of potential candidates for negative pairing at such scales. The comparative performance of the model across different positive-to-negative ratios is detailed in Table S7. Although the performance at a 1:100 ratio is comparable, the 1:10 ratio was found to be optimal in terms of the training time and performance trade-off. This ratio is ten times faster in both dataset creation and training processes, without compromising the quality of predictions.

Ten datasets from model and non-model organisms

We used the same methodology employed for the construction of the STRING11.5_neighbor_exclusion human dataset to build test sets on four non-model organisms, *E. caballus*, *B. taurus*, *N. scutatus*, and *A. pisum*, and six model organisms, *M. musculus*, *D. melanogaster*, *C. elegans*, *G. gallus*, *S. cerevisiae* and *E. coli*. All test sets contain 15,000 positive and 150,000 negative entries that were extracted from the STRING v12.0.⁶⁰ The ensemble of these datasets can be found in the SENSE-PPI repository under the name STRING12.0_species.

To conduct additional testing on non-model organisms to study how evolutionary distance affects the model performance, we have also constructed two separate training datasets. The first one is composed of human STRING11.5_neighbor_exclusion dataset augmented by data from *D. melanogaster* and *C. elegans*. This augmented portion contains 23,324 positive and 233,240 negative interactions coming from *D. melanogaster* and 17,835 positive and 178,350 negative interactions from *C. elegans*, extracted from STRING v12.0. Additionally, the second set comprises all the data from the first one as well as the PPI extracted from *G. gallus* (15,000 positives and 150,000 negatives).

The IDPpi dataset describing the interactions between structured proteins and IDPs

The IDPpi dataset³⁶ contains interactions with intrinsically disordered proteins (IDPs). It is organized in five distinguished datasets corresponding to the testing sets in.³⁶ Each subset comprises 3,500-4,000 pairs of proteins from *H. sapiens* where one of the partners is intrinsically disordered. Pairs were extracted from the Human Integrated Protein-Protein Interaction rEference (HIPPIE) database.⁵⁵ Negative interactions were defined through random sampling: two proteins (one that is an IDP and one that has a stable structure) were selected randomly and were assumed to be non-interacting unless they were already known as interacting. The positive-to-negative ratio of each testing set is 1:1.

Human-virus PPI dataset

This dataset²⁴ was designed in order to predict interactions between *H. sapiens* and different viruses. Interacting protein pairs were extracted from the Human-Virus Protein-Protein Interactions database (HVPPPI)⁵¹ and negative pairs were constructed using the protocol of the dissimilarity negative sampling, which uses a sequence similarity-based method to explore protein pairs that are unlikely to interact.²⁴

The training set includes 7,371 positive and 117,326 negative interactions between human and viral proteins and contains 16,933 human and 676 viral proteins in total. Viruses span over many families, and the full list of species included in training can be found in the SENSE-PPI repository (data/human_virus folder).

We constructed two test sets. The first set consists of interactions involving proteins from the *Epstein-Barr* virus. This set encompasses 99 viral proteins and 9,426 human proteins, of which only 596 human proteins demonstrate interaction with *Epstein-Barr* proteins. The set contains a total of 1,308 positive interactions and 14,428 negative interactions. The second dataset exclusively comprises interactions featuring proteins associated with *Influenza* viruses of type A, B, and C. This dataset includes 4,080 positive pairs and 15,724 negative pairs, incorporating a total of 95 viral proteins and 10,594 human proteins. Among these, 1,533 human proteins are involved in interactions with *Influenza* virus proteins. Notably, all proteins from both *Influenza* and *Epstein-Barr* viruses are excluded from the training process. However, it’s important to mention that all human proteins present in the test data have been included in the training set at least once.

Sequence identity scores for proteins in testing sets

Given a pair of proteins in a testing set, we define the 'mean pair sequence identity' as the mean of the sequence identities of the two proteins relative to the training set. Note that for each protein, we consider the maximum of all sequence identity values between the protein and all proteins in the training set. The computations were performed using the `mmseqs` search command of the MMseqs2 suite.⁴⁵ In this particular case, in order to calculate the sequence identity MMseqs2 calculates the ratio of identical aligned residues to the total number of aligned columns, which includes columns containing a gap in either sequence. This can be achieved by using `-alignment-mode 3`. If the search fails to find the similarity with a protein in the training set, the sequence identity value is set to 0 by default.

Design of the SENSE-PPI architecture

Recent advancements in deep learning for protein science have led to the introduction of pre-trained language models,⁶² which have been developed for natural language processing to learn compressed, informed, and abstract data representations that are later transferred to proteins by adapting them to amino acid sequences. PLMs have demonstrated tremendous potential for a broad range of protein-related problems, such as predicting 3D shapes, interactions, mutational outcomes, and subcellular localizations.^{63–66} They have been trained on large amounts of sequence data, such as the UniProt database.⁵⁸ Here, we use ESM2³⁵ which is a state-of-the-art general-purpose PLM that has been challenged to predict structure, contact sites, and other protein properties. We used the ESM2 version with 3B parameters, trained on the UniRef 50 dataset.

SENSE-PPI is designed as a Siamese architecture (Figure 4) consisting of two identical modules with shared weights that are merged together to provide a single output value $p \in [0, 1]$ describing the confidence in the interaction, where 0 and 1 stand for low and high confidence respectively. The Siamese design is needed to ensure the commutativity of the model: the output score has to be the same for pairs A-B and B-A.

The model takes two amino acid sequences of variable size (L, M) as input. Sequences longer than 800 amino acids are trimmed accordingly to this maximum size. The 36 layers of the PLM ESM2 are used to define two embeddings, of size $L \times 2560$ and $M \times 2560$, for the two amino acid sequences, respectively. Each sequence is processed by one of the Siamese modules sharing the weights with the second module. The Siamese module is composed of 3 layers of gated recurrent units (GRUs) which process the input sequence from one end to another. Each GRU layer is bidirectional, where half of the units are used for one direction of processing, and the other half for the opposite direction. The final layer produces a vector of shape 256 - it takes only the last output of GRU to perform a many-to-one type of processing, where a sequence is "projected" to a single number for each unit. The gated recurrent units process sequences that are masked beforehand: even though all the sequences are padded to the same length to fit the model, the GRU processes only the part that contains the actual sequence and skips the padding. After the GRU layers, the two output vectors are combined together via the Hadamard product. The resulting vector is passed through two linear layers in order to compute the final score. A dropout of 0.5 is applied at every layer (GRU and linear) except for the ESM2 modules.

Deep learning architectures used for comparison

We conducted a comparative analysis between SENSE-PPI and several existing sequence-based deep learning models, namely PIPR,²¹ D-SCRIPT,²³ Topsy-Turvy²² and STEP.³⁷

PIPR is an end-to-end Siamese architecture that incorporates a deep residual recurrent convolutional neural network. This architecture integrates multiple occurrences of convolution layers and residual gated recurrent units. Each amino acid is represented by an embedding that captures its contextual and physico-chemical relationship within the sequence. PIPR uses a multigranular feature aggregation process, effectively leveraging the sequential and robust local information present in protein sequences.

D-SCRIPT is a deep learning method specifically designed to predict physical PPIs. It features an LSTM-based PLM model for sequence embedding. This is followed by the computation of an outer product between two vector representations, resulting in a three-dimensional tensor. Linear and convolutional layers are then applied to preprocess the tensor, enabling a contact map to be predicted and the probability of interaction to be deduced. D-SCRIPT addresses the challenges associated with the weak cross-species generalization exhibited by previous state-of-the-art models. In addition, it establishes a significant overlap with contact maps for protein complexes with known 3D structures.

The Topsy-Turvy model combines a traditional sequence-based approach (D-SCRIPT) with a network-based approach (GLIDE⁶⁷) for PPI inference. While relying exclusively on sequence data for predictions, Topsy-Turvy employs a transfer-learning strategy during training, assimilating insights from global and molecular views of protein interactions. The model achieves state-of-the-art results, enabling genome-scale PPI predictions. It has been applied to non-model organisms.

STEP has been developed primarily for the prediction of virus–host PPIs. The core component of the STEP architecture is a BERT-based PLM, which generates sequence embeddings. These embeddings are combined through elementwise multiplication and further processed by linear layers.

Guo's yeast dataset: Comparative results not reported in the main text

The Guo dataset⁹ comprises 2,497 proteins and 11,188 interactions, with an equal number of positive and negative samples. The 5,943 positive interactions were extracted from the DIP 20070219 database.⁶⁸ Sequences are more than 50 amino acids long and the majority of them exhibit pairwise sequence identities <40%.⁹ The negative samples were generated through a random pairing of proteins that appeared in the

positive data but lacked any evidence of interaction. Pairs were subsequently filtered based on the subcellular localization of the proteins, thereby excluding non-interacting pairs residing in the same location. No further filtering was realized in this dataset to allow a fair comparison with models already evaluated on it. Namely, we compared SENSE-PPI with Guo's performance, PIPR and STEP. All approaches provide significantly high scores, and these scores are very close to each other. The average Matthews' correlation coefficient of SENSE-PPI is 95.46, which represents an improvement on the other deep learning models (94.77 and 94.17 for STEP and PIPR respectively), although this difference is not substantial (see Table S8). The exceptional success of all methods is explained by the characteristics of the computational experiment. Therefore, we have decided not to include it in the main text. All models were tested by performing 5-fold cross-validation on a dataset of protein pairs where the same protein appears in several pairs in the full dataset. Since the testing data represent a randomly chosen subset of sequence pairs at each cross-validating fold, in the validation step, each model mostly evaluates proteins already "seen" during training. Such pairs are known to show better performance than those excluded from training^{10,44} and remain particularly easy to predict.

Implementation details

The implementation of SENSE-PPI was carried out using Python 3 and the PyTorch deep learning framework.⁶⁹ We trained two versions of the model: a light version, used only for intermediate testing (see Figure 1E), and a regular version that produced all other results presented here. The regular version is presented in two final forms: the one that was trained on the STRING11.0 human dataset and the one that was trained on the STRING11.5_neighbor_exclusion human dataset. The light version of SENSE-PPI has a smaller model size, with an ESM2 embedding block containing 35M parameters and a classification head of 1M parameters. In contrast, the regular version of SENSE-PPI has an ESM2 embedding block comprising 3B parameters and a classification head of 2M parameters. The ESM2 module is not retrained, with the only trainable component being the classification head of the model. Each GRU layer contains 256 units followed by a dropout of 0.5. The three final linear layers have 256, 32, and 1 neurons respectively. The training was performed using AdamW optimizer, a batch size of 32, and a learning rate of 10^{-4} with the exception of the Guo's yeast dataset, for which we used a batch size of 64 and a learning rate of 10^{-3} .

The model was trained on eight NVIDIA A100 80Go PCIe GPUs, with the training process for a dataset comprising approximately 950,000 interactions (involving 9,648 unique proteins) taking about one day to complete. Predictions on a dataset of the same size were executed in 135 minutes using a single NVIDIA A100 80GB PCIe GPU.

SENSE-PPI framework

The *senseppi* package contains 5 primary commands: "train", "test", "predict", "predict_string", and "create_dataset". The first three commands are used to train the model, obtain metrics on test data, and perform predictions respectively.

The "predict_string" command can be used to calculate the predictions for a protein and its known partners according to STRING. The script takes a given number of proteins that are known to interact with a protein of interest, then it performs all versus all predictions and returns the results along with the visualization where one can easily compare the model's score with real data in STRING. Multiple proteins of interest can be used as input at the same time.

The "create_dataset" command creates a dataset for a given species from STRING which is based on the algorithm described in Section "The STRING11.5_neighbor_exclusion human dataset". The taxon ID is used here as the main input.

The package also contains an internal script to compute ESM2 embeddings. This script is a modification of the original code found in the ESM repository,³⁵ it was slightly edited in order to automate processes during training and testing with SENSE-PPI. All datasets used in this study are also provided, including both those newly created and those taken from other publications.

The two trained versions of SENSE-PPI, on the STRING11.0 and STRING11.5_neighbor_exclusion human datasets, are available in the SENSE-PPI repository in the pretrained_models folder, under the names dsript.ckpt and senseppi.ckpt respectively. In addition, this folder also contains pretrained versions used for other tests presented in this work such as the human-virus model, and models based on a combination of model species (fly/worm as well as fly/worm/human).

QUANTIFICATION AND STATISTICAL ANALYSIS

To properly evaluate SENSE-PPI performance we rely on a ground truth set by reference databases such as the STRING database and the following quantities: known interactions that are identified by SENSE-PPI (true positives, TP), interactions identified by SENSE-PPI which are not known (false positives, FP), interactions that are not found by SENSE-PPI but are known (false negatives, FN), and interactions that are not known and are not detected by SENSE-PPI (true negatives, TN). We use four standard metrics of performance: recall (sensitivity) = $TP/(TP + FN)$, precision (positive predictive value) = $TP/(TP + FP)$, F1-Score = $2 TP/(2 TP + FP + FN)$, and the Matthews' correlation coefficient $MCC = (TP \cdot TN - FP \cdot FN)/K$ where $K = \sqrt{(TP+FP)(TP+FN)(TN+FP)(TN+FN)}$.

We also use the precision-recall metric showing the tradeoff between precision and recall for different thresholds. A high area under the curve (AUPRC) represents both high recall and high precision, where high precision relates to a low false positive rate, and high recall relates to a low false negative rate. High scores for both show that the classifier is returning accurate results (high precision), as well as returning a majority of all positive results (high recall). Another measure is the AUROC, calculated as the area under the receiver operating characteristic (ROC) curve, showing the trade-off between true positive rate (TPR) and false positive rate (FPR) across different decision thresholds.