

Project Paper - Split Protein Prediction

Group 2 - Chuqiao Feng, Yuyan Cai, Yuanhang Li, Sohail Mohammod

December 8, 2024

Abstract

Fluorescent proteins (FPs) are widely used in biomolecular and cellular studies and assay developments, especially in biomolecular complementation assays. However, molecular splitting for complementary purposes in FPs usually require massive efforts of structural validation and experimentation. With advancement of Protein Language Models (PLMs), structural prediction for molecular splitting becomes possible. We propose leveraging PLM, such as ProteinBERT, combined with a Multilayer Perceptron (MLP), to predict optimal protein split sites. By fine-tuning ProteinBERT on labeled datasets for secondary structure and stability, we preserve pre-trained knowledge while adapting to fluorescence-specific tasks. Our pipeline includes encoding protein sequences, extracting attention weights to interpret the model, and visualizing attention distributions. The fine-tuned model predicts split sites by evaluating protein stability and secondary structure features, improving model precision. Results demonstrate robust performance on stability prediction and protein split site identification, validating the model's ability to enhance the creation of split variants for FPs.

1 Introduction

Since their discovery, FPs have revolutionized the entire field of molecular and cell biology by tracking various biological processes in real time. [7] FPs, which originate from marine organisms such as *Aequorea victoria*, can emit natural fluorescences that can be applied to study protein location tracking, protein-protein interactions (PPI), and dynamic changes in cellular systems [26]. After decades of development, advancements in FP engineering has led to various variants with increased brightness, higher stability, or expanded spectral features [7]. The capacity to visualize and measure biomolecular processes with high resolution has expedited not only biomedical research, but also drug discovery and clinical diagnostics [26]. Among all applications, Bimolecular Fluorescence Complementation (BiFC), which utilizes split FPs to visualize protein-protein interactions, exemplifies this potential by facilitating the study of intracellular signaling pathways, intercellular communications, and complex modulation networks in different cell types, thus enhancing current understanding of biological dynamics at molecular level [26]. In real world settings, increasing demands for high-throughput and precise diagnostics have emerged. The traditional way of choosing split site to generate BiFC tool utilize structure information of FPs. Split sites are typically chosen to avoid secondary structures and conserved regions, but the sheer number of suitable candidates makes this task challenging. [7]

To tackle challenges in protein split site searching, the goal of this project is to leverages PLMs combined with a MLP to predict optimal positions for complementation. The proposed workflow involves encoding protein sequence information using PLMs such as ProteinBERT and training the model on BiFC data gathered from the literature. This approach aims to enhance the BiFC toolkit by creating split variants for FPs like mCardinal and mNeptune, expanding the BiFC toolbox, which can empower more accurate and stronger visualization of PPIs and promote the molecular diagnostics and drug discovery

This innovative application of PLMs represents a significant step forward in computational protein engineering, providing a data-driven approach to enable high-throughput and efficient design of BiFC.

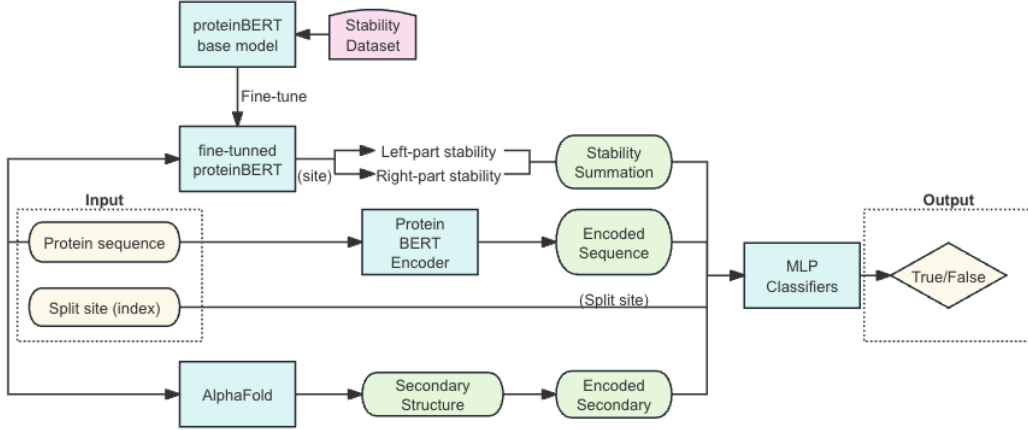


Figure 1: Pipeline of protein split site prediction

2 Materials and Methods

2.1 Graphical overview of methods

2.2 Dataset

2.2.1 Fluorescence split site

The split site data were collected from experimental results, where we looked through over thirty published papers. Fluorescence protein were split and attached to two different proteins, and when the two proteins are adjacent, the two split fluorescence fragments will assemble and glow. Sequences were collected from Protein Data Bank (PDB), UniProt, Fluorescent Protein Data Base (FPbase), and through Addgene plasmid sequences. You can find the data and their resources/ references we used at this link: https://github.com/Carrie1013/DS596-ProteinProject/blob/main/raw_sequence.csv. [1, 4–6, 8–15, 17–25, 27–31]

2.2.2 Protein stability

The stabilities of protein sequence that we used for fine-tuning proteinBERT were taken from the Tasks Assessing Protein Embeddings (TAPE) [?], this is a benchmark for evaluating protein sequence models, and they use data generated by a novel combination of parallel DNA synthesis and protein stability measurements [?]. The datasets include protein sequences labeled with their stabilities.

Protein Stability	Train	Valid	Test
mean	0.18	0.18	1.00
std	0.57	0.57	0.41
min	-1.97	-1.78	-0.27
max	3.40	3.24	2.66
data count	48251	5362	12851

Table 1: Data description of protein stability

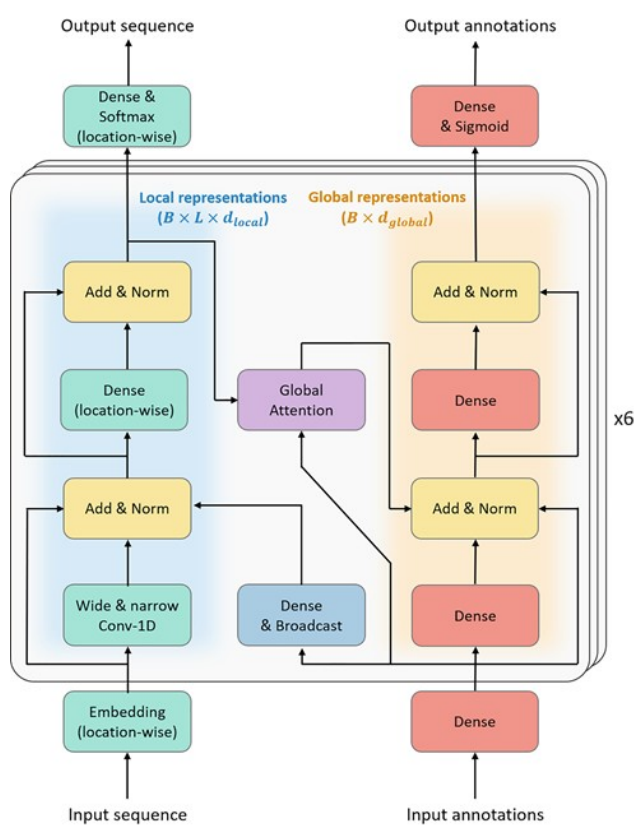
2.3 Principles underlying method

Sequence Encoding with proteinBERT

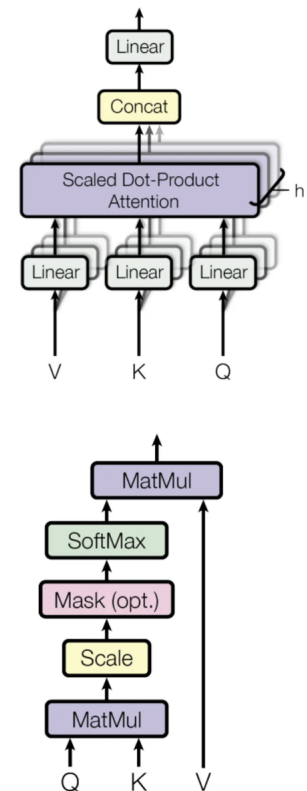
Sequence-based large-scale language models such as BERT are often used in natural language processing (NLP) tasks. There is a significant similarity between NLP and understanding biological sequences, protein sequences are made up of amino acids (e.g. M, V, S, K), which is similar to words in a sentence. Also, just like

the meaning of words depends on their context in the sentence, the role of amino acids in a protein depends on the residues around them. BERT’s architecture is designed to capture context in text and is essentially suitable for understanding context in protein sequences. Therefore, many scientists have tried to apply similar language model structures to protein sequence-related tasks.

ProteinBert [2] is one of the SOTA work in this field. In this paper, ProteinBert is pre-trained on 106M proteins on two simultaneous tasks: bidirectional language modeling of protein sequences and gene ontology (GO) annotation prediction, the pre-training makes the model easily transferable between various protein-related tasks through downstream data sets. In addition, the transformer architecture of proteinBERT employs a self-attention mechanism to capture dependencies between amino acids in a protein sequence. Protein sequences often have long-distance dependencies due to their complex spatial structures, and ProteinBERT is able to capture these implicit relationships through its attention mechanism.



(a) Model structure for ProteinBERT [2].



(b) Structure of multi-head attention layer and scaled dot-product attention.

Figure 2: Model architecture and mechanism

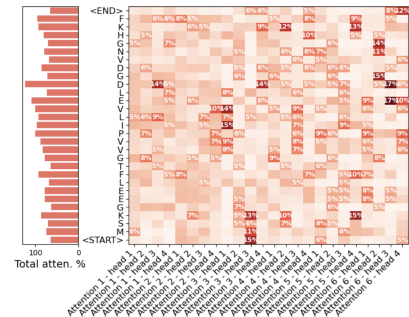
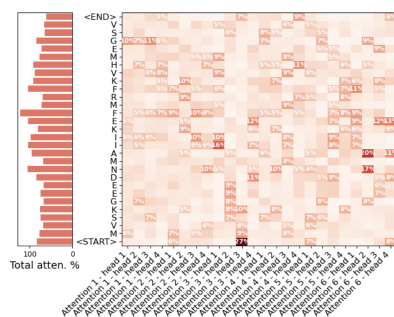


Figure 3: Visualization of attention mechanism in protein sequences

Although inspired by BERT, proteinBERT’s architecture is different and includes several innovations. ProteinBERT is a denoising autoencoder [2], the pre-trained model can be directly used for a variety of tasks, but due to the purpose of our project and the complexity of the input, we only use proteinBERT’s encoder here in our project. Input a protein sequence, and the model outputs a contextual representation of the amino acid (local embedding) or a feature representation of the entire sequence (global embedding). ProteinBERT’s design allows it to learn biologically meaningful patterns from sequences and adapt to a variety of protein-related tasks.

In our project, when we input the protein sequence (1,), the encoder will output an array in shape (1, *seq_len*), where *seq_len* is a hyper-parameter. Below is an example of input and output of model encoder.

Sampled DataFrame:

Encoded Sequence Array:

```
0    MSKGEELFTGVVPILVELDG...    array([[23, 10, 15, ..., 25, 25, 25],
1    MSKGEELFTGVVPILVELDG...    [23, 10, 15, ..., 25, 25, 25],
2    MSKGEELFTGVVPILVELDG...    [23, 10, 15, ..., 25, 25, 25],
3    MSKGEELFTGVVPILVELDG...    ...,
4    MSKGEELFTGVVPILVELDG...    [23, 10, 18, ..., 25, 25, 25],
...    [23, 10, 18, ..., 25, 25, 25],
2940 MVSELKENMPMKLYMEGTVN...    [23, 10, 18, ..., 25, 25, 25]],
2941 MVSELKENMPMKLYMEGTVN...
2942 MVSELKENMPMKLYMEGTVN...    dtype=int32)
2943 MVSELKENMPMKLYMEGTVN...
2944 MVSELKENMPMKLYMEGTVN...
Name: Sequence, Length: 2945, dtype: object
```

Secondary Structure Annotation

The proteinBERT provided a benchmark for secondary structure prediction, however, we found that their task was trained by various of short sequences, while our sequences are too long for this model. Therefore, for better performance and higher accuracy in the next step, we choose another two methods for this task.

We implemented secondary structure prediction with both the PSIPRED [3] and AlphaFold [16] method. Both of the methods have their goods and limitations. PSIPRED is a tool specifically designed to predict the secondary structure of proteins, relying primarily on lightweight models based on neural networks. The input sequence length and compute resources affect its run time, but it is usually done at the minute level and can be written to the code pipeline using APIs. AlphaFold is a 3D structure prediction tool that provides accurate protein models. While generating these models could require significant computational resources and time, AlphaFold could deliver results with high level of accuracy. The output from AlphaFold is provided in .cif format. To extract secondary structure information, the DSSP (Dictionary of Secondary Structure of Proteins) algorithm can be applied, which analyzes the structure to determine elements such as alpha helices, beta sheets, and other structural features.

The output secondary structures of these models are in string format, so we need to encode this sequence. Since all elements of secondary structure are letters representing structure categorical information, we choose to use one-hot encoding in this step. Following is a part of GFP protein’s secondary structure prediction output.

MSKGEELFTGVVPILVELDGDVNGHKFSVSGEGDATYGKLT...
sample protein sequence

CCCCHHHHCCEEEEEEEEEEEEEEECCCEEEEEEECCCCC...
sample secondary structure output

Sequence Stability Calculation

We suspect that when a position on the protein sequence is a split protein, the splitted protein is more likely to reassemble because the two subsequent strands generated by the split may have a lower stability state. With this in mind, we found a stability dataset for proteins and used it to fine-tune proteinBERT so that the model can predict the approximate stability of arbitrary protein sequences. Each split site is divided into two sequences around the site, and the stability of the two sub-sequences is predicted respectively, and the sum is used as the final split site prediction feature. Due to computational complexity, this feature can be chosen not to be used in long sequence prediction scenarios.

2.4 Models and Parameters

An important feature of ProteinBERT is sequence length flexibility. To avoid the risk of overfitting the model to a particular constant length, we found the method by switching the coding length of the protein sequence, using 128, 512, or 1024 labels. The most appropriate $seq = 512$. The loss function minimized by ProteinBERT during pretraining

$$\text{Loss} = - \sum_{i=1}^l \log(\hat{S}_{i,S_i}) - \sum_{j=1}^{8943} (A_j \cdot \log(\hat{A}_j) + (1 - A_j) \cdot \log(1 - \hat{A}_j))$$

where l is the sequence length, $S_i \in \{1, \dots, 26\}$ is the sequence's true token at position i , $\hat{S}_{i,k} \in [0, 1]$: the predicted probability that position i has the token k (for any $k \in \{1, \dots, 26\}$), $A_j \in \{0, 1\}$ is the true indicator for annotation j (for any $j \in \{1, \dots, 8943\}$), $\hat{A}_j \in [0, 1]$ is the predicted probability that the protein has annotation j .

And the the loss function minimized by ProteinBERT that we use during fine-tuning on stability

$$\text{Loss} = \frac{1}{N} \sum_{i=1}^N (y_i - \hat{y}_i)^2$$

where N is the sample size, y_i is the true label of protein stability, and \hat{y}_i is the predicted label of protein stability. As for our final MLP classifier to define whether the site can be a split site or not, we use log loss function and ADAM optimizer, where

$$\text{Log Loss} = -\frac{1}{N} \sum_{i=1}^N \sum_{j=1}^M y_{ij} \log(\hat{y}_{ij})$$

So, through the pipeline, we fine-tuned the proteinBERT and trained the MLP classifier, to get a higher performance, data standardization and dimensional reduction were implemented before the features were added to the classifier. The parameter details of models are here below:

Param Name	Task	Default Value
seq_len	Fine-tune the maximum length of the input sequence.	512
batch_size	Fine-tune batch size for each training.	32
max_epochs_per_stage	Fine-tune maximum number of rounds for freezing layer training.	40
lr	Fine-tune learning rate (for unfrozen model layers).	$1e^{-4}$
max_iter	Max MLP classifier training iteration.	500
seq_len	Length of encoded sequence to input MLP.	512
hidden_layer_sizes	MLP hidden layer size.	(128, 64)
test_size	Test size for MLP training.	0.2
pca_components	The dimensions added to the model were reduced using PCA.	80

Table 2: Parameter description and default values for models.

3 Results

Stability prediction task

Although the loss function is defined above, since the label is numerical but not binary and we can not evaluate it by metrics like accuracy, we use Spearman's rank correlation here to show how our model is, where $\rho = 1 - \frac{6 \sum d_i^2}{n(n^2-1)}$. Spearman's rank correlation coefficient is a rank correlation coefficient that measures a monotone relationship between two variables, rather than just a linear relationship. This means that even if the relationship between variables is not strictly linear, the Spearman coefficient captures this trend change well. From the figure below, we can see the correlation is very high and close to 1, so the model is definitely performing well.

Split site prediction task

Due to the sparsity and unbalance of the data, running the model classification directly may lead to serious overfitting problems, so we used sampling technology before running, respectively fitting the classifier on the

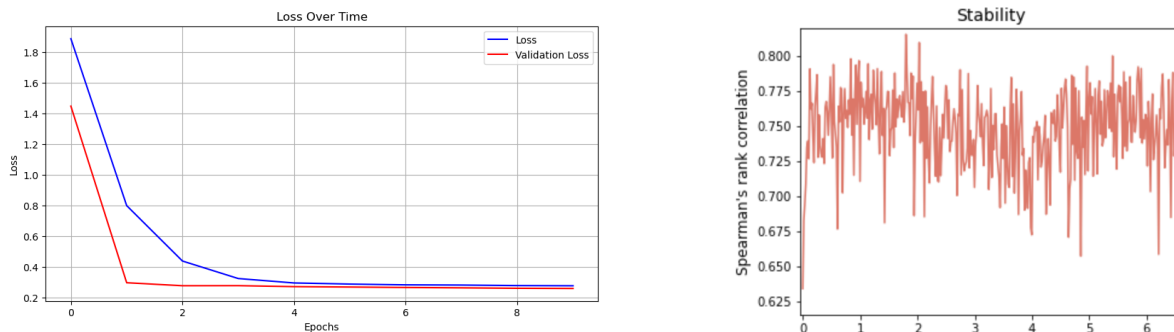


Figure 4: evaluation for stability prediction

unsampled data, resampled data and undersampled data. Unsampled is used as the raw data, and the over-sampled method copies or generates new data points to increase the sampled data volume. down-sampled decreases the number of most samples, and randomly removes most samples to approximate the sampled data volume to a few. Both of these approaches can solve the problem of raw data imbalance. In addition, we also conducted training on data that did not contain stability features, and the results obtained are shown in the following table 3 and 4.

Dataset	Class	Precision	Recall	F1-score	Support	Accuracy	Macro Avg F1
unsampled	0	0.98	1.00	0.99	795	0.98	0.49
	1	0.00	0.00	0.00	17		
over-sampled	0	1.00	0.84	0.91	794	0.92	0.92
	1	0.87	1.00	0.93	810		
down-sampled	0	0.88	0.64	0.74	11	0.76	0.77
	1	0.69	0.90	0.78	10		

Table 3: Classification Report Comparison for Datasets without stability feature.

Dataset	Class	Precision	Recall	F1-score	Support	Accuracy	Macro Avg F1
unsampled	0	0.97	1.00	0.99	573	0.97	0.49
	1	0.00	0.00	0.00	16		
over-sampled	0	1.00	0.90	0.95	579	0.96	0.95
	1	0.92	1.00	0.96	677		
down-sampled	0	0.40	0.50	0.44	4	0.71	0.62
	1	0.83	0.77	0.80	13		

Table 4: Classification Report Comparison for Datasets with stability

As can be seen from the table, it is obvious that the model shows overfitting phenomenon on the unsampled data set. Finally, we used the model fitted on the undersampled data set with stability characteristics to evaluate the test set, and the prediction accuracy on 20 pieces of data was about 75%, indicating that it was an accurate model.

4 Discussion

Result Analysis It can be seen from the results that we can use a large language model to encode proteins, and after combining spatial structure, stability, and other features, the model can have a good performance for splitting site prediction of FPs.

Limitation of the Project One major limitation of the project is the small size of the dataset. A larger and more diverse dataset could have been generated to improve the robustness and generalizability of the results. While gathering data from literature could be time-consuming, which more time and effort beyond the scope of this project. The second limitation is that, in real-world scenarios, sequences may include overlapping regions or deleted regions that can recombine, but these aspects should have been accounted for in this study, potentially

limiting the applicability of the findings. Furthermore, there is a need for additional features to be studied. For example, the stability of the sequences has not been studied well from a biochemistry perspective, which requires more investigation into this important aspect.

Possible Directions Future directions based on current work will focus not just on FPs, but also other protein types. A broader dataset would facilitate a more comprehensive understanding in complementary studies of proteins. Additionally, exploration of additional features and parameters, such as those impacting split fragments stability will be included for more perspective of analysis. Moreover, as the data were all collected from experiments which could have variations from experimental design. More standardization of the data source could also be done.

5 Conclusion

In this work, we developed and validated a computational pipeline by leveraging a protein language model called ProteinBERT combined with MLP for accurate prediction of optimal splitting sites in FPs. The ProteinBERT model was pretrained with stability dataset, to encode the protein sequence. Secondary structure annotation was added as another feature into the MLP classifier to provide a robust framework for splitting site recognition. Our results first demonstrate the feasibility of the model with a high Spearman’s rank correlation close to 1. As for splitting task prediction, we chose to use oversampled and undersampled dataset to solve the overfitting problems of the raw data and oversampled data had the best performance while the undersampled data also had a good level of accuracy. These strategies improved the performance in terms of the overfitting problems and enhance generalizability of the model under a small data size. Our findings indicated the potential of the pipeline for predicting more protein splitting mechanisms for advancement in biomolecular complementation assay. When combined with other benchmarks, our model potentially can be integrated into splitting tasks of other proteins like gene editing effectors, leading to more orthogonal and flexible biomolecular tools in biomedical research. In the future, with expanded dataset in other protein types, higher accuracy and performance of our pipeline can be expected.

Code Availability

Our code repository can be found in <https://github.com/Carrie1013/DS596-ProteinProject>

References

- [1] Francesca Anson, Pintu Kanjilal, S. Thayumanavan, and Jeanne A. Hardy. Tracking exogenous intracellular casp-3 using split GFP. *Protein Science*, 30(2):366–380, February 2021.
- [2] N. Brandes, D. Ofer, Y. Peleg, N. Rappoport, and M. Linial. Proteinbert: a universal deep-learning model of protein sequence and function. *Bioinformatics*, 38(8):2102–2110, 2022.
- [3] D. W. A. Buchan and D. T. Jones. The psipred protein analysis workbench: 20 years on. *Nucleic Acids Research*, 47(W1):W402–W407, 2019.
- [4] Stéphanie Cabantous, Hau B. Nguyen, Jean-Denis Pedelacq, Faten Koraïchi, Anu Chaudhary, Kumkum Ganguly, Meghan A. Lockard, Gilles Favre, Thomas C. Terwilliger, and Geoffrey S. Waldo. A new protein-protein interaction sensor based on tripartite split-gfp association. *Scientific Reports*, 3(1):2854, October 2013.
- [5] Stéphanie Cabantous, Thomas C. Terwilliger, and Geoffrey S. Waldo. Protein tagging and detection with engineered self-assembling fragments of green fluorescent protein. *Nature Biotechnology*, 23(1):102–107, January 2005.
- [6] Brian P. Callahan, Matthew J. Stanger, and Marlene Belfort. Cut and glow: Protease activation of split green fluorescent protein. *Chembiochem : a European journal of chemical biology*, 11(16):2259–2263, November 2010.
- [7] Sam Duwé and Peter Dedecker. Optimizing the fluorescent protein toolbox and its use. *Current opinion in biotechnology*, 58:183–191, 2019.

- [8] Jin-Yu Fan, Zong-Qiang Cui, Hong-Ping Wei, Zhi-Ping Zhang, Ya-Feng Zhou, Yun-Peng Wang, and Xian-En Zhang. Split mCherry as a new red bimolecular fluorescence complementation system for visualizing protein–protein interactions in living cells. *Biochemical and Biophysical Research Communications*, 367(1):47–53, February 2008.
- [9] Siyu Feng, Sayaka Sekine, Veronica Pessino, Han Li, Manuel D. Leonetti, and Bo Huang. Improved split fluorescent proteins for endogenous protein labeling. *Nature Communications*, 8(1):370, August 2017.
- [10] Siyu Feng, Sayaka Sekine, Veronica Pessino, Han Li, Manuel D. Leonetti, and Bo Huang. Improved split fluorescent proteins for endogenous protein labeling. *Nature Communications*, 8(1):370, August 2017.
- [11] Indraneel Ghosh, Andrew D. Hamilton, and Lynne Regan. Antiparallel leucine zipper-directed protein reassembly: application to the green fluorescent protein. *Journal of the American Chemical Society*, 122(23):5658–5659, June 2000.
- [12] Fabian Hertel, Gary C. H. Mo, Sam Duwé, Peter Dedecker, and Jin Zhang. Refsofi for mapping nanoscale organization of protein-protein interactions in living cells. *Cell reports*, 14(2):390–400, January 2016.
- [13] Chang-Deng Hu, Yurii Chinenov, and Tom K. Kerppola. Visualization of interactions among bzip and rel family proteins in living cells using bimolecular fluorescence complementation. *Molecular Cell*, 9(4):789–798, April 2002.
- [14] Chang-Deng Hu and Tom K. Kerppola. Simultaneous visualization of multiple protein interactions in living cells using multicolor fluorescence complementation analysis. *Nature Biotechnology*, 21(5):539–545, May 2003.
- [15] Chang-Deng Hu and Tom K. Kerppola. Simultaneous visualization of multiple protein interactions in living cells using multicolor fluorescence complementation analysis. *Nature biotechnology*, 21(5):539–545, May 2003.
- [16] J. M. Jumper et al. Accurate structure prediction of biomolecular interactions with alphafold 3. *Nature*, 2024.
- [17] Yuriko Kakimoto, Shinya Tashiro, Rieko Kojima, Yuki Morozumi, Toshiya Endo, and Yasushi Tamura. Visualizing multiple inter-organellar contact sites using the organellar-targeted split-GFP system. *Scientific Reports*, 8(1):6175, April 2018.
- [18] Daichi Kamiyama, Sayaka Sekine, Benjamin Barsi-Rhine, Jeffrey Hu, Baohui Chen, Luke A. Gilbert, Hiroaki Ishikawa, Manuel D. Leonetti, Wallace F. Marshall, Jonathan S. Weissman, and Bo Huang. Versatile protein tagging in cells with split fluorescent protein. *Nature Communications*, 7(1):11046, March 2016.
- [19] Takaaki Kojima, Satoshi Karasawa, Atsushi Miyawaki, Takeshi Tsumuraya, and Ikuo Fujii. Novel screening system for protein–protein interactions by bimolecular fluorescence complementation in *Saccharomyces cerevisiae*. *Journal of Bioscience and Bioengineering*, 111(4):397–401, April 2011.
- [20] L. A. Kost, E. V. Putintseva, A. R. Pereverzeva, D. M. Chudakov, K. A. Lukyanov, and A. M. Bogdanov. Bimolecular fluorescence complementation based on the red fluorescent protein FusionRed. *Russian Journal of Bioorganic Chemistry*, 42(6):619–623, November 2016.
- [21] Tuğba Köker, Anthony Fernandez, and Fabien Pinaud. Characterization of split fluorescent protein variants and quantitative analyses of their self-assembly process. *Scientific Reports*, 8(1):5344, March 2018.
- [22] Tuğba Köker, Nathalie Tang, Chao Tian, Wei Zhang, Xueding Wang, Richard Martel, and Fabien Pinaud. Cellular imaging by targeted assembly of hot-spot SERS and photoacoustic nanoprobe using split-fluorescent protein scaffolds. *Nature Communications*, 9(1):607, February 2018.
- [23] You Ri Lee, Jong-Hwa Park, Soo-Hyun Hahm, Lin-Woo Kang, Ji Hyung Chung, Ki-Hyun Nam, Kwang Yeon Hwang, Ick Chan Kwon, and Ye Sun Han. Development of bimolecular fluorescence complementation using dronpa for visualization of protein–protein interactions in cells. *Molecular Imaging and Biology*, 12(5):468–478, October 2010.
- [24] Magnus Lundqvist, Niklas Thalén, Anna-Luisa Volk, Henning Gram Hansen, Eric von Otter, Per-Åke Nygren, Mathias Uhlen, and Johan Rockberg. Chromophore pre-maturation for improved speed and sensitivity of split-GFP monitoring of protein secretion. *Scientific Reports*, 9(1):310, January 2019.

- [25] Cécile Polge, Stéphanie Cabantous, and Daniel Taillandier. Tripartite split-gfp for high throughput screening of small molecules: a powerful strategy for targeting transient/labile interactors like e2-e3 ubiquitination enzymes. *ChemBioChem*, 25(6):e202300723, March 2024.
- [26] Houming Ren, Qingshan Ou, Qian Pu, Yuqi Lou, Xiaolin Yang, Yujiao Han, and Shiping Liu. Comprehensive review on bimolecular fluorescence complementation and its application in deciphering protein–protein interactions in cell signaling pathways. *Biomolecules*, 14(7):859, 2024.
- [27] Daria M. Shcherbakova, Mikhail Baloban, Alexander V. Emelyanov, Michael Brenowitz, Peng Guo, and Vladislav V. Verkhusha. Bright monomeric near-infrared fluorescent proteins as tags and biosensors for multiscale imaging. *Nature Communications*, 7(1):12405, August 2016.
- [28] Daria M. Shcherbakova, Mikhail Baloban, Alexander V. Emelyanov, Michael Brenowitz, Peng Guo, and Vladislav V. Verkhusha. Bright monomeric near-infrared fluorescent proteins as tags and biosensors for multiscale imaging. *Nature Communications*, 7(1):12405, August 2016.
- [29] Y. John Shyu, Han Liu, Xuehong Deng, and Chang-Deng Hu. Identification of new fluorescent protein fragments for bimolecular fluorescence complementation analysis under physiological conditions. *BioTechniques*, 40(1):61–66, January 2006.
- [30] Michael Walter, Christina Chaban, Katia Schütze, Oliver Batistic, Katrin Weckermann, Christian Näke, Dragica Blazevic, Christopher Grefen, Karin Schumacher, Claudia Oecking, Klaus Harter, and Jörg Kudla. Visualization of protein interactions in living plant cells using bimolecular fluorescence complementation. *The Plant Journal*, 40(3):428–438, November 2004.
- [31] Qingyan Yuan, Minhui Wu, Yibo Liao, Shuli Liang, Yuan Lu, and Ying Lin. Rapid prototyping enzyme homologs to improve titer of nicotinamide mononucleotide using a strategy combining cell-free protein synthesis with split GFP. *Biotechnology and Bioengineering*, 120(4):1133–1146, April 2023.