

REALM: Region-Empowered Antibody Language Model for Antibody Property Prediction

Toru Nishino

Bio Science & Engineering Laboratories Bio Science & Engineering Laboratories Bio Science & Engineering Laboratories
FUJIFILM Corporation FUJIFILM Corporation FUJIFILM Corporation
Tokyo, Japan Tokyo, Japan Kaisei, Japan
toru.nishino@fujifilm.com noriji.kato@fujifilm.com takuya.tsutaoka@fujifilm.com

Noriji Kato

Takuya Tsutaoka

Yuanzhong Li

Bio Science & Engineering Laboratories
FUJIFILM Corporation
Kaisei, Japan
yuanzhong.li@fujifilm.com

Masahito Ohue

School of Computing
Institute of Science Tokyo
Yokohama, Japan
ohue@c.titech.ac.jp

Abstract— Protein language models (pLM) are beneficial to build antibody property prediction models. However, current pLMs lack the ability to understand antibody properties because region and structure information is not effectively embedded. We propose the Region-Empowered Antibody Language Model (REALM), a pLM built by multi-task pretraining strategy of residue prediction and region prediction tasks in antibodies, to incorporate not only co-evolution but also region information of antibodies. We demonstrate that our REALM improves the understanding of antibody properties, including hydrophobicity and thermo-stability.

Index Terms—Antibody Property Prediction, Protein Language Model, Property Prediction, Biopharmaceutical

I. INTRODUCTION

To reduce the manufacturing costs of antibody drugs, it is crucial to predict physicochemical properties such as hydrophobicity and thermo-stability from antibody sequences. Recent emerging protein language models (pLMs) are beneficial for predicting antibody properties. As the pLM is pretrained on a large amount of protein sequence, the model can predict antibody properties even though it is fine-tuned using only a small amount of data regarding the target task.

However, current pLMs face challenges due to their focus on learning antibody co-evolution rather than antibody properties. Therefore, although existing pLMs are useful for understanding the mutation effects, they are still insufficient for understanding antibody properties.

This study aims to build a more effective pLM to understand antibody properties by using region information of antibodies. The region information of antibodies, including loops and turns, significantly influences their properties, so the accuracy of the antibody property prediction improves.

II. METHOD

We propose Region-Empowered Antibody Language Model (REALM), which embeds not only antibody sequences but

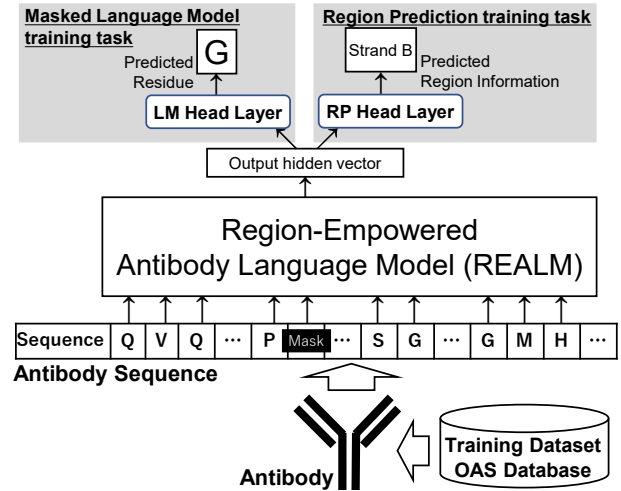


Fig. 1. An overview of REALM pretraining. REALM is pretrained with two tasks: masked language model and region prediction tasks.

also region information into the pLM. Fig. 1 shows an overview of our pretraining strategy. We employ a multi-task pretraining to embed both the amino acid residues in the antibody sequence and the region information they belong to.

REALM uses two language model head layers; language model head layer (LM Head) and region prediction head layer (RP Head). We employed the region prediction task as an auxiliary task for embedding region information of the residues in the antibody sequence. The LM Head outputs the probabilities of the masked residues, and the RP Head outputs the probabilities of the region each residue belongs in the input antibody sequence. The total loss is the sum of the losses for both the residue prediction task and the region prediction task.

We use nine strands (from strand A to strand G) and three complementarity determining regions (from CDR1 to

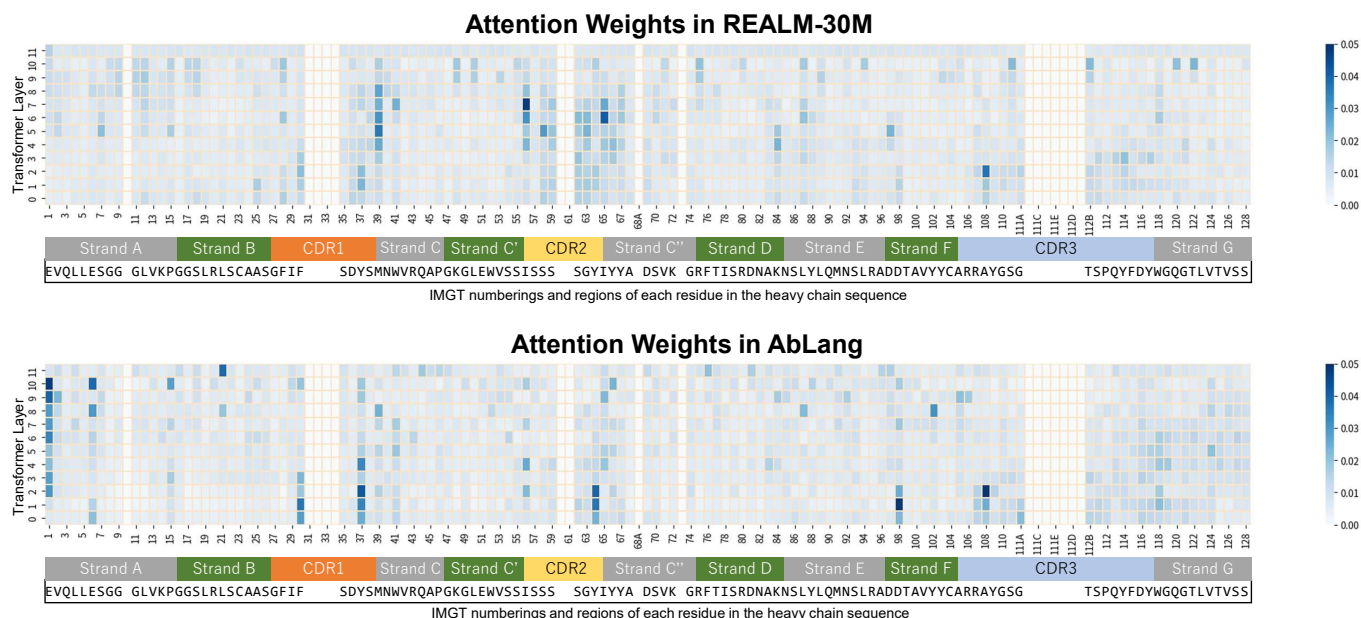


Fig. 3. An example of normalized attention weights in our REALM and AbLang [10] models. The horizontal axis shows the number of each residue and the region it belongs to in the input antibody sequence in the IMGT numbering system, and the vertical axis indicates the layers of the transformer layer in the antibody language model. The attention weights are all normalized along with the x-axis, that is, the residues of the input sequence. The positions where amino acid residues are missing correspond to the missing numbers in IMGT’s unique numbering system.

task of region prediction, pays more attention to the region boundaries, and in turn, the region information is embedded in the model.

In AbLang, the weight of attention for the first token is higher. This is because AbLang uses absolute position embedding; therefore, it is necessary to measure the relative distance to the first token in order to determine the relative position of the residues within the antibody. In REALM, the attention weights of the tokens in the middle of the sequence are higher than those in the first token. This is due to the application of RoPE positional embedding [5]. In antibody sequences, the middle part of the sequence, especially the area around the CDR, is more important for understanding antibody characteristics than the end of the sequence. Therefore, we assume that our REALM is able to focus on the important parts of antibody sequences more appropriately than AbLang.

VI. CONCLUSION

We propose a Region-Empowered Antibody Language Model (REALM) that uses multi-task learning with token prediction and the region prediction task. The evaluation results showed that REALM improves the accuracy of the two assays, hydrophobicity and thermal stability. This result and analysis of attention weights demonstrate that the region information is embedded effectively. In the future, we will additionally combine the pre-training task to embed physico-chemical information to our antibody language model.

REFERENCES

[1] M.-P. Lefranc, V. Giudicelli, C. Ginestoux, J. Bodmer, W. Müller, R. Bontrop, M. Lemaître, A. Malik, V. Barbié, and D. Chaume, “IMGT,

the international ImMunoGeneTics database,” *Nucleic Acids Research*, vol. 27, no. 1, pp. 209–212, 1999.

[2] M. Ruiz and M.-P. Lefranc, “IMGT gene identification and Colliers de Perles of human immunoglobulins with known 3D structures,” *Immunogenetics*, vol. 53, pp. 857–883, 2002.

[3] T. H. Olsen, F. Boyles, and C. M. Deane, “Observed antibody space: A diverse database of cleaned, annotated, and translated unpaired and paired antibody sequences,” *Protein Science*, vol. 31, no. 1, pp. 141–146, 2022.

[4] L. Shehata, D. P. Maurer, A. Z. Wec, A. Lilov, E. Champney, T. Sun, K. Archambault, I. Burnina, H. Lynaugh, X. Zhi *et al.*, “Affinity maturation enhances antibody specificity but compromises conformational stability,” *Cell Reports*, vol. 28, no. 13, pp. 3300–3308, 2019.

[5] J. Su, M. Ahmed, Y. Lu, S. Pan, W. Bo, and Y. Liu, “Roformer: Enhanced transformer with rotary position embedding,” *Neurocomputing*, vol. 568, p. 127063, 2024.

[6] N. Shazeer, “GLU variants improve transformer,” *arXiv preprint arXiv:2002.05202*, 2020.

[7] B. Zhang and R. Sennrich, “Root mean square layer normalization,” *Advances in Neural Information Processing Systems*, vol. 32, 2019.

[8] R. Tibshirani, “Regression shrinkage and selection via the lasso,” *Journal of the Royal Statistical Society Series B: Statistical Methodology*, vol. 58, no. 1, pp. 267–288, 1996.

[9] T. Akiba, S. Sano, T. Yanase, T. Ohta, and M. Koyama, “Optuna: A next-generation hyperparameter optimization framework,” in *Proceedings of the 25th ACM SIGKDD international conference on Knowledge Discovery & Data mining*, 2019, pp. 2623–2631.

[10] T. H. Olsen, I. H. Moal, and C. M. Deane, “AbLang: an antibody language model for completing antibody sequences,” *Bioinformatics Advances*, vol. 2, no. 1, p. vbac046, 2022.

[11] H. Jing, Z. Gao, S. Xu, T. Shen, Z. Peng, S. He, T. You, S. Ye, W. Lin, and S. Sun, “Accurate prediction of antibody function and structure using bio-inspired antibody language model,” *Briefings in Bioinformatics*, vol. 25, no. 4, p. bbae245, 2024.