

A Deep Learning Model for Predicting the Thermal Stability of Collagen Triple Helices

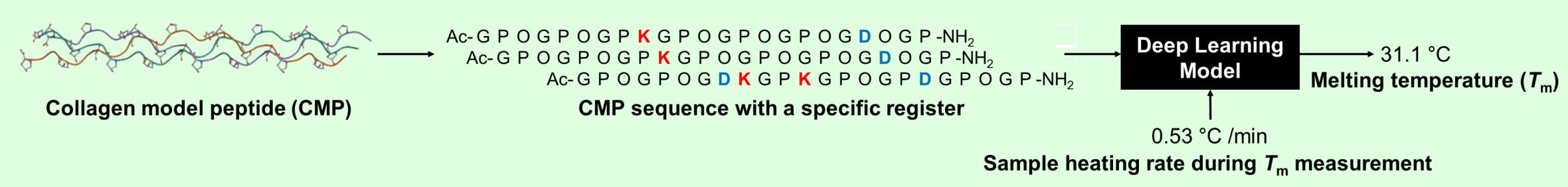
Kiseop Im, Ian Warm, Tomas Fiala,* and Helma Wennemers*

Laboratory of Organic Chemistry, ETH Zürich, D-CHAB, Vladimir-Prelog-Weg 3, 8093 Zürich, Switzerland.

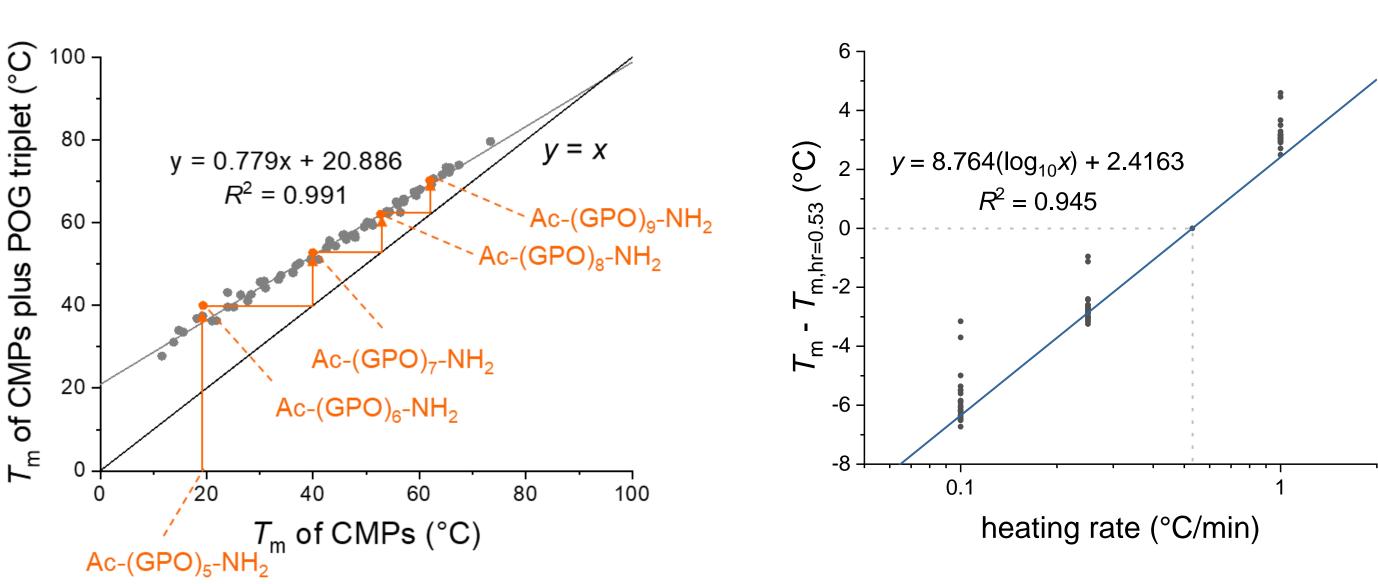


1. Introduction

Collagen model peptides (CMPs) are key tools for investigating the physico-chemical properties of collagen, the most abundant mammalian protein, as well as designing synthetic collagenbased materials and probes. Like collagen, CMPs self-assemble into a triple-helical structure that is characterized by its melting temperature (T_m). The T_m is highly dependent on the CMP sequence.^{1,4} Because of the colossal number of possible sequence variations, computational methods are necessary for deciphering the trends in the thermal stability of collagen triple helices. Despite several attempts to predict T_m values from CMP sequences, 5-7 no computational models integrate all key aspects that influence the T_m to achieve accurate predictions. Here, we design a deep-learning model that dissects the pairwise intra- and interstrand interactions in CMPs and uses them to accurately predict T_m values from CMP sequences. Our approach comprehensively deals with critical aspects affecting the experimental $T_{\rm m}$ such as frame shifts, appling groups, and heating rates.



2. Considered Trends in $T_{\rm m}$



- Any sequence modifications or pairwise-interaction additions change the T_m of the CMP.
- The change in $T_{\rm m}$ linearly depends on the basal $T_{\rm m}$ of the parental sequence but is not affected by the existing sequence or pairwise interactions.

 $4,694 \text{ CMP-} T_{\text{m}} \text{ pairs}$

 $T_{\rm m}$ increases with a higher heating rate logarithmically

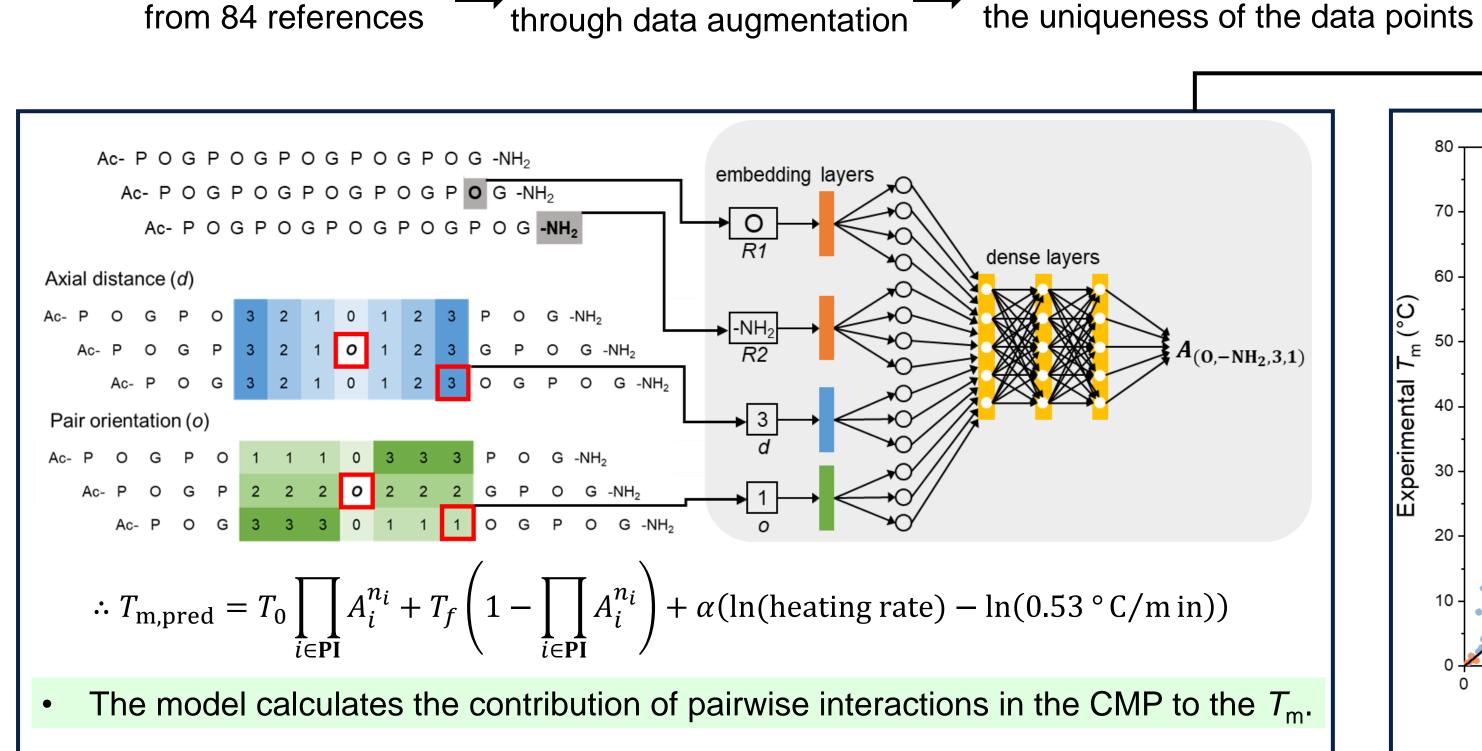
$T_{\rm m}$ (°C) Sequence →Existing sequence-T_m pairs H-GPOGPOGPOGPOGPOGPO-OH 27.8 c-GPOGPOGPOGPOGPOGPO-OH 35.6 Augmented pairs for decoupling c-GPOGPOGPOGPOGPOGPO-NH2 52.7 c-GPOGPOGPOGPOGPOGPOG-NH2 45.5 Ac-GPOGPOGPOGDRGPOGPOGPOGG-NH 37.1 Using previously identified general trends, 8,9 3,695 more sequence- $T_{\rm m}$ $R^2 = 0.994$ pairs were mathematically generated. 214 augmented data cross-validated with existing experimental data $T_{\rm m}$ of augmented data (°C)

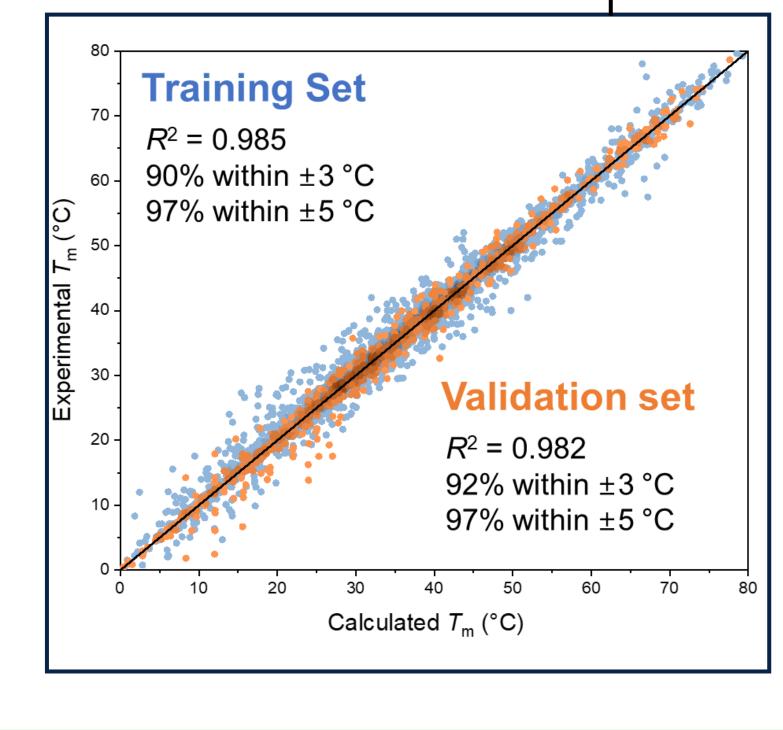
→ deep-learning model construction → training and optimization → model testing

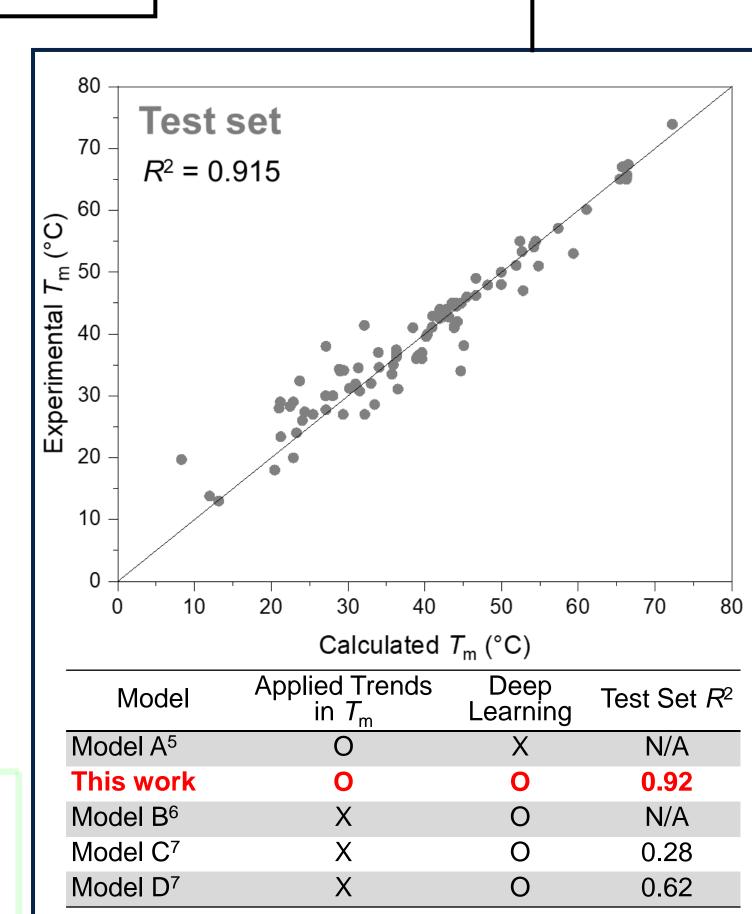
3. Model Training and Testing

1,099 CMP- $T_{\rm m}$ pairs ____

from 84 references







4. Conclusion

- We constructed a deep learning model that considers key trends in collagen $T_{\rm m}$ and successfully trained the model with reported and augmented data.
- The deep learning model achieves unprecedented accuracy and generalization in predicting collagen $T_{\rm m}$ values.
- Significance: our model represents a general tool that will contribute to a better understanding of CMP sequence-property relationships and aid the design of collagen triple helices with desired stability.

References

1. Shoulders, M. D.; Raines, R. T., Collagen Structure and Stability. Annu. Rev. Biochem. 2009, 78 (1), 929-958. 2. Okuyama, K., Miyama, K., Mizuno, K., & Bächinger, H. P. Crystal structure of (Gly-Pro-Hyp) 9: Implications for the collagen

the mechanical behavior of biomedical materials. 2022, 125, 104921.

Terminal Residues and Caps. Angew. Chem., Int. Ed. 2023, 62 (3), e202214728

- molecular model. *Biopolymers*, **2012**, *97*(8), 607-616. 3. Privalov, P. Stability of proteins: proteins which do not present a single cooperative system. Advances in protein chemistry. 1982,
- 4. Persikov, A. V., Ramshaw, J. A., Kirkpatrick, A. & Brodsky, B. Amino acid propensities for the collagen triple-helix. Biochemistry.
- 5. Walker, D. R. et al. Predicting the stability of homotrimeric and heterotrimeric collagen helices. Nature Chemistry. 2021, 13, 260-
- 6. Yu, C.-H. et al. ColGen: An end-to-end deep learning model to predict thermal stability of de novo collagen sequences. Journal of
- 7. Khare, E., Gonzalez-Obeso, C., Kaplan, D. L. & Buehler, M. J. CollagenTransformer: end-to-end transformer model to predict thermal stability of collagen triple helices using an NLP approach. ACS Biomaterials Science & Engineering. 2022, 8, 4301-4310.
- 8. Fiala, T.; Barros, E. P.; Ebert, M.-O.; Ruijsenaars, E.; Riniker, S.; Wennemers, H., Frame Shifts Affect the Stability of Collagen Triple Helices. J. Am. Chem. Soc. 2022, 144 (40), 18642-18649. 9. Fiala, T.; Barros, E. P.; Heeb, R.; Riniker, S.; Wennemers, H., Predicting Collagen Triple Helix Stability through Additive Effects of





dataset split according to