

## Petabase-scale Homology Search for Structure Prediction

Sewon Lee<sup>1,\*</sup>, Gyuri Kim<sup>1,\*</sup>, Eli Levy Karin<sup>2</sup>, Milot Mirdita<sup>1</sup>, Sukhwan Park<sup>3</sup>, Rayan Chikhi<sup>4</sup>,  
Artem Babaian<sup>5,6</sup>, Andriy Kryshchak<sup>7</sup>, and Martin Steinegger<sup>1,3,8,9,†</sup>

<sup>1</sup>*School of Biological Sciences, Seoul National University*

<sup>2</sup>*ELKMO*

<sup>3</sup>*Interdisciplinary Program in Bioinformatics, Seoul National University*

<sup>4</sup>*Institut Pasteur, Université Paris Cité, G5 Sequence Bioinformatics*

<sup>5</sup>*Department of Molecular Genetics, University of Toronto, Toronto*

<sup>6</sup>*Donnelly Centre for Cellular and Biomolecular Research, University of Toronto*

<sup>7</sup>*Genome Center, University of California, Davis*

<sup>8</sup>*Artificial Intelligence Institute, Seoul National University*

<sup>9</sup>*Institute of Molecular Biology and Genetics, Seoul National University*

*\*These authors contributed equally*

*†Corresponding author: [martin.steinegger@snu.ac.kr](mailto:martin.steinegger@snu.ac.kr)*

The CASP15 competition highlighted the critical role of multiple sequence alignments (MSAs) in protein structure prediction, as demonstrated by the success of the top AlphaFold2-based prediction methods. To push the boundaries of MSA utilization, we conducted a petabase-scale search of the Sequence Read Archive (SRA), resulting in gigabytes of aligned homologs for CASP15 targets. These were merged with default MSAs produced by ColabFold-search and provided to ColabFold-predict. By using SRA data, we achieved highly accurate predictions (GDT\_TS > 70) for 66% of the non-easy targets, whereas using ColabFold-search default MSAs scored highly in only 52%. Next, we tested the effect of deep homology search and ColabFold's advanced features, such as more recycles, on prediction accuracy. While SRA homologs were most significant for improving ColabFold's CASP15 ranking from 11th to 3rd place, other strategies contributed too. We analyze these in the context of existing strategies to improve prediction.