

# 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 Krysfafovich<sup>7</sup> and Martin Steinegger<sup>1,3,8,9,†</sup>

<sup>1</sup>School of Biological Sciences, Seoul National University, South Korea; <sup>2</sup>ELKMO, Denmark; <sup>3</sup>Interdisciplinary Program in Bioinformatics, Seoul National University, South Korea; <sup>4</sup>Institut Pasteur, Université Paris Cité, France; <sup>5</sup>Department of Molecular Genetics, <sup>6</sup>Donnelly Centre for Cellular and Biomolecular Research, University of Toronto, Canada; <sup>7</sup>Genome Center, University of California, Davis, USA; <sup>8</sup>Artificial Intelligence Institute, <sup>9</sup>Institute of Molecular Biology and Genetics, Seoul National University, South Korea  
\*Equal contributors; †Correspondence: martin.steinegger@snu.ac.kr

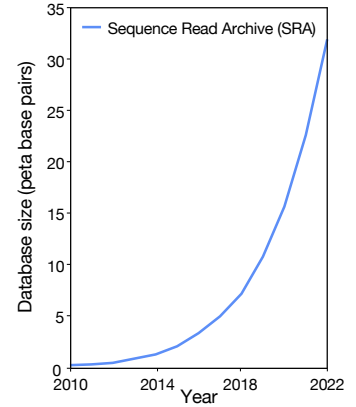
## Abstract

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.

## SRA: the largest public sequence database

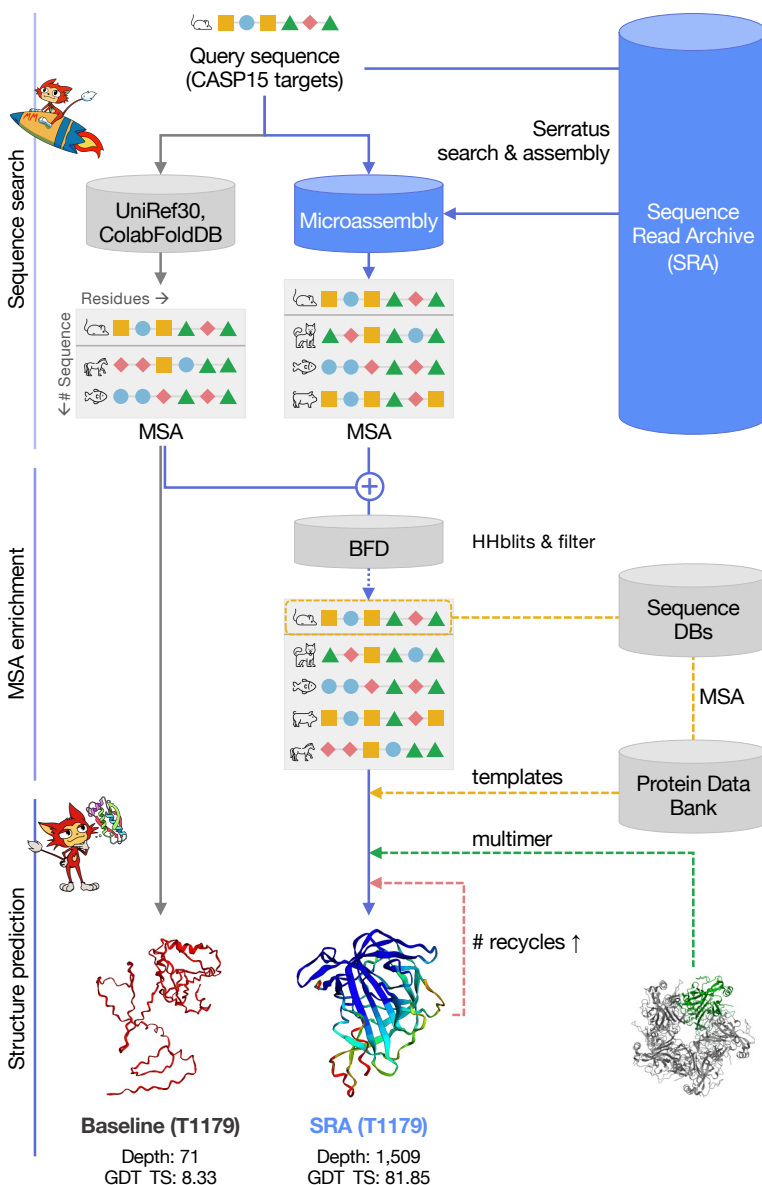
Database	Giga base pairs
<b>SRA</b>	<b>32,000,000</b>
IMG/M	27,000
MGnify*	3,500
BFD*	2,600
Metaclust*	1,700
<b>Microassembly</b>	<b>1,500</b>
ColabFoldDB*	800
UniRef30*	30

\*Estimated based on # sequences

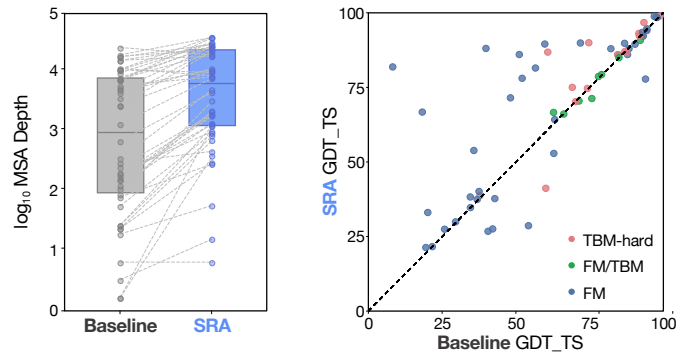


**Microassembly:** search results from SRA assembled with maviralSpades

## Deep & wide sequence search for structure prediction



## SRA-enriched MSAs lead to better predictions

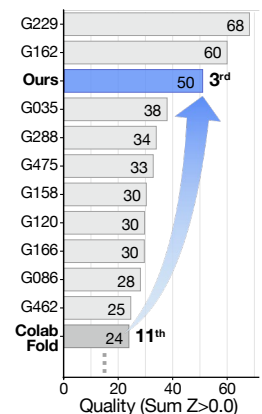


**TBM:** templated-based modeling; **FM:** free-modeling

**GDT\_TS (global distance test):** prediction accuracy

## Factors contributing to improved prediction

	GDT_TS	Sum Z (>0.0)
CFDB (Baseline)	65.8 ± 25.6	17.1
CFDB + HH	67.4 ± 25.2	15.8
<b>CFDB + SRA</b>	<b>71.0 ± 24.4</b>	<b>30.3</b>
CFDB + SRA + HH	<b>71.4 ± 23.3</b>	26.3
CFDB + rec	68.1 ± 25.3	17.7
CFDB + SRA + temp	70.9 ± 24.0	31.6
CFDB + SRA + mul	71.0 ± 24.1	31.0
<b>CFDB + SRA + rec</b>	<b>75.5 ± 22.2</b>	<b>40.6</b>
<b>Model 1</b>	<b>77.9 ± 20.6</b>	<b>50.3</b>



**Abbr.** HH: HHblits; rec: # recycles ↑; temp: templates; mul: multimer

**Sum Z:** sum of GDT\_TS (accuracy) Z-scores among server groups over 0

**Model1:** predicted-best strategies for each target

## References

Lee & Kim, et al. *bioRxiv*, 2023.07.10.548308 (2023).  
Accepted in *CSHL Perspect. Biol.*  
Edgar, et al. *Nature*, 602, 142-147 (2022).  
Jumper, et al. *Nature*, 596, 583-589 (2021).  
Krysfafovich, et al. *Proteins*, 89, 1607-1617 (2021).  
Mirdita, et al. *Nat. Methods*, 19, 679-682 (2022).  
Ovchinnikov, et al. *Science*, 355, 294-298 (2017).  
Steinegger & Söding. *Nat. Biotechnol.*, 35, 1026-1028 (2017).



Download Poster (PDF)