

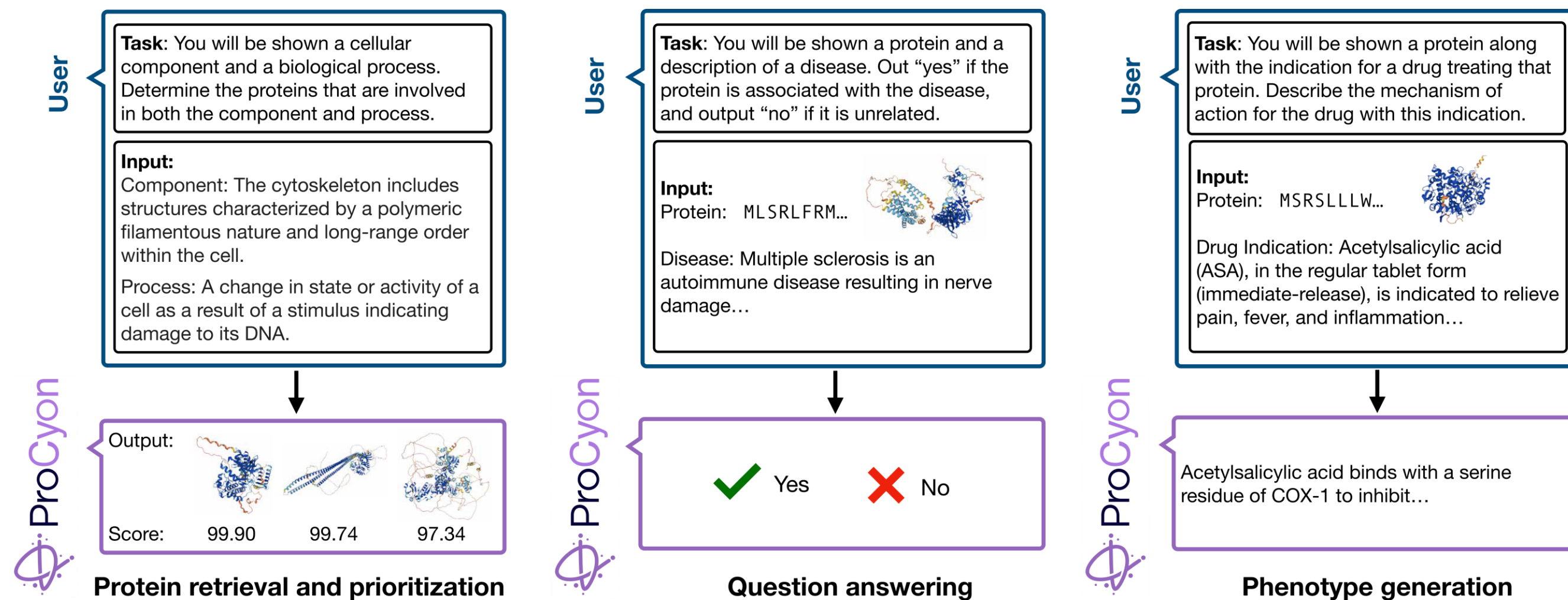


Owen Queen<sup>1,\*</sup>, Yepeng Huang<sup>1,\*</sup>, Robert Calef<sup>1,2,\*</sup>, Valentina Giunchiglia<sup>1,3,4</sup>, Tianlong Chen<sup>1,2</sup>, George Dasoulas<sup>1</sup>, LeAnn Tai<sup>2</sup>, Yasha Ektefaie<sup>1</sup>, Ayush Noori<sup>1</sup>, Joseph Brown<sup>5</sup>, Tom Cogley<sup>2,6</sup>, Karin Hrovatin<sup>7,8</sup>, Tom Hartvigsen<sup>9</sup>, Fabian J. Theis<sup>7,10</sup>, Bradley L. Pentelute<sup>5,14</sup>, Vikram Khurana<sup>11,12,14</sup>, Manolis Kellis<sup>2,14</sup>, Marinka Zitnik<sup>1,13,14,15,‡</sup>

<sup>1</sup>Department of Biomedical Informatics, Harvard Medical School; <sup>2</sup>Computer Science and Artificial Intelligence Laboratory, MIT; <sup>3</sup>Department of Brain Sciences, Imperial College London; <sup>4</sup>Centre for Neuroimaging Sciences, King's College London; <sup>5</sup>Department of Chemistry, MIT; <sup>6</sup>Department of Computing, Imperial College London; <sup>7</sup>Institute of Computational Biology, Computational Health Center, Helmholtz Munich; <sup>8</sup>TUM School of Life Sciences Weihenstephan, Technical University of Munich; <sup>9</sup>School of Data Science, University of Virginia; <sup>10</sup>School of Computation, Information and Technology, Technical University of Munich; <sup>11</sup>Department of Neurology, Brigham and Women's Hospital; <sup>12</sup>Harvard Stem Cell Institute; <sup>13</sup>Kempner Institute for the Study of Natural and Artificial Intelligence, Harvard University; <sup>14</sup>Broad Institute of MIT and Harvard; <sup>15</sup>Harvard Data Science Initiative; \*Co-first authors; ‡Corresponding author: marinka@hms.harvard.edu

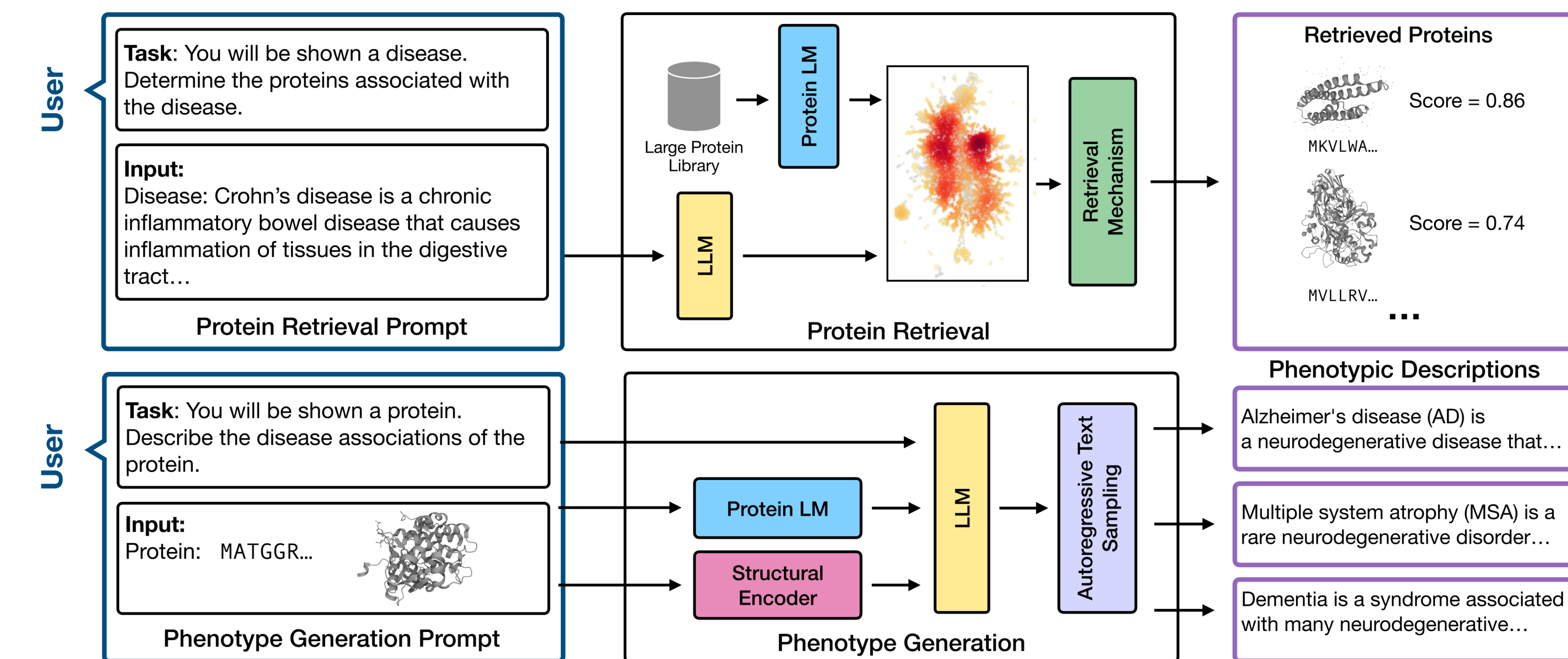
## Motivation and Model Overview

- Characterization of protein sequence and structure has rapidly outpaced functional characterization
- 20% of human protein-coding genes lack known functions
- 95% of life science publications focus on ~5000 well-characterized proteins



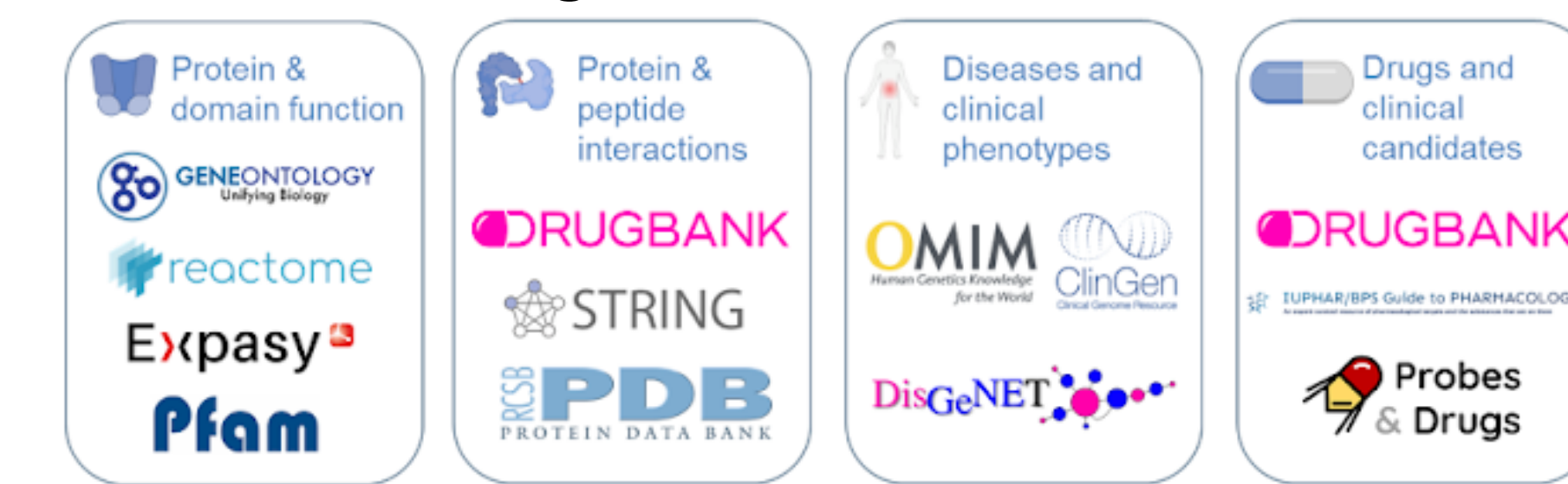
- **Phenotype:** Observable characteristic of protein function, e.g. enzymatic activity, pathway involvement, disease association
- Comprehensive phenotype annotation requires:
  - Integration of knowledge across disparate domains
  - Generalizing beyond pre-defined categories (GO)
- Natural language as a unifying modality to solve these problems

- Solution: Co-train large language model and protein representation models
- Result: Model that operates over **interleaved natural language and protein inputs**, enabling **universal annotation of proteins**



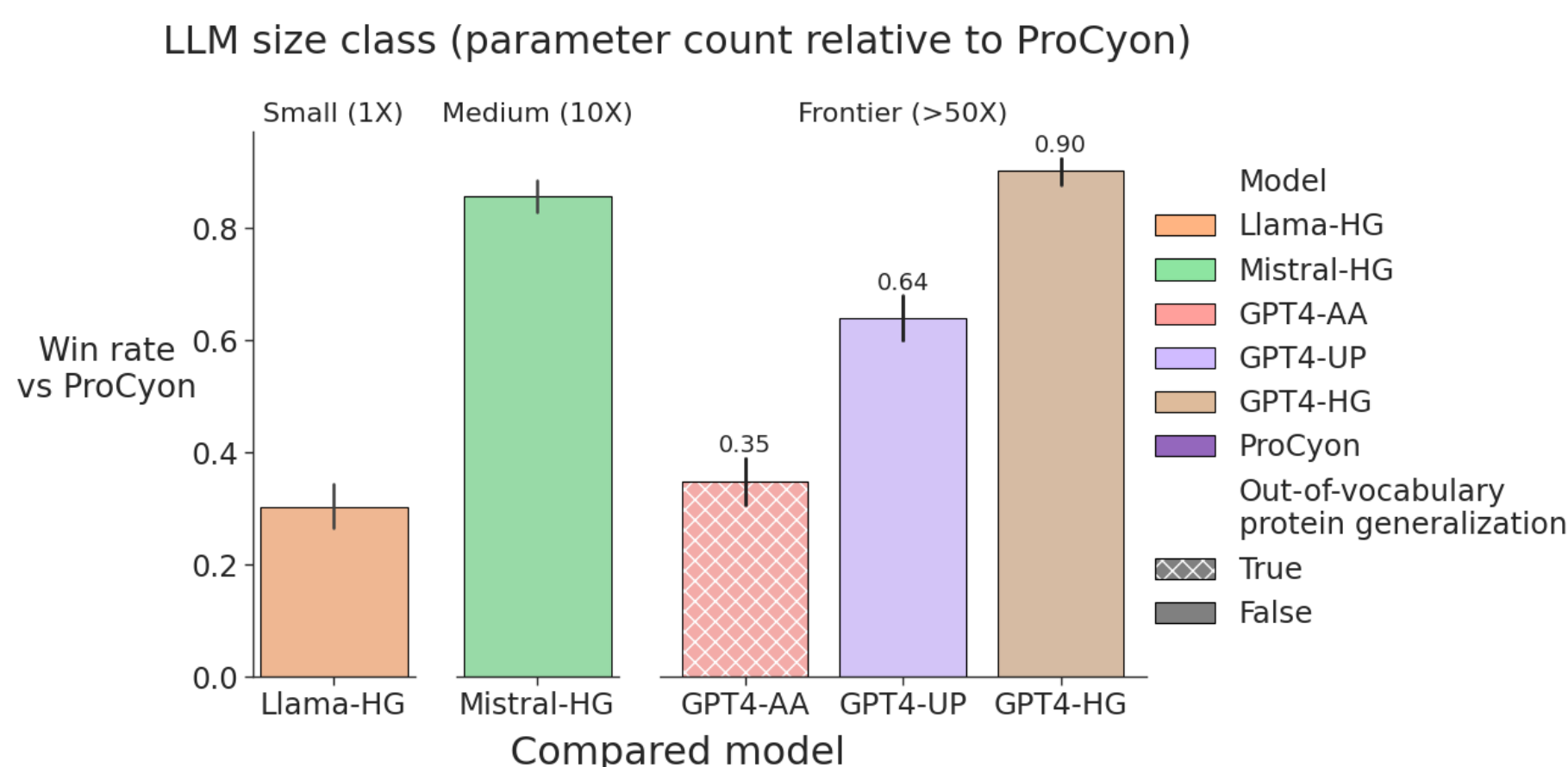
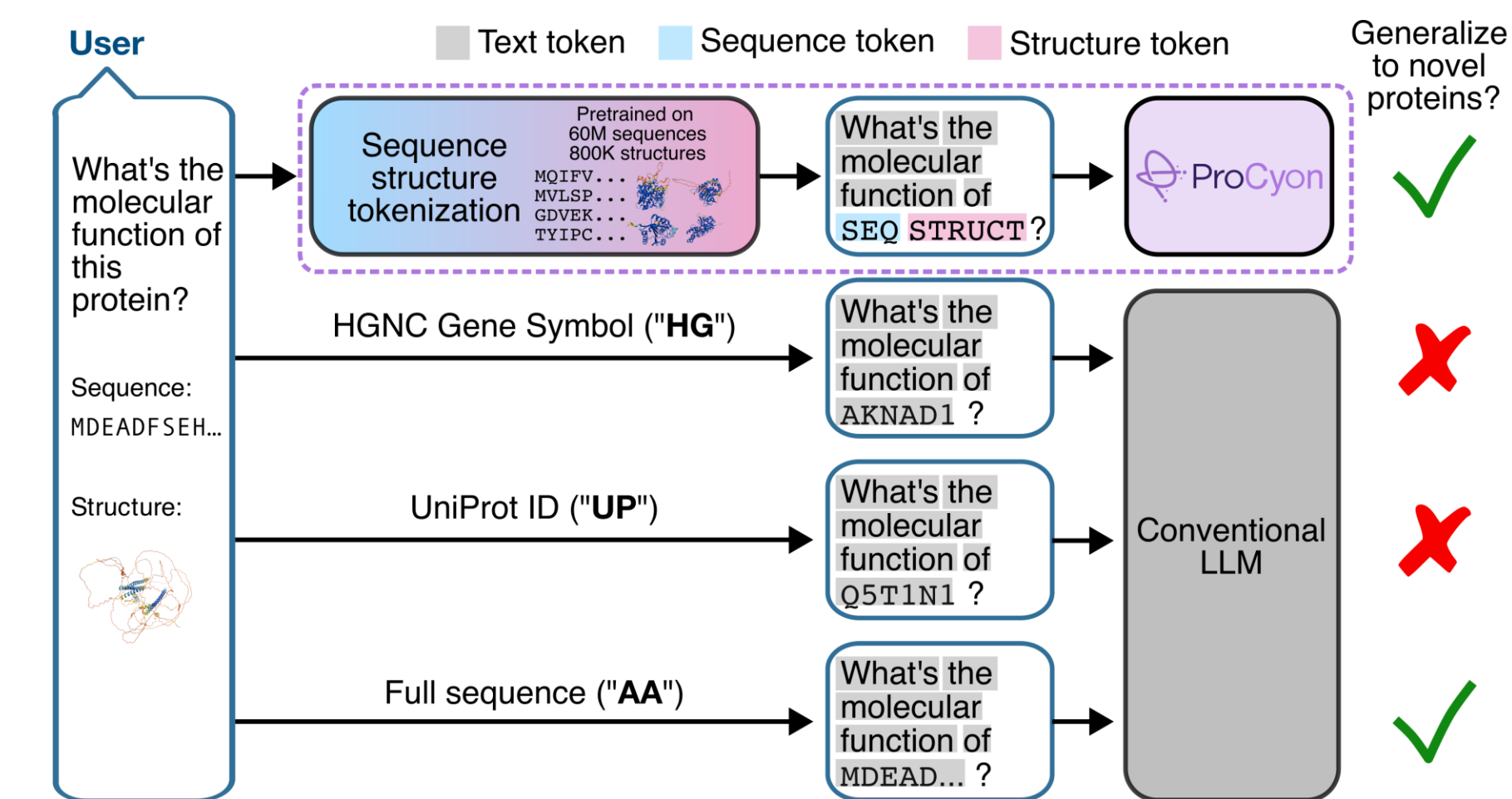
## ProCyon-Instruct Dataset

- Source protein phenotypes from broad range of sources
- 18k human proteins + 48k phenotypes → 1.85M protein-phenotype pairs → LLM rephrasing → **33M samples**
- Protein and phenotype are transformed into prompts for **instruction tuning** of the model

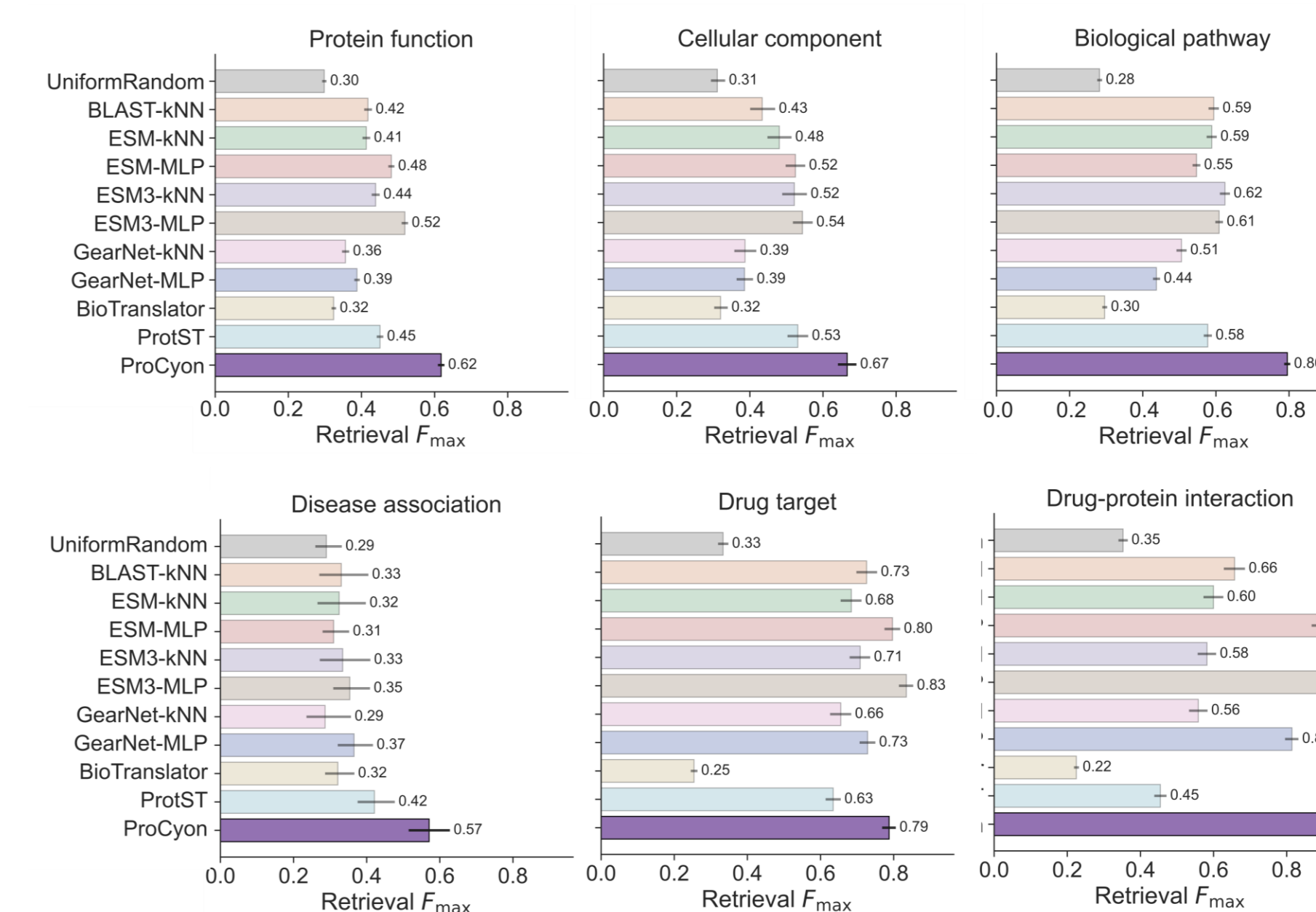


**Task definition:** ...  
**Positive example:**  
 Input: [phenotype\_description]  
 Protein: [protein\_embedding]  
 Output: yes  
**Negative example:**  
 Input: [phenotype\_description]  
 Protein: [protein\_embedding]  
 Output: no  
**Instance:**  
 Input: [phenotype\_description]  
 Protein: [protein\_embedding]  
 Output: ?

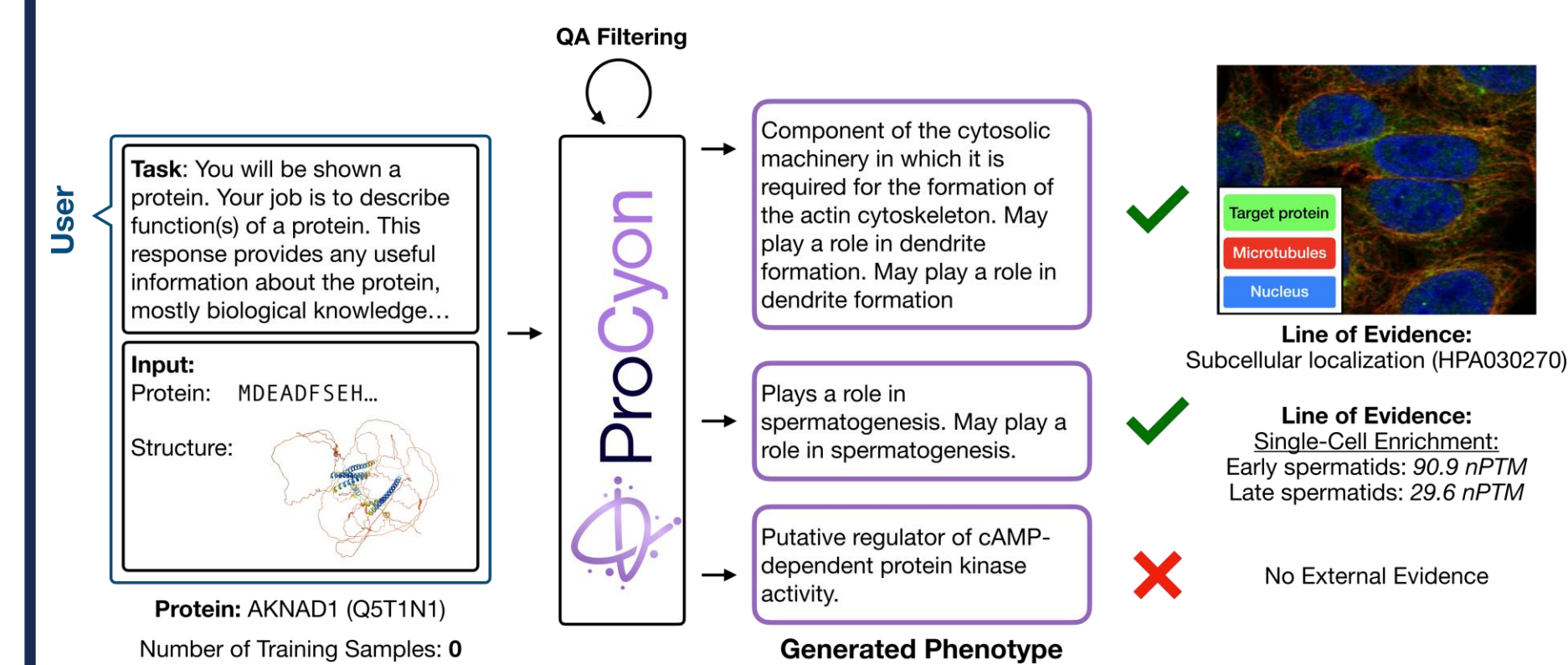
## Comparison to Conventional LLMs



## Benchmarking



## Poorly Characterized Protein Annotation



ProCyon generated phenotypes for the protein AKNAD1. Supporting external evidence found for 2 out of 3 predicted phenotypes.