

Interleaved Tool-Call Reasoning for Protein Function Understanding

Anonymous ACL submission

Abstract

Recent advances in large language models (LLMs) have highlighted the effectiveness of chain-of-thought reasoning in symbolic domains such as mathematics and programming. However, our study shows that directly transferring such text-based reasoning paradigms to protein function understanding is ineffective: reinforcement learning mainly amplifies superficial keyword patterns while failing to introduce new biological knowledge, resulting in limited generalization. We argue that protein function prediction is a knowledge-intensive scientific task that fundamentally relies on external biological priors and computational tools rather than purely internal reasoning. To address this gap, we propose PFUA, a tool-augmented protein reasoning agent that unifies problem decomposition, tool invocation, and grounded answer generation. Instead of relying on long unconstrained reasoning traces, PFUA integrates domain-specific tools to produce verifiable intermediate evidence. Experiments on four benchmarks demonstrate that PFUA consistently outperforms text-only reasoning models with an average performance improvement of 103%.

1 Introduction

Understanding protein function is a fundamental task in computational biology with broad implications in drug discovery, disease understanding, and synthetic biology. Despite the exponential growth in protein sequence databases, a significant portion of proteins lack reliable functional annotations (Zhou et al., 2019). Wet-lab experiments to determine protein function are time-consuming and resource-intensive, motivating scalable computational pipelines for automated functional annotation from sequence information. The task of protein function understanding aims to automatically predict the biological roles of proteins using computational models (Fang et al., 2024). These

include predictions of catalytic reactions, cellular functions based on Gene Ontology (GO) terms, and the identification of conserved domains or sequence motifs. Traditional supervised fine-tuning (SFT) approaches learn a direct mapping from protein sequences to functional outputs, often achieving competitive performance. These models encode protein modality features effectively, but their predictions remain largely uninterpretable. As a result, understanding the reasoning process behind functional predictions, and enabling models to generalize beyond pattern matching remains a persistent challenge.

Recently, the DeepSeek R1 model has demonstrated remarkable reasoning capabilities enabled by cold-start reasoning data construction and GRPO-based reinforcement learning. These methods significantly enhance reasoning performance in symbolic domains such as mathematics and code generation. Inspired by R1-style text-based reasoning, we construct a cold-start protein reasoning dataset using *kimi-k2-0905-preview*. We then train the Qwen2.5-3B model with SFT followed by reinforcement learning, using a mixture of format rewards and accuracy rewards computed via ROUGE_L and F1. However, our early explorations reveal a notable discrepancy between protein reasoning tasks and symbolic reasoning tasks. Without SFT on Kimi-generated reasoning data, the Qwen2.5-3B-R1-Zero model primarily receives format rewards while failing to achieve sufficient accuracy rewards. After cold-start supervised fine-tuning, the initial reward increases substantially, but subsequent improvement plateaus quickly, converging around 0.4. Inspection of the model’s generated rationales shows that the model tends to rely on repetitive high-frequency keywords to accumulate partial rewards, rather than identifying biologically meaningful functional cues. This failure mode highlights the science knowledge-intensive nature of protein function understanding: unlike

084 mathematics or programming, protein-related ques-
 085 tions cannot be solved through symbolic reasoning
 086 alone. This finding aligns with conclusions from
 087 prior work (Yue et al., 2025), which shows that
 088 reinforcement learning mainly improves the sam-
 089 pling probability of correct reasoning trajectories,
 090 but the reasoning capability itself is largely deter-
 091 mined during pre-training. RL does not grant the
 092 model new knowledge, nor can it compensate for
 093 missing domain expertise.

094 Based on these observations, we posit that long
 095 CoT symbolic reasoning training, which is highly
 096 effective for mathematical or code-generation tasks,
 097 is not directly applicable to protein function under-
 098 standing. Protein function prediction fundamen-
 099 tally depends on domain knowledge, structural
 100 priors, and evolutionary constraints, rather than
 101 purely abstract deduction. The integration of large
 102 language models with external knowledge sources
 103 and computational tools has emerged as a promis-
 104 ing approach for scientific applications requiring
 105 both reasoning and domain expertise, grounding
 106 LLM outputs in external corpora can reduce hallu-
 107 cination and enable knowledge-intensive question
 108 answering. Therefore, We argue that agent-style
 109 reasoning frameworks with domain-specific tools
 110 are better aligned with the intrinsic demands of
 111 protein function understanding task.

112 To this end, we propose PFUA, a tool-powered
 113 protein reasoning agent that couples an online
 114 LLM with computational biology tools. Instead
 115 of relying on unconstrained long-chain symbolic
 116 CoT, PFUA decomposes the query, invokes tools
 117 only when needed, and iteratively updates hy-
 118 potheses based on verifiable tool outputs, yield-
 119 ing grounded reasoning traces and more reliable
 120 answers. Extensive experiments across four bench-
 121 marks demonstrate that PFUA consistently out-
 122 performs BioMedGPT-R1, improving the average
 123 ROUGE-L recall by 98.20% on Mol-Instructions.
 124 On UniProtQA, PDB-QA, and CAFA, PFUA
 125 further surpasses BioMedGPT-R1 by 233.53%,
 126 24.97%, and 55.57%, respectively. Our main con-
 127 tributions are as follows:

- 128 • We empirically characterize the mismatch be-
 129 tween internal text-based reasoning and pro-
 130 tein function understanding.
- 131 • we introduce PFUA, a new inference
 132 paradigm for protein function understanding
 133 that explicitly incorporates biological tools
 134 into the reasoning process.

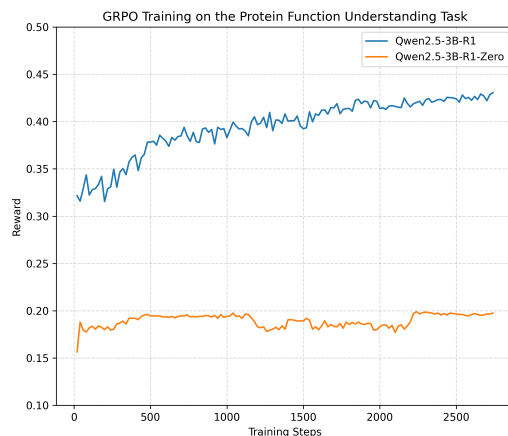


Figure 1: GRPO Training of the Protein Function Understanding Task.

- 135 • To our knowledge, we release the first multi-
 136 turn, tool-interleaved long thinking corpus for
 137 protein function understanding.

2 Related Work 138

2.1 LLMs for Protein Understanding 139

140 Recent protein large language models (LLMs) re-
 141 formulate protein understanding tasks into text gen-
 142 eration problems (Abdine et al., 2024; Fei et al.,
 143 2025). These methods typically align protein se-
 144 quence or structure representations with pretrained
 145 LLMs via query-based compression (Liu et al.,
 146 2024b), cross-attention (Qiu et al., 2024), projec-
 147 tion (Liu et al., 2024a), or discrete tokenization
 148 strategies (Ma et al., 2025). While achieve strong
 149 empirical performance, these data-driven LLMs
 150 largely operate as black-box predictors, relying on
 151 statistical correlations rather than explicit biochem-
 152 ical reasoning. The core challenge lies in inter-
 153 pretability and domain-specific reasoning.

2.2 Reasoning LLMs 154

155 Chain-of-Thought (CoT) prompting enables LLMs
 156 to perform multi-step reasoning by generating inter-
 157 mediate natural language explanations (Wei et al.,
 158 2022; Kojima et al., 2022; Wang et al., 2022). Re-
 159 cent advances further improve reasoning through
 160 test-time scaling (Snell et al., 2024) and reinforc-
 161 e-ment learning (Guo et al., 2025), leading to strong
 162 results in domains such as mathematics (Shao et al.,
 163 2024). However, in scientific applications of pro-
 164 tein understanding, such reasoning often remains
 165 purely text-based, as exemplified by BioMedGPT-
 166 R1 (Luo et al., 2024). The generated rationales

167 may reflect surface-level verbalization rather than
168 grounded mechanistic inference, as the model lacks
169 access to structured biological knowledge and com-
170 putational validation. This limits the applicabil-
171 ity of standalone reasoning LLMs to complex bio-
172 chemical problems.

173 2.3 Tool-Powered LLMs

174 The integration of large language models with ex-
175 ternal knowledge sources and computational tools
176 has emerged as a promising approach for scientific
177 applications requiring both reasoning and domain
178 expertise (Chen et al., 2023; Jin et al., 2025; Song
179 et al., 2025; Li et al., 2025). Retrieval Augmented
180 Generation (RAG) demonstrated how grounding
181 LLM outputs in external corpora can reduce hallu-
182 cination and enable knowledge-intensive question
183 answering (Lewis et al., 2020). Building on this
184 foundation, ReAct (Yao et al., 2023) interleaves
185 reasoning with action execution, allowing mod-
186 els to incorporate tool outputs and observations
187 into their reasoning process. ReTool (Feng et al.,
188 2025) proposed to leverage reinforcement learning
189 to strategically determine when and how to invoke
190 the code interpreter.

191 3 Methods

192 3.1 Tool Pool

193 In order to equip our agent with robust capabil-
194 ities for autonomous protein function investiga-
195 tion, we construct a curated pool of computational
196 tools (Cheskis et al., 2024). We prioritize tools
197 that are programmatically accessible, provide rapid
198 responses, and offer high evidential value for func-
199 tion prediction. As shown in Figure 2, these tools
200 are executed through a unified executor, allowing
201 the model to call them seamlessly within its reason-
202 ing loop.

203 **Sequence basic properties.** As a fast,
204 mechanism-agnostic sanity check, which
205 computes lightweight descriptors directly from the
206 amino acid sequence. The tool reports (i) sequence
207 length, (ii) the maximum hydrophobic run length
208 as a proxy for transmembrane propensity, and (iii)
209 a low-complexity index to flag highly repetitive or
210 compositionally biased sequences. These features
211 support early-stage triage: for example, extremely
212 long hydrophobic runs suggest membrane proteins
213 (for which soluble-enzyme assumptions may not
214 hold), while high low-complexity scores are often
215 associated with intrinsically disordered regions

216 typical of regulatory proteins. In the pipeline of
217 catalytic activity task, for example, this tool can
218 be used to establish a baseline hypothesis about
219 whether a query looks enzyme-like and to prevent
220 overconfident downstream interpretation when the
221 sequence strongly indicates a non-enzymatic class.

222 **Homology search with MMseqs2.** To ground
223 predictions in curated biological knowledge, we
224 use MMseqs2 for rapid sequence similarity search
225 against a high-quality reference database. Specif-
226 ically, we select Swiss-Prot as the target database
227 due to its strong curation standards and rich func-
228 tional annotations. The current Swiss-Prot snap-
229 shot used in our experiments contains 573,661 en-
230 tries. Given a query sequence, this tool performs
231 an MMseqs2 search and selects the best hit us-
232 ing a deterministic ranking criterion (e.g., lowest
233 E-value, highest bit score). It then extracts struc-
234 tured evidence from the corresponding Swiss-Prot
235 record, including protein name, FUNCTION text,
236 catalytic activity statements (reaction equations),
237 EC numbers, cofactors, subcellular locations, and
238 GO terms. The resulting evidence JSON provides
239 an auditable bridge from homology to functional
240 inference, enabling the agent to (i) constrain the
241 hypothesis space to a specific protein family/mech-
242 anism class and (ii) select the most appropriate
243 catalytic reaction when multiple reactions or side
244 activities are listed in the annotation.

245 **Pfam domain analysis.** Pfam is a widely used
246 protein domain database that represents conserved
247 protein families as profile hidden Markov models
248 (HMMs). In our tool pool, we employ Pfam HMM
249 scanning as a primary mechanism-level analysis
250 step. Given a query protein sequence, we scan it
251 against the Pfam-A HMM library to identify sta-
252 tistically significant domain hits, along with their
253 alignment boundaries, coverage, and confidence
254 scores. Importantly, Pfam analysis constrains the
255 functional hypothesis space at the domain and fold
256 level before any protein-level annotation is consid-
257 ered. By anchoring predictions in conserved do-
258 main families (e.g., transferase folds, oxidoreduc-
259 tase domains, or regulatory modules), the agent can
260 reason about plausible biochemical mechanisms
261 while avoiding premature commitment to overly
262 specific functions. Many Pfam families encompass
263 multiple related enzymes that share a conserved
264 fold but differ in substrate specificity or biological
265 role. By treating Pfam hits as mechanistic and ar-
266 chitectural evidence, rather than direct functional

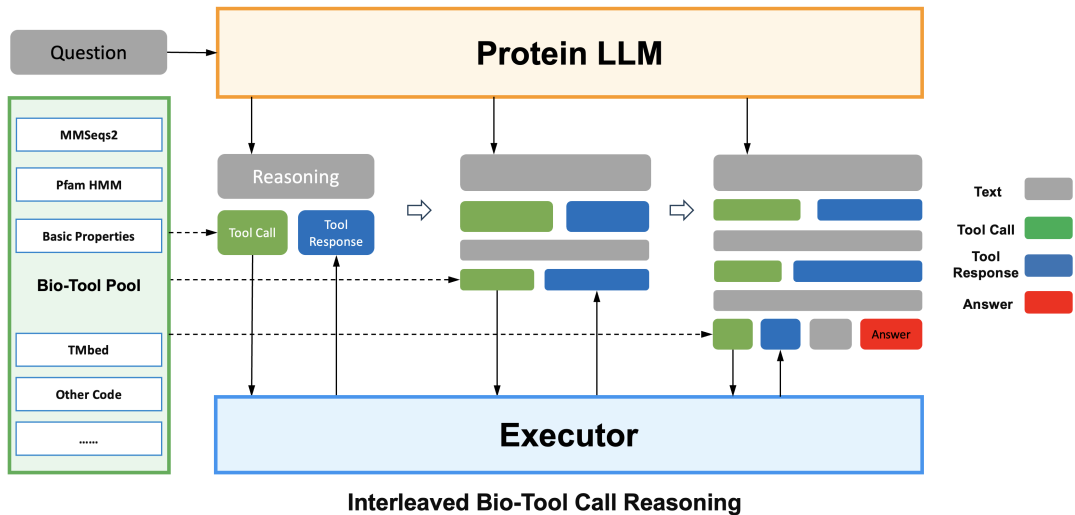


Figure 2: Overview of the interleaved tool call reasoning pipeline for protein function understanding.

267 labels, the model avoids overinterpretation and pre- 300
 268 serves flexibility for subsequent disambiguation 301
 269 steps. Subsequent tools, such as sequence homol- 302
 270 ogy search are then used to refine and disambiguate 303
 271 these domain-level hypotheses. 304
 305

272 TMbed transmembrane topology prediction. 306

273 TMbed is a protein transmembrane topology pre- 307
 274 diction tool based on large-scale protein language 308
 275 model embeddings. Instead of relying on hand- 309
 276 crafted hydrophobicity rules or shallow statisti- 310
 277 cal models, TMbed leverages contextualized se- 311
 278 quence embeddings derived from pretrained trans- 312
 279 former models to infer transmembrane helices and 313
 280 membrane-associated regions with high sensitivity 314
 281 and robustness, particularly for short or atypical 315
 282 sequences. In our tool pool, TMbed is used as a 316
 283 structure-aware localization discriminator that op- 317
 284 erates prior to homology-based annotation. Given 318
 285 a query amino acid sequence, TMbed predicts 319
 286 residue-level topology states (transmembrane helix 320
 287 versus non-membrane) and provides a global as- 321
 288 sessment of whether the protein is likely membra- 322
 289 ne-associated. This information is especially critical 323
 290 for GO annotation tasks, where cellular component 324
 291 (CC) terms such as membrane, endoplasmic retic- 325
 292 ulum membrane, or membrane-associated com- 326
 293 plexes fundamentally constrain the plausible func- 327
 294 tional hypotheses. We emphasize TMbed as a high- 328
 295 value intermediate tool rather than a standalone 329
 296 annotator. Its predictions are not interpreted in is- 330
 297 olation; instead, they are integrated with sequen- 331
 298 ce-level sanity checks and homology-based evidence
 299 to ensure consistency between predicted topology,

known protein families, and curated functional an- 300
 notations. In particular, TMbed is highly informa- 301
 tive for short proteins and small complex subunits, 302
 where traditional domain databases may provide 303
 limited coverage but membrane topology remains 304
 a decisive biological signal. 305

4 Experiments 306

4.1 Experimental Setup 307

4.1.1 Baselines 308

309 We compare with four categories of baselines. 310
 (1) **SFT Methods:** ProfT3 (Liu et al., 2024b), 311
 Prot2Text (Abdine et al., 2024), BioMedGPT (Luo 312
 et al., 2024), and Qwen2.5-3B-SFT (Hui et al., 313
 2024), which rely on parametric knowledge. (2) 314
Text-based reasoning (R1-style): BioMedGPT- 315
 R1 (Luo et al., 2024) and Qwen2.5-3B-R1, which 316
 are trained to produce intermediate long reasoning 317
 traces before the final answer. (3) **Online LLMs:** 318
 we directly prompt closed-source online LLMs to 319
 answer protein oriented queries. (4) **Multi-source 320
 RAG:** online LLMs equipped with retrieval, where 321
 tool results from multiple sources are appended to 322
 the query as additional context. (5) **Tool-powered 323
 protein agents:** online LLMs that interleave rea- 324
 soning with explicit tool calls during inference, 325
 enabling multi-step, tool-aware decision making 326
 for protein QA.

327 For all online-LLM-based baselines, we choose 328
 the *Kimi-K2-Thinking* (Team et al., 2025), *Qwen3- 329
 Max-Preview* (Yang et al., 2025), and *DeepSeek- 330
 Reasoner* (Guo et al., 2025) models as the back- 331
 bones. To ensure deterministic decoding, we set the

Methods	Func.	Cat.	Dom.	Desc.	Avg.
<i>Supervised Finetuning</i>					
BioMedGPT	5.98 / 4.28	7.97 / 6.28	1.81 / 1.81	4.84 / 4.12	5.15 / 4.12
ProtT3	15.46 / 10.81	17.36 / 12.41	16.83 / 12.22	19.49 / 15.96	17.28 / 12.85
Prot2Text	16.56 / 11.61	18.24 / 13.05	11.49 / 9.37	49.14 / 47.38	23.86 / 20.35
Qwen2.5-3B-SFT	40.74 / 30.85	41.45 / 34.52	42.60 / 32.13	33.95 / 25.68	39.69 / 30.79
<i>Text-based Reasoning</i>					
BioMedGPT-R1	35.16 / 26.80	27.64 / 22.22	30.60 / 23.33	27.78 / 20.82	30.30 / 23.29
Qwen2.5-3B-R1	49.10 / 38.04	61.33 / 46.54	51.01 / 41.10	42.56 / 32.37	51.00 / 39.51
<i>Online LLM Baseline</i>					
DeepSeek [†]	26.65 / 19/68	21.86 / 18.02	26.09 / 19.70	16.53 / 10.76	22.78 / 17.04
Kimi [‡]	25.08 / 17.32	22.99 / 17.97	27.28 / 20.07	22.30 / 14.80	24.41 / 17.54
Qwen [♠]	18.53 / 12.62	15.93 / 12.27	15.13 / 11.28	18.59 / 11.94	17.05 / 12.03
<i>Multi-Source RAG</i>					
DeepSeek [†]	38.23 / 25.68	32.36 / 23.75	30.35 / 21.46	28.54 / 16.73	32.37 / 21.91
Kimi [‡]	22.30 / 14.14	18.56 / 13.15	21.64 / 17.12	45.41 / 26.44	26.98 / 17.71
Qwen [♠]	37.42 / 24.52	50.66 / 40.67	16.44 / 12.44	42.61 / 26.34	36.78 / 25.99
<i>Tool-Powered Reasoning</i>					
DeepSeek [†]	59.71 / 38.31	47.95 / 34.15	48.36 / 34.36	58.11 / 35.70	53.53 / 35.63
Kimi [‡]	57.98 / 35.73	67.68 / 46.90	39.76 / 28.56	60.08 / 36.78	56.38 / 36.99
Qwen [♠] (PFUA)	66.43 / 44.29	72.32 / 54.18	54.98 / 44.26	63.60 / 41.90	64.33 / 46.16

Table 1: Main results on protein-oriented tasks from the Mol-Instructions dataset (Fang et al., 2024). The tasks include protein function prediction (Func.), catalytic activity prediction (Cat.), domain and motif recognition (Dom.), and general textual description generation (Desc.). For each task, performance is reported using ROUGE-1 and ROUGE-L recall (ROUGE-1 / ROUGE-L). For all online-LLM-based settings, we use [†]DeepSeek-Reasoner (Guo et al., 2025), [‡]Kimi-K2-Thinking (Team et al., 2025), and [♠]Qwen3-Max-Preview (Yang et al., 2025) as the backbones.

sampling temperature to 0.0 for all models (including both SFT-based and online-LLM-based baselines). The prompt templates for online LLM baselines, multi-source RAG, and tool-powered settings are provided in Appendix A. For the Qwen2.5-3B-R1 baseline, the cold-start SFT data are synthesized using *kimi-k2-0905-preview* following the template in Appendix C.

4.1.2 Benchmarks

We evaluate on four protein QA benchmarks covering complementary knowledge sources and reasoning demands: Mol-Instructions (instruction-following protein/molecule tasks) (Fang et al., 2024), UniProtQA (curated UniProt annotation-based questions on function, process, and localization) (Luo et al., 2024), PDB-QA (queries grounded in PDB entries, emphasizing domain/topology reasoning) (Liu et al., 2024b), and CAFA (GO-centric function inference in the CAFA setting) (Zhou et al., 2019). Together, they assess instruction understanding, factual recall, and functional reasoning.

4.1.3 Metrics

Since our models generate explicit reasoning, predictions are often longer and more descriptive than the typically concise gold references. In this case, precision-oriented metrics can be overly sensitive to additional but non-contradictory details. Therefore, we use recall-based metrics as our main results (Lin, 2004). ROUGE-1 Recall measures unigram coverage of the reference. Let R be the reference and P the prediction. ROUGE-L Recall is computed from the longest common subsequence (LCS) and normalized by the reference length:

$$\text{ROUGE-L} = \frac{\text{LCS}(R, P)}{|R|}. \quad (1)$$

These recall-based metrics emphasize whether the prediction covers the key information in the gold reference, which better matches reasoning-enhanced generation.

4.2 Main Results

Table 1 summarizes the results on Mol-Instructions. Overall, PFUA achieves the best performance across all four tasks and both ROUGE-1/ROUGE-L

[Question]

Analyze the protein with the following sequence and describe its properties:

PLAKKHNVKILPADSEHSAIFQCIQGLPEGALRRILTASGGAFRDLPVEKLKEVKVADALKHPNWNMGKKITVDSATLFN
KGLEVIEAHYLFGAHEYDDIEIVIHPPQSIHSMVETQDSSVLAQLGWPMRLPILYTLWPERVYCSEITWPRLDLCNVDLTF
KKPDHVKYPSMDLAYAAGRAGGTMTGVLSAANEKAVEMFIDEKISYLD

[Inference]

To analyze the given protein sequence

The segment 'KK...K...R...R...KK' near C-terminus (especially 'TFKK') and central 'K...R...R' suggests bipartite or monopartite NLS. Strong indicator of nuclear localization. DNA-binding motifs? Basic region + potential helix -> could resemble a basic helix-loop-helix (bHLH) or zinc finger, but only 2 Cys - not enough for classical Zn fingers. BLAST/Pfam homology (simulated): This sequence closely resembles human High Mobility Group proteins, particularly HMGB1 or HMGB2. HMGB1 signature: Two HMG-box domains (DNA-binding) - Acidic C-terminal tail - Nuclear localization - Sequence length 200-240 aa - High in K, A, P, E. Our sequence matches this profile.

Function: DNA bending/chromatin remodeling; pro-inflammatory DAMP when extracellular. Homology: Strong similarity to HMGB1 protein. Key Features: High lysine/arginine content -> DNA binding & NLS. Acidic patch (e.g., 'EYDDIE') -> characteristic of HMGB acidic tail. Redox-sensitive cysteines. No signal peptide -> non-classical secretion. HMGB1-like protein involved in DNA bending and inflammatory signaling.

[Ground Truth]

A brief overview of the protein with the provided amino acid sequence is as follows: Enzyme of the plastid non-mevalonate pathway for isoprenoid biosynthesis that catalyzes the NADPH-dependent rearrangement and reduction of 1-deoxy-D-xylulose-5-phosphate (DXP) to 2-C-methyl-D- erythritol 4-phosphate (MEP). Required for chloroplast development.

Table 2: Illustration of online LLM (Qwen3-Max-Preview) fabricating tool-call traces and purported tool outputs without actual tool execution during long chain-of-thought reasoning on Mol-Instruction (UniProtKB accession A0A1U9X9N3).

recall, substantially outperforming SFT, text-based reasoning, and multi-source RAG baselines. This indicates that explicitly tool-aware reasoning is highly effective for protein-oriented QA, particularly when the thinking process requires integrating heterogeneous biological signals.

Compared with SFT models that mainly rely on parametric knowledge, R1-style text reasoning provides notable gains (ROUGE-L +28.32% on Qwen2.5-3B backbone), suggesting that intermediate reasoning traces improve answer structuring. However, simply appending multi-source tool outputs as context (RAG) yields uneven benefits and remains limited in tasks such as domain/motif recognition. In contrast, PFUA delivers consistent improvements across Func., Cat., Dom., and Desc., with ROUGE-L +16.83% against Qwen2.5-3B-R1 on average, supporting the advantage of actively interleaving reasoning with explicit auto tool calls to query and consolidate evidence during inference. Moreover, the gains are especially pronounced on tasks that require precise mechanistic evidence (e.g., Dom. and Func.), where homology, domain boundaries, and topology signals must be jointly verified rather than heuristically inferred from text alone. These results suggest that PFUA

improves not only surface-form generation but also the reliability of evidence grounding by reducing uncertainty through targeted tool queries. Finally, the strong and stable improvements across tasks imply better generalization to diverse protein QA intents, highlighting the robustness of tool-mediated reasoning under heterogeneous biological contexts.

5 Analysis

5.1 Comparison of Three Inference Paradigms

Table 1 compares three inference paradigms under the same online-LLM backbones: (i) direct prompting (Online LLM Baseline), (ii) passive evidence injection (Multi-Source RAG, where tool outputs are appended as context), and (iii) tool-powered reasoning (where the model interleaves reasoning with explicit tool calls and evidence updates). Two consistent trends emerge. First, direct prompting yields uniformly low recall on all tasks, indicating that parametric knowledge alone is insufficient for protein-oriented questions that require precise functional, catalytic, and domain-level evidence. Moreover, as shown in Table 2, we observe that online LLMs tend to produce tool-like but unver-

375
376
377
378
379
380
381
382
383
384
385
386
387
388
389
390
391
392
393
394
395
396
397
398
399
400

401
402
403
404
405
406
407
408
409
410
411
412
413
414
415
416
417
418
419
420
421
422
423
424

ifiable statements (e.g., fabricated domain names, invented hits, or arbitrary physicochemical properties), which inflates narrative plausibility but harms evidence faithfulness and downstream answer correctness.

5.2 More Benchmarks

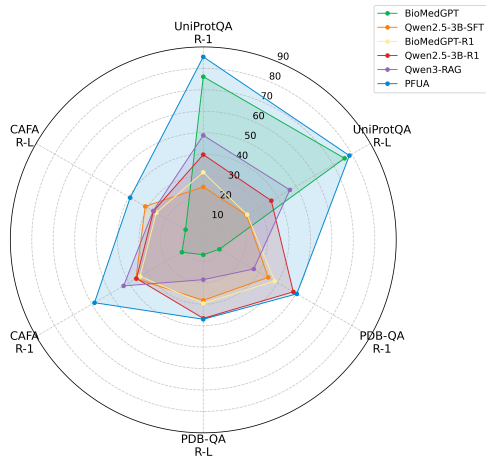


Figure 3: Results on three additional protein QA benchmarks. Performance is reported using ROUGE-1 and ROUGE-L recall (ROUGE-1 / ROUGE-L). The backbone online LLM of Qwen3-RAG and PFUA is *Qwen3-Max-Preview*.

RAG helps, but remains inconsistent and brittle. Multi-Source RAG improves over direct prompting for some backbones (e.g., Qwen: 17.05/12.03 \rightarrow 36.78/25.99 on Avg.), suggesting that providing external evidence is necessary. However, the gains are not uniform: Kimi shows only marginal improvement on Avg. (24.41/17.54 \rightarrow 26.98/17.71), and per-task performance can remain unstable, especially for domain/motif recognition. This indicates a key limitation of passive context augmentation: simply concatenating heterogeneous tool outputs does not guarantee that the model will *select, prioritize, and integrate* the right evidence, nor does it prevent partial misinterpretation of tool results.

Tool-powered reasoning yields robust, backbone-agnostic gains. In contrast, tool-powered reasoning consistently delivers large improvements across all four tasks for every backbone. Relative to direct prompting, this corresponds to substantial overall gains on Avg. (e.g., +277% in ROUGE-1 and +284% in ROUGE-L on *Qwen3-Max-Preview* backbone). The consistent uplift suggests that the key factor

is not merely *access* to external evidence, but *how* evidence is operationalized during inference: explicit tool calls enforce grounded intermediate states, enable iterative hypothesis revision, and reduce the tendency to hallucinate tool-derived facts. Overall, these results support tool-powered protein agents as a more reliable and scalable paradigm than either parametric-only prompting or passive multi-source RAG.

Figure 3 extends evaluation to UniProtQA, PDB-QA, and CAFA, which emphasize curated annotation recall and GO-centric function inference. Overall, PFUA achieves the best performance on all three benchmarks, indicating that tool-powered reasoning generalizes beyond Mol-Instructions and remains effective under diverse evidence types and question styles. A closer look shows that Multi-Source RAG brings only moderate and sometimes unstable gains, especially on PDB-QA. In contrast, PFUA yields substantial improvements over the strongest non-tool baseline BioMedGPT on UniProtQA (+12.2% / +3.4% relatively in ROUGE-1/ROUGE-L).

5.3 Case Study

Table 3 presents a representative example from Mol-Instruction (*general_function*) to compare text-only reasoning (BioMedGPT-R1) with tool-powered reasoning (PFUA). Given the sequence, BioMedGPT-R1 generates a long narrative that explicitly references “function” and “similarity” sections and then copies a porin-centric description (e.g., OmpC/OmpF, bacteriophage lambda receptor), which is not supported by any sequence-derived evidence. The reasoning also contains self-referential, low-information statements (highlighted in red) and repeated fragments, indicating that the model is hallucinating a UniProt-like annotation page and filling it with plausible-looking content rather than inferring attributes from the input sequence.

In contrast, PFUA produces an answer that matches the expected label by grounding the summary in tool evidence (full tool-call traces in Appendix B). The tool-powered summary correctly identifies the protein as a bacterial mechanosensitive channel (MscL) and states the key functional property required by the task—a stretch-activated membrane channel involved in osmotic downshock protection. Importantly, PFUA’s final answer is concise and aligned with the sequence-level signals (membrane protein with characteristic MscL fami-

Question: Examine the given protein sequence and share a brief overview of its attributes: MLKEFKEFALKGNVLDLAIIVVM-GAAFNKIVTSLVYIIMPLIGKIFGSVDFAKDWEFWGKIYGLFIQSIIDFIIVAIALFIFV KIANTLVKKEEPEEEIEENTVLL-TEIRDLLRAK

Text-based Reasoning: Okay, I'm trying to examine the protein sequence and share an attribute. **Hmm, that looks a bit complicated, but I can try to break it down.** \n\nFirst, I'm looking at the function section. **Oh, wait, that's a lot of text. I think I can understand it.** \n\nThen there's the similarity section. Oh, that's a lot to take in. I can try to look up some of the mentioned proteins. \n\nThe pore-forming subunit of a porin which assembles as an undecamer (10 subunits form the central pore, 1 subunit forms the periplasmic part) to form a channel that is 2 nm wide in the outer membrane and 1.5 nm wide in the inner membrane. The channel is non-specific and allows the passive diffusion of small molecules with a molecular weight below 600-800 Daltons through both membranes. \n\nThe homologous porins OmpC and OmpF are involved in the transport of a variety of molecules including water, ions, glucose, aminoacids, nucleoside monophosphates, lipopolysaccharide and antibiotics. \n**In addition, OmpF is involved in the recognition of the bacteriophage lambda receptor.** \n\nThe homologous porins OmpC and OmpF are involved in the transport of a variety of molecules including water, ions, glucose, aminoacids, nucleoside monophosphates, lipopolysaccharide and antibiotics. \n**In addition, OmpF is involved in the recognition of the bacteriophage lambda receptor.**

Tool-powered Reasoning: <answer> This protein is a ****large-conductance mechanosensitive channel (MscL)**** from bacteria, functioning as a critical osmotic safety valve. Key attributes:
****Structural Features:**** - 117 amino acids in length - Contains a major transmembrane domain spanning approximately residues 14-50 (predicted by TMbed) - Belongs to the MscL protein family (Pfam domain PF01741.24) with 98% coverage
****Biological Function:**** - **Acts as a stretch-activated ion channel that opens in response to membrane tension during osmotic shock** - Serves as an **emergency release valve** to prevent cell lysis by allowing rapid efflux of solutes when cells experience sudden hypoosmotic stress - Located in the cell membrane as a multi-pass transmembrane protein
****Key Evidence:**** - Extremely significant Pfam hit to the MscL domain (e-value: 8e-39) - Transmembrane prediction confirms membrane integration despite initial heuristic suggesting otherwise
This protein represents a fundamental bacterial survival mechanism, allowing cells to **rapidly respond to osmotic downshock by releasing cytoplasmic solutes** before membrane rupture occurs. </answer>

Answer: A brief overview of the protein with the provided amino acid sequence is as follows: **Channel that opens in response to stretch forces in the membrane lipid bilayer.** May participate in the **regulation of osmotic pressure changes within the cell.**

Table 3: Case study of tool-powered reasoning (PFUA) vs. text-only reasoning (BioMedGPT-R1) on Mol-Instruction (Task: *general_function*, UniProtKB accession: *A0A1U9X9N3*). The full tool-call traces and intermediate reasoning of PFUA are in Appendix B.

ly/domain support), whereas the text-only baseline drifts to an unrelated porin narrative despite its lengthy chain-of-thought.

6 Conclusion

In this work, we investigated the limitations of purely text-based reasoning paradigms for protein function understanding and showed that R1-style chain-of-thought reinforcement learning fails to deliver substantial gains in this knowledge-intensive scientific task. Our analyses demonstrate that protein function prediction relies critically on external biological priors and verifiable computational evidence, which cannot be acquired through internal reasoning alone. To bridge this gap, we proposed PFUA, a tool-augmented protein reasoning agent that integrates large language models with domain-specific biological tools, enabling grounded intermediate evidence and more reliable predictions. Extensive experiments across multiple benchmarks confirm that PFUA significantly outperforms text-only reasoning models while providing interpretable reasoning traces. This work highlights tool-integrated agents as a promising paradigm for scientific AI systems, and we antic-

ipate that future research will extend this framework to broader bioinformatics tasks and richer tool ecosystems.

Limitations

First, our study mainly demonstrates the effectiveness of tool-augmented reasoning under a fixed tool pool; the design and optimization of the tool set itself are not explored and may further affect performance. Second, although the datasets are constructed following established protocols, certain samples may still require additional manual verification to reduce potential annotation noise. Third, our experiments focus on protein QA, and the effectiveness of tool augmentation on broader protein-related tasks remains to be systematically validated, such as enzyme optimization, protein design, protein-protein interaction prediction, and molecular docking. Finally, the current evaluation relies on concise gold answers, which may not fully reflect the quality of more elaborate reasoning traces; designing metrics that better align concise references with complex reasoning outputs is left for future work.

References

- Hadi Abdine, Michail Chatzianastasis, Costas Bouyioukos, and Michalis Vazirgiannis. 2024. Prot2text: Multimodal protein’s function generation with gnns and transformers. In *Proceedings of the AAAI Conference on Artificial Intelligence*, volume 38, pages 10757–10765.
- Wenhu Chen, Xueguang Ma, Xinyi Wang, and William W. Cohen. 2023. Program of thoughts prompting: Disentangling computation from reasoning for numerical reasoning tasks. *Preprint*, arXiv:2211.12588.
- Shani Cheskis, Avital Akerman, and Asaf Levy. 2024. Deciphering bacterial protein functions with innovative computational methods. *Trends in Microbiology*.
- Yin Fang, Xiaozhuan Liang, Ningyu Zhang, Kangwei Liu, Rui Huang, Zhuo Chen, Xiaohui Fan, and Hua-jun Chen. 2024. Mol-instructions: A large-scale biomolecular instruction dataset for large language models. *Preprint*, arXiv:2306.08018.
- Xiao Fei, Michail Chatzianastasis, Sarah Almeida Carneiro, Hadi Abdine, Lawrence P Petalidis, and Michalis Vazirgiannis. 2025. Prot2text-v2: Protein function prediction with multimodal contrastive alignment. *arXiv preprint arXiv:2505.11194*.
- Jiazhan Feng, Shijue Huang, Xingwei Qu, Ge Zhang, Yujia Qin, Baoquan Zhong, Chengquan Jiang, Jinxin Chi, and Wanjun Zhong. 2025. Retool: Reinforcement learning for strategic tool use in llms. *Preprint*, arXiv:2504.11536.
- Daya Guo, Dejian Yang, Haowei Zhang, Junxiao Song, Ruoyu Zhang, Runxin Xu, Qihao Zhu, Shirong Ma, Peiyi Wang, Xiao Bi, and 1 others. 2025. Deepseek-r1: Incentivizing reasoning capability in llms via reinforcement learning. *arXiv preprint arXiv:2501.12948*.
- Binyuan Hui, Jian Yang, Zeyu Cui, Jiayi Yang, Dayiheng Liu, Lei Zhang, Tianyu Liu, Jiajun Zhang, Bowen Yu, Keming Lu, and 1 others. 2024. Qwen2. 5-coder technical report. *arXiv preprint arXiv:2409.12186*.
- Bowen Jin, Hansi Zeng, Zhenrui Yue, Jinsung Yoon, Sercan Arik, Dong Wang, Hamed Zamani, and Jiawei Han. 2025. Search-r1: Training llms to reason and leverage search engines with reinforcement learning. *Preprint*, arXiv:2503.09516.
- Takeshi Kojima, Shixiang Shane Gu, Machel Reid, Yutaka Matsuo, and Yusuke Iwasawa. 2022. Large language models are zero-shot reasoners. *Advances in neural information processing systems*, 35:22199–22213.
- Patrick Lewis, Ethan Perez, Aleksandra Piktus, Fabio Petroni, Vladimir Karpukhin, Naman Goyal, Heinrich Küttler, Mike Lewis, Wen-tau Yih, Tim Rocktäschel, and 1 others. 2020. Retrieval-augmented generation for knowledge-intensive nlp tasks. *Advances in neural information processing systems*, 33:9459–9474.
- Xuefeng Li, Haoyang Zou, and Pengfei Liu. 2025. Torl: Scaling tool-integrated rl. *Preprint*, arXiv:2503.23383.
- Chin-Yew Lin. 2004. Rouge: A package for automatic evaluation of summaries. In *Text summarization branches out*, pages 74–81.
- Nuowei Liu, Changzhi Sun, Tao Ji, Junfeng Tian, Jianxin Tang, Yuanbin Wu, and Man Lan. 2024a. Evollama: Enhancing llms’ understanding of proteins via multimodal structure and sequence representations. *Preprint*, arXiv:2412.11618.
- Zhiyuan Liu, An Zhang, Hao Fei, Enzhi Zhang, Xiang Wang, Kenji Kawaguchi, and Tat-Seng Chua. 2024b. Prott3: Protein-to-text generation for text-based protein understanding. *Preprint*, arXiv:2405.12564.
- Yizhen Luo, Jiahuan Zhang, Siqi Fan, Kai Yang, Massimo Hong, Yushuai Wu, Mu Qiao, and Zaiqing Nie. 2024. Biomedgpt: An open multimodal large language model for biomedicine. *IEEE Journal of Biomedical and Health Informatics*.
- Zicheng Ma, Chuanliu Fan, Zhicong Wang, Zhenyu Chen, Xiaohan Lin, Yanheng Li, Shihao Feng, Ziqiang Cao, Jun Zhang, and Yi Qin Gao. 2025. Prottex: Structure-in-context reasoning and editing of proteins with large language models. *Journal of Chemical Information and Modeling*.
- Jiezhong Qiu, Junde Xu, Jie Hu, Hanqun Cao, Liya Hou, Zijun Gao, Xinyi Zhou, Anni Li, Xiujuan Li, Bin Cui, Fei Yang, Shuang Peng, Ning Sun, Fangyu Wang, Aimin Pan, Jie Tang, Jieping Ye, Junyang Lin, Jin Tang, and 3 others. 2024. Instructplm: Aligning protein language models to follow protein structure instructions. *bioRxiv*.
- Zhihong Shao, Peiyi Wang, Qihao Zhu, Runxin Xu, Junxiao Song, Xiao Bi, Haowei Zhang, Mingchuan Zhang, YK Li, Yang Wu, and 1 others. 2024. Deepseekmath: Pushing the limits of mathematical reasoning in open language models. *arXiv preprint arXiv:2402.03300*.
- Charlie Snell, Jaehoon Lee, Kelvin Xu, and Aviral Kumar. 2024. Scaling llm test-time compute optimally can be more effective than scaling model parameters. *arXiv preprint arXiv:2408.03314*.
- Huatong Song, Jinhao Jiang, Yingqian Min, Jie Chen, Zhipeng Chen, Wayne Xin Zhao, Lei Fang, and Ji-Rong Wen. 2025. R1-searcher: Incentivizing the search capability in llms via reinforcement learning. *Preprint*, arXiv:2503.05592.
- Kimi Team, Yifan Bai, Yiping Bao, Guanduo Chen, Jiahao Chen, Ningxin Chen, Ruijue Chen, Yanru Chen, Yuankun Chen, Yutian Chen, and 1 others. 2025. Kimi k2: Open agentic intelligence. *arXiv preprint arXiv:2507.20534*.

663 Xuezhi Wang, Jason Wei, Dale Schuurmans, Quoc Le,
664 Ed Chi, Sharan Narang, Aakanksha Chowdhery, and
665 Denny Zhou. 2022. Self-consistency improves chain
666 of thought reasoning in language models. *arXiv*
667 *preprint arXiv:2203.11171*.

668 Jason Wei, Xuezhi Wang, Dale Schuurmans, Maarten
669 Bosma, Fei Xia, Ed Chi, Quoc V Le, Denny Zhou,
670 and 1 others. 2022. Chain-of-thought prompting elic-
671 its reasoning in large language models. *Advances*
672 *in neural information processing systems*, 35:24824–
673 24837.

674 An Yang, Anfeng Li, Baosong Yang, Beichen Zhang,
675 Binyuan Hui, Bo Zheng, Bowen Yu, Chang
676 Gao, Chengen Huang, Chenxu Lv, and 1 others.
677 2025. Qwen3 technical report. *arXiv preprint*
678 *arXiv:2505.09388*.

679 Shunyu Yao, Jeffrey Zhao, Dian Yu, Nan Du, Izhak
680 Shafraan, Karthik Narasimhan, and Yuan Cao. 2023.
681 [React: Synergizing reasoning and acting in language](#)
682 [models](#). *Preprint*, arXiv:2210.03629.

683 Yang Yue, Zhiqi Chen, Rui Lu, Andrew Zhao, Zhaokai
684 Wang, Shiji Song, and Gao Huang. 2025. Does re-
685 inforcement learning really incentivize reasoning ca-
686 pacity in llms beyond the base model? *arXiv preprint*
687 *arXiv:2504.13837*.

688 Naihui Zhou, Yuxiang Jiang, Timothy R Bergquist,
689 Alexandra J Lee, Balint Z Kacsóh, Alex W Crocker,
690 Kimberley A Lewis, George Georghiou, Huy N
691 Nguyen, Md Nafiz Hamid, Larry Davis, Tunca Do-
692 gan, Volkan Atalay, Ahmet S Rifaioğlu, Alperen
693 Dalkiran, Rengul Cetin-Atalay, Chengxin Zhang, Re-
694becca L Hurto, Peter L Freddolino, and 149 others.
695 2019. [The cafa challenge reports improved protein](#)
696 [function prediction and new functional annotations](#)
697 [for hundreds of genes through experimental screens](#).
698 *bioRxiv*.

A Templates of Three Inference Paradigms

Table 4 summarizes the prompt templates used to instantiate three representative inference paradigms for protein understanding. The *online LLM baseline* adopts a minimal instruction-only setting, where the model relies solely on internal parametric knowledge to produce a step-by-step analysis followed by a concise answer. The *multi-source RAG* template augments the model with pre-collected evidence from heterogeneous tools (e.g., computed properties, homology search, domain scanning, and topology prediction), and explicitly constrains the model to ground its reasoning in the provided outputs without requesting any additional external calls. In contrast, the *tool-powered reasoning* template frames the model as an agent that can actively decide *when* and *why* to invoke specific bioinformatics tools, requiring hypothesis-driven reasoning, uncertainty tracking, and iterative belief updates after each tool result. Collectively, these templates establish a controlled comparison from static, tool-free inference, to evidence-conditioned RAG, and finally to adaptive, decision-centric tool use for more reliable and interpretable protein analysis.

B Case Study of Full Tool-Call Reasoning

As shown in Table 5, 6, 7, to provide a concrete and fully transparent view of how our tool-powered paradigm operates in practice, we include a complete case study that records the agent’s reasoning trajectory, including all intermediate hypotheses, tool-invocation decisions, and evidence-driven updates. This example is intentionally placed in the appendix because it is substantially longer than typical main-text examples, but it serves as an important qualitative supplement to the quantitative results. Specifically, the case study demonstrates how the agent (i) starts from sequence-level cues to form initial functional hypotheses, (ii) identifies key uncertainties that cannot be resolved reliably from parametric knowledge alone, (iii) selects appropriate tools (e.g., basic physicochemical profiling, domain scanning, homology search, and topology prediction) with explicit expectations of the evidence each tool should provide, and (iv) iteratively revises its interpretation after observing tool outputs. Overall, the full trace illustrates the core advantage of tool-call reasoning: rather than producing a single-shot explanation, the agent performs hypothesis-driven evidence acquisition and

belief updating, yielding a more grounded and auditable protein interpretation.

C Synthesizing R1-Style Reasoning Traces

For the baseline Qwen2.5-3B-R1 model, which adopts a standard DeepSeek-R1-style training pipeline, we construct synthetic reasoning traces as cold-start supervision. Specifically, we prompt an online LLM, Kimi (*kimi-k2-0905-preview*), with the original Mol-Instruction inputs (i.e., the question and the protein sequence) using the template in Table 8. This procedure yields 2,000 synthetic cold-start SFT examples for training Qwen2.5-3B-R1.

Template for Online LLM Baseline

[ROLE]

You are a professional bioinformatics assistant.

[TASK]

Please first provide detailed reasoning and analysis.

[CONSTRAINTS]

Then give a concise final answer wrapped strictly inside `<answer></answer>` tags.

Template for Multi-Source RAG

[ROLE]

You are an expert protein analysis assistant.

[TASK]

Analyze the given protein sequence. You are provided with external tool outputs (computed properties, homology search, domain scan, and topology prediction).

Use these tool results as evidence to reasoning and answer the question.

[CONSTRAINTS]

- Do NOT request additional tools or external calls. Everything you need is already included below.
- The final answer MUST be wrapped in `<answer>...</answer>`.

Template for Tool-Powered Reasoning

[ROLE]

You are an expert protein analysis agent.

[TASK]

Your goal is to analyze the protein sequence and produce a biologically meaningful interpretation.

You should reason step-by-step, form hypotheses, and use tools only when they help reduce uncertainty.

[REASONING REQUIREMENTS]

Before calling tools, you MUST:

- propose hypotheses about the protein
- explain which uncertainties still remain

For EVERY tool call, you MUST:

- explicitly explain WHY this tool is needed
- describe WHAT evidence you expect it to provide

After each tool result, you MUST:

- summarize what new evidence was obtained
- update or revise your hypothesis
- decide whether additional tools are needed

[TOOLS]

You may call the following tools through function calling:

- *seq_basic_props*: basic physicochemical properties
- *pfam_hmm_scan*: domain and family inference
- *mmseqs2_besthit_uniprot*: homolog search and functional annotation
- *tmbcd_predict*: transmembrane and topology prediction
- In addition to the tools listed, the assistant is allowed to compute any other relevant properties that can be derived using Python code.
- Prefer "sequence_ref": "query" instead of pasting long sequences.

[OUTPUT]

When finished, wrap the final answer in `<answer>...</answer>`.

Table 4: Inference templates for online LLM, multi-source RAG and tool-powered protein reasoning.

```

<lim_start>user
[ROLE]
You are an expert protein analysis agent.
[TASK]
Your goal is to analyze the protein sequence and produce a biologically meaningful interpretation.
You should reason step-by-step, form hypotheses, and use tools only when they help reduce uncertainty.
[REASONING REQUIREMENTS]
Before calling tools, you MUST:
- propose hypotheses about the protein
- explain which uncertainties still remain
For EVERY tool call, you MUST:
- explicitly explain WHY this tool is needed
- describe WHAT evidence you expect it to provide
After each tool result, you MUST:
- summarize what new evidence was obtained
- update or revise your hypothesis
- decide whether additional tools are needed
[TOOLS]
You may call the following tools through function calling:
- seq_basic_props: basic physicochemical properties
- pfam_hmmscan: domain and family inference
- mmseqs2_besthit_uniprot: homolog search and functional annotation
- tmbed_predict: transmembrane and topology prediction
- In addition to the tools listed, the assistant is allowed to compute any other relevant properties that can be derived using Python code.
- Prefer "sequence_ref": "query" instead of pasting long sequences.
[OUTPUT]
When finished, wrap the final answer in <answer>...</answer>.
Examine the given protein sequence and share a brief overview of its attributes:
Protein sequence (id=query):
""MLKEFKEFALKGNVLDLAIAVVMGAAFNKIVTSLVTYIIMPLIGKIFGSVDFAKDWEFWGIKYGLFIQSIIDFIIVAIAL
FIFVKIANTLVKKEPEEEIEENTVLLTEIRDLLRAK""
<lim_end>

<lim_start>assistant
I need to analyze this protein sequence step by step. Let me start by examining the sequence and forming initial hypotheses.
The sequence is:
MLKEFKEFALKGNVLDLAIAVVMGAAFNKIVTSLVTYIIMPLIGKIFGSVDFAKDWEFWGIKYGLFIQSIIDFIIVAIAL
FIFVKIANTLVKKEPEEEIEENTVLLTEIRDLLRAK
Initial observations: 1. The sequence starts with "MLKEFKEFALK" - this looks like it could be a signal peptide or N-terminal region 2. There are several hydrophobic regions (like "LAIAVVMGAAFNKIVTSLVTYIIMPLIGKIFG") 3. There are charged regions with many glutamates (E) and lysines (K) at the end 4. The length appears to be around 120-130 amino acids
Hypotheses: 1. This could be a membrane protein due to the hydrophobic stretches 2. It might have transmembrane domains 3. The C-terminal region with many charged residues might be cytoplasmic 4. It could be involved in transport or signaling based on the domain structure
Uncertainties that remain: 1. Exact length and composition 2. Whether it has transmembrane domains and how many 3. What protein family it belongs to 4. Its specific function 5. Its subcellular localization
I should start with basic properties to get fundamental information about the sequence. This will help me understand the basic characteristics before diving into more specific analyses.
Let me call seq_basic_props first to get the basic properties.
<tool_call>
{"name": "seq_basic_props", "arguments": {"sequence_ref": "query"}}
</tool_call>
<lim_end>

<lim_start>user
<tool_response>
{"length": 117, "hydrophobic_run_max": 12, "low_complexity_index_0to1": 0.1171, "heuristics": {"looks_membrane_like": false, "looks_low_complexity_like": false}}
</tool_response>
<lim_end>

```

Table 5: Full reasoning traces of PFUA on Mol-Instruction (Task: *general_function*, UniProtKB accession: *A0A1U9X9N3*) Part 1.

<lim_start>assistant

Now I have the transmembrane prediction results. Let me analyze what this shows:

The raw prediction shows: - A long stretch of "H" characters (representing helical transmembrane regions) starting around position 14 and continuing for about 36 residues - Then a stretch of "h" characters (likely representing another type of transmembrane or membrane-associated region) later in the sequence - The prediction indicates 36 tm_signal_letter_hits and has_tm_signal_heuristic: true

This confirms that this is indeed a transmembrane protein, which aligns with the UniProt annotation that it's a "Multi-pass membrane protein". The initial heuristic from seq_basic_props that said "looks_membrane_like: false" was incorrect - this is definitely a membrane protein.

Let me also check for Pfam domains to get more detailed information about the domain architecture and confirm this is indeed an MscL channel.

<tool_call>

```
{"name": "pfam_hmmscan", "arguments": {"sequence_ref": "query"}}
```

</tool_call>

<lim_end>

<lim_start>user

<tool_response>

```
{"hits": [{"pfam_id": "MscL", "pfam_acc": "PF01741.24", "query": "query", "evalue": 8e-39, "score": 133.3, "hmm_from": 1, "hmm_to": 124, "ali_from": 1, "ali_to": 115, "coverage_query": 0.9829, "desc": "Large-conductance mechanosensitive channel, MscL"}, {"pfam_id": "Anoctamin", "pfam_acc": "PF04547.18", "query": "query", "evalue": 0.045, "score": 13.1, "hmm_from": 129, "hmm_to": 199, "ali_from": 23, "ali_to": 102, "coverage_query": 0.6838, "desc": "Calcium-activated chloride channel"}, {"pfam_id": "Mim2", "pfam_acc": "PF19117.6", "query": "query", "evalue": 0.047, "score": 13.8, "hmm_from": 15, "hmm_to": 31, "ali_from": 32, "ali_to": 48, "coverage_query": 0.1453, "desc": "Mitochondrial import 2"}, {"pfam_id": "UPF0239", "pfam_acc": "PF06783.17", "query": "query", "evalue": 0.16, "score": 12.8, "hmm_from": 21, "hmm_to": 60, "ali_from": 61, "ali_to": 99, "coverage_query": 0.3333, "desc": "Uncharacterised protein family (UPF0239)"}, {"pfam_id": "DUF3917", "pfam_acc": "PF13055.11", "query": "query", "evalue": 0.32, "score": 12.0, "hmm_from": 21, "hmm_to": 68, "ali_from": 33, "ali_to": 84, "coverage_query": 0.4444, "desc": "Protein of unknown function (DUF3917)"}], "selected_domains": [{"pfam_id": "MscL", "pfam_acc": "PF01741.24", "query": "query", "evalue": 8e-39, "score": 133.3, "hmm_from": 1, "hmm_to": 124, "ali_from": 1, "ali_to": 115, "coverage_query": 0.9829, "desc": "Large-conductance mechanosensitive channel, MscL"}]}
```

</tool_response>

<lim_end>

<lim_start>assistant

Now I have all the key information to provide a comprehensive analysis. Let me synthesize the findings:

1. Basic properties: - Length: 117 amino acids - Maximum hydrophobic run: 12 residues - Low complexity index: 0.1171 (relatively low) - Initial heuristic suggested it didn't look membrane-like, but this was incorrect
2. UniProt homology search: - 100% identity to Q4L656 (Large-conductance mechanosensitive channel) - Function: Channel that opens in response to stretch forces in membrane lipid bilayer - Role: Regulation of osmotic pressure changes within the cell - Subcellular location: Cell membrane, multi-pass membrane protein - GO terms: plasma membrane, mechanosensitive monoatomic ion channel activity
3. Transmembrane prediction: - Confirmed transmembrane signal with 36 hits - Shows a long hydrophobic transmembrane region (the HHHHH... stretch) - Confirms it's a membrane protein despite initial heuristic
4. Pfam domain analysis: - Strong hit to MscL domain (PF01741.24) with extremely low e-value (8e-39) - Coverage of 98.29% of the query sequence - Other weaker hits are likely false positives or overlapping domains - The MscL domain is definitive for large-conductance mechanosensitive channels

This is clearly a bacterial mechanosensitive channel protein that acts as an emergency release valve during osmotic shock. The transmembrane prediction shows the characteristic hydrophobic region that spans the membrane, and the Pfam analysis confirms it belongs to the MscL family.

The initial heuristic that said it didn't look membrane-like was incorrect - this is a classic membrane protein. The hydrophobic run of 12 residues was actually part of a larger transmembrane domain.

I can now provide a comprehensive final answer.

<answer>

This protein is a **large-conductance mechanosensitive channel (MscL)** from bacteria, functioning as a critical osmotic safety valve. Key attributes:

Structural Features: - 117 amino acids in length - Contains a major transmembrane domain spanning approximately residues 14-50 (predicted by TMbed) - Belongs to the MscL protein family (Pfam domain PF01741.24) with 98

Biological Function: - Acts as a stretch-activated ion channel that opens in response to membrane tension during osmotic shock - Serves as an emergency release valve to prevent cell lysis by allowing rapid efflux of solutes when cells experience sudden hypoosmotic stress - Located in the cell membrane as a multi-pass transmembrane protein

Key Evidence: - Extremely significant Pfam hit to the MscL domain (e-value: 8e-39) - Transmembrane prediction confirms membrane integration despite initial heuristic suggesting otherwise

This protein represents a fundamental bacterial survival mechanism, allowing cells to rapidly respond to osmotic downshock by releasing cytoplasmic solutes before membrane rupture occurs.

</answer>

<lim_end>

Table 7: Full reasoning traces of PFUA on Mol-Instruction (Task: *general_function*, UniProtKB accession: A0A1U9X9N3) Part 3.

Template for Synthesizing Reasoning Trace

You are an expert AI in bioinformatics and computational biology. Your task is to generate a detailed, step-by-step intermediate reasoning process that connects the given input to the given output. The reasoning process should explain how one might logically derive the output from the input.

Input:

question: "question"

protein sequence: "sequence"

Output:

answer: "answer"

Your response should consist solely of the reasoning process enclosed in <think> ... </think>.

Table 8: Cold-start template for synthesizing reasoning traces used to train Qwen2.5-3B-R1 (generated via *kimi-k2-0905-preview*).