

# ProtoLLM: Automatic Evaluation Framework of LLMs on Domain-Specific Scientific Protocol Formulation Tasks

Anonymous EMNLP submission

## Abstract

Automated generation of scientific protocols executable by robots can significantly accelerate scientific research processes. Large Language Models (LLMs) excel at Scientific Protocol Formulation Tasks (SPFT), but the evaluation of their capabilities rely on human evaluation. Here, we propose a flexible, automatic framework to evaluate LLMs' capability on SPFT: *ProtoLLM*<sup>1</sup>. This framework prompts the target model and GPT-4 to extract pseudocode from biology protocols using only predefined lab actions and evaluates the output of target model using LLAM-EVAL, the pseudocode generated by GPT-4 serving as a baseline and Llama-3 acting as the evaluator. Our adaptable prompt-based evaluation method, LLAM-EVAL, offers significant flexibility in terms of evaluation model, material, criteria, and is free of cost. We evaluate GPT variations, Llama, Mixtral, Gemma, Cohere, and Gemini. Overall, we find that GPT and Cohere is a powerful scientific protocol formulators. We also introduce BIOPROT 2.0, a dataset with biology protocols and corresponding pseudocodes, which can aid LLMs in formulation and evaluation of SPFT. Our work is extensible to assess LLMs on SPFT across various domains and other fields that require protocol generation for specific goals.

## 1 Introduction

Laboratory automation is essential for accelerating scientific research processes. However, most contemporary laboratories use manual labor, especially in the field of biology. This not only constrains the scope for scalability, but also introduces potential vulnerabilities in reproducibility (Kwok, 2010).

One of the barriers for automation in biology is the reliance on manual experiments when validating scientific protocols. Traditionally, trial-and-error approach has been employed to formulate

<sup>1</sup>The dataset and code are available [here](#).

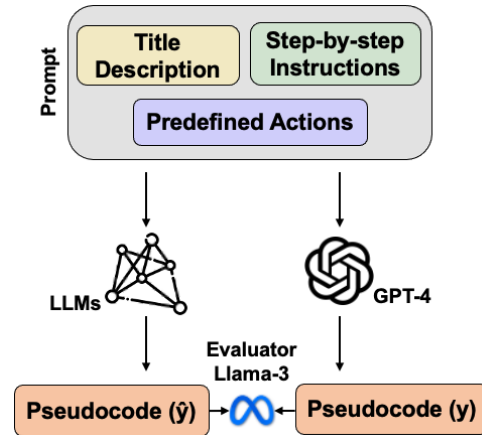


Figure 1: **Overview of the *ProtoLLM* Framework.** A protocol containing a title, descriptions, step-by-step instructions, and predefined biology lab actions is given to both a target model and GPT-4 for pseudocode generation. Then, Llama-3 evaluates these outputs considering the target model's pseudocode as the prediction ( $\hat{y}$ ) and GPT-4's as a baseline ( $y$ ).

a protocol to achieve a certain goal. As a breakthrough, LLMs have demonstrated remarkable capabilities in formulating precise experimental protocols across diverse fields (White et al., 2023; Jablonka et al., 2023). These protocols comprise pseudocodes with actionable sequences that can be executed by machines which can be automated. Yet, efforts in biology to utilize LLMs for pseudocode formulation are yet to achieve desired outcomes (Inagaki et al., 2023). These works rely on human evaluations, and objective evaluation methods for protocol formulation are nonexistent. Therefore, it is necessary to establish an automated evaluation framework on formulating protocols to move beyond manual labor.

Previous work suggests a framework to assess the capabilities of LLMs on SPFT: *BioPlanner* (O'Donoghue et al., 2023). This method outlines three primary steps: (i) extracting pseudofunctions

060 and pseudocode<sup>2</sup> from a protocol using an evalu- 107  
061 ator, (ii) using the target model to produce pseu- 108  
062 docode given the pseudofunctions, and (iii) evalu- 109  
063 ating the pseudocode generated in step (ii) against 110  
064 the original pseudocode in (i). Using this frame- 111  
065 work, they performed evaluation exclusively on 112  
066 GPTs (Brown et al., 2020; OpenAI, 2023). 113

067 We highlight the following key observations: (1) 114  
068 Various representations of pseudofunctions corre- 115  
069 sponding to identical experimental actions, causes 116  
070 performance degradation and inconsistency of the 117  
071 evaluation framework. (2) The repertoire of ac-  
072 tions executed in biology labs is confined to a fi-  
073 nite set of actions. (3) High values in traditional  
074 automatic metrics (i) does not necessarily imply  
075 human-perceived good quality in scientific proto-  
076 cols. (4) The use of automatic metrics (i) requires  
077 manual labor, which limits the transition to fully  
078 automatic evaluation.

079 Here, we propose an evaluation framework that  
080 evaluates the capabilities of LLMs in SPFT: *Proto-*  
081 *coLLM* (Figure 1). First, we define a set of actions  
082 in advance (Table 1), which eliminates individual  
083 action (pseudofunction) extraction step and vari-  
084 ations of actions on each occasion. Second, we  
085 independently zero-shot prompted the target model  
086 and GPT-4 (OpenAI, 2023) to extract pseudocode  
087 from biology protocols, only using predefined ac-  
088 tions as pseudofunctions. Lastly, we use LLAM-  
089 EVAL to evaluate the response, treating the target  
090 model’s pseudocode as a prediction ( $\hat{y}$ ) and that  
091 of GPT-4’s as a baseline ( $y$ ). LLAM-EVAL offers  
092 significant flexibility in terms of evaluation model,  
093 material, and criteria. This approach is inspired  
094 by the automated extraction of chemical synthesis  
095 actions from experimental procedures<sup>3</sup> (Vaucher  
096 et al., 2020). We compared multiple LLMs to our  
097 framework, including GPT variations (Brown et al.,  
098 2020; OpenAI, 2023), Llama, Mixtral, Gemma,  
099 Cohere, and Gemini (Google, 2024). We find that  
100 GPT-4o and Cohere+ is a powerful scientific proto-  
101 col formulator.

102 We also introduce BIOPROT 2.0, a larger dataset  
103 with scientific protocols and the corresponding  
104 pseudocodes that can aid LLMs in formulation and  
105 evaluation of SPFT.

106 Overall, we make the following contributions:

<sup>2</sup>Pseudofunctions represent laboratory actions, while pseu-  
docode embodies protocols composed of these pseudofunc-  
tions.

<sup>3</sup>A set of actions in chemistry labs were defined prior to  
the pseudocode extraction process.

1. We propose *ProtoLLM*: a flexible, auto-  
matic framework for evaluating LLMs on  
SPFT using domain knowledge and LLAM-  
EVAL.
2. We propose LLAM-EVAL, an evaluation  
method that uses a form-filling paradigm of-  
fering significant flexibility in terms of evalu-  
ation model, material, and criteria.
3. We introduce the BIOPROT 2.0 dataset, featur-  
ing protocols and corresponding pseudocode  
for evaluating and aiding LLMs on SPFT.

## 2 Related Works

**Task-specific Evaluation** LLMs have been evalu-  
ated based on their performance in specific tasks.  
Information extraction abilities were measured  
by the generated quality of summaries (Durmus  
et al., 2020; Wang et al., 2020), paper reviews  
(Zhou et al., 2024), question correction (Fan et al.,  
2024), or combination of a few tasks (Labrak et al.,  
2024). However, these studies do not provide  
comprehensive evaluations and only assess very  
limited aspects, thus limiting their generalizability  
to other abilities or tasks.

**LLM Evaluation on SPFT** Recent work pro-  
poses a three-step framework (Section 1) for  
the evaluation of scientific protocols in biology:  
*BioPlanner* (O’Donoghue et al., 2023). This  
work evaluates GPT’s performance in three tasks:  
next-step prediction, pseudocode generation, and  
pseudofunction retrieval. It employs statistical  
scoring methods including Levenshtein distance  
( $\mathcal{L}_d$ ) and BLEU (Papineni et al., 2002) to measure  
the relevance between a baseline and generated  
protocols, despite their modest correlation with  
human judgments.

**Domain-specific LLMs in Science** A Large  
number of LLMs have been trained, fine-  
tuned, or augmented for domain-specific uses.  
ChemBERTa/-2 (Chithrananda et al., 2020; Ahmad  
et al., 2022), MatSciBERT (Gupta et al., 2021),  
MaterialsBERT (Shetty et al., 2023), Chem-  
crow (Bran et al., 2023), and LLM augmentation  
methods for various experiment-related tasks (Guo  
et al., 2023) has been introduced in chemistry.  
BioGPT (Luo et al., 2022), BioBERT (Lee et al.,  
2019), CamemBERT-bio (Touchent et al., 2024),  
BlueBERT (Peng et al., 2019), PubmedBERT (Gu  
et al., 2020), BioMegatron (Shin et al., 2020), and

Action Name	Description
Transfer	Move substances between containers using lab equipment, such as pipettes.
Centrifuge	Spin at high speed to separate mixture components by density.
Vortex	Mix solutions by creating a vortex for even distribution.
SetTemp	Set specific temperatures for reactions or processes.
Wait	Period of inactivity to allow reactions or condition stabilization.
Wash	Rinse materials, often with solvents to remove contaminants.
Measure	Quantify substances or properties using instruments.
Microscopy	Use a microscope to observe and analyze cell morphology and structures.
CellDetachment	Release adherent cells from a culture surface using enzymatic or mechanical methods.
CellCount	Determine the number of cells in a sample using a hemocytometer or automated counter.
InvalidAction	Undefined action due to documentation error or ambiguity.
OtherLanguage	Text in non-English, indicating translation need.
NoAction	Text not corresponding to any defined action.
PCR	Amplify DNA segments through Polymerase Chain Reaction.
Gel	Separate molecules by size in a gel with electric field.
Culture	Grow cells in lab to study behavior or for experimentation.
Dilute	Reducing the concentration of a solution by adding solvent.

Table 1: **Predefined Set of Actions.** List of actions performed in biological experiments and the corresponding descriptions. Actions above the line represent the basic actions, with the last three specifically designated for instances where a new protocol introduces an undefined action. Actions below represent the coarse-grained actions.

157 ProtoCode (Jiang et al., 2024) has been introduced  
158 in biology.

159 **Evaluating LLMs with LLMs** Evaluation  
160 of LLMs encompasses a dual-method approach:  
161

- 162 (i) Statistical scoring: BLEU (Papineni et al.,  
163 2002), ROUGE (Lin, 2004), METEOR  
164 (Banerjee and Lavie, 2005), Levenshtein Dis-  
165 tance
- 166 (ii) Model-based scoring: G-Eval (Liu et al.,  
167 2023), Prometheus (Kim et al., 2023),  
168 BLEURT (Sellam et al., 2020), Natural Lan-  
169 guage Inference (NLI)
- 170 (iii) Combination of (i) and (ii): GPTScore (Fu  
171 et al., 2023), SelfCheckGPT (Manakul et al.,  
172 2023), BERTScore (Zhang et al., 2020),  
173 SciBERTScore (O’Donoghue et al., 2023),  
174 WMD (Kusner et al., 2015), MoverScore  
175 (Zhao et al., 2019), Question Answer Gen-  
176 eration (QAG) Score

177 In tasks where reasoning is involved, (ii)(iii) out-  
178 performs (i). Previous work adopted (i) with (iii)  
179 being minimal (O’Donoghue et al., 2023). In this  
180 work, we adopt the notion of G-Eval (Liu et al.,  
181 2023), a framework for evaluating LLM-generated  
182 text, which prompts GPT with text and criteria,  
183 then scores based on its output.

### 184 3 Methods

185 The *ProtoCoLLM* framework can evaluate the capa-  
186 bility of LLMs on SPFT in three steps (Figure 2):

(1) prompt the target LLM to generate pseudocode  
187 based on the given protocol, (2) repeat previous  
188 step for GPT-4, and (3) LLAM-EVAL for evalua-  
189 tion. To utilize this framework, we curated proto-  
190 cols in biology (Section 3.1), predefined actions  
191 performed in biology labs (Section 3.2), prompted  
192 LLMs for pseudocode generation (Section 3.3), and  
193 prompted Llama-3 for evaluation (LLAM-EVAL)  
194 (Section 3.6).  
195

#### 196 3.1 Data Curation of Protocols in Biology

197 Each protocol is composed of three core elements:  
198 a title, description, and experimental steps. We cu-  
199 rated the dataset through a process of collection and  
200 refinement. We collected a set of keywords relevant  
201 to biology. Then, we used a scoring system based  
202 on the number of keywords included in the descrip-  
203 tion of each protocol from *protocols.io*<sup>4</sup> (Teytelman  
204 et al., 2016). We refined the dataset collected in the  
205 previous step using automated and manual methods.  
206 (Appendix A.1.)

#### 207 3.2 Defining Actions

208 The defined actions are composed of two parts:  
209 (i) **basic actions** corresponding to a single action  
210 which can be performed directly in biology  
211 labs, and (ii) **coarse-grained actions** which  
212 corresponds to a large set of basic actions repeated  
213 throughout various protocols. Defined actions were  
214 reviewed by experts with intensive experiences

<sup>4</sup>A platform for reproducible protocol sharing provides access to more than 15k publicly available protocols, and has no limitations regarding the use of LLMs.

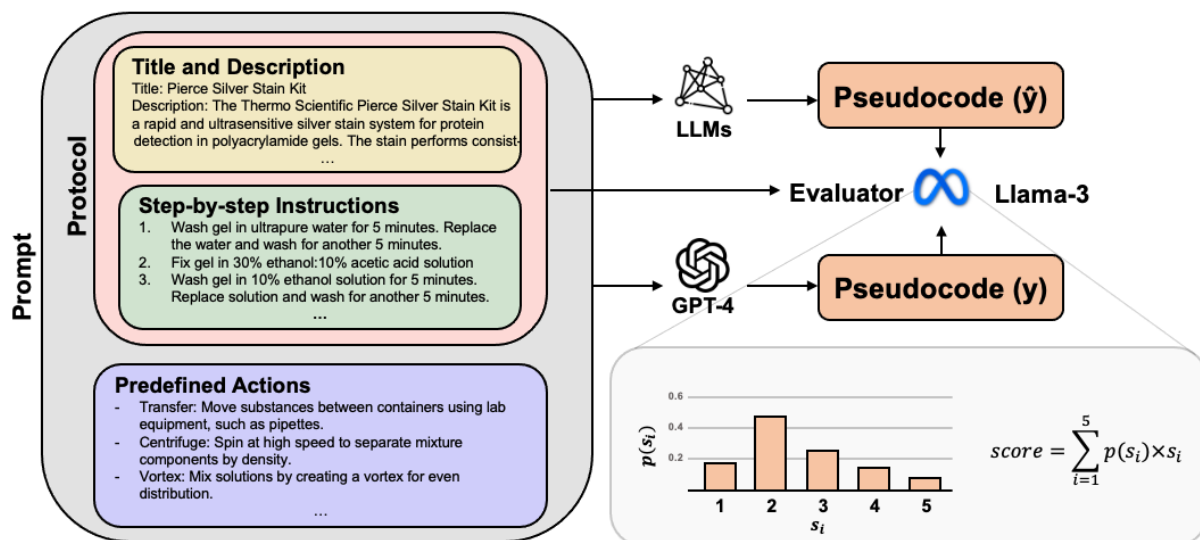


Figure 2: The *ProtoLLM* Framework.

in biology experiments. The target model specifies the arguments for each action on each occasion.

**Basic Actions** Since the repertoire of actions executed in biology labs is confined to a finite set of actions, we defined a set of actions performed in biology labs prior to the extraction of pseudocode from protocols (Table 1). We performed a comprehensive literature review to define the set of basic actions performed in biology labs.

**Coarse-grained Actions** We observed that a series of complex, repetitive actions can be effectively encapsulated and described by a single, comprehensive action. For instance, the process of diluting a solution is conceptually straightforward and can possibly be defined by basic actions. However, this involves intricate calculations and logical reasoning, which can result in performance degradation by calculation mistakes and posing variations in representations of an identical process. To this end, we coarse-grained these complex set of actions into a singular action.

### 3.3 Prompting Pseudocode Generation

To evaluate the target LLMs on SPFT, we prompted the models to generate pseudocode based on a protocol collected at Section 3.1. Models are instructed to use only the actions defined in Section 3.2 as the function name. However, they were allowed to define the arguments for each pseudofunction as needed for each occasion. If applicable, the fixed prompt,

including the instructions and predefined actions, was provided in the system message, while the protocol was included in the user message. In this work, we prompted GPT-3 (Brown et al., 2020), GPT-4 (OpenAI, 2023), Gemini (Google, 2024), Claude3 (Anthropic, 2023), and Cohere. Below is the prompt for generating pseudocode based on the given protocol. Note that actions and corresponding descriptions presented in Table 1 are placed at *{actions}*.

*You are an AI that generates Python pseudocode for biology protocols. This pseudocode must accurately describe a complete scientific protocol to obtain a result. You will be provided with the title, description, and steps of the biology protocol, and your task is to convert it to Python pseudocode.*

*You may define the arguments on your own. You must ONLY use these functions.*

*{actions}*

*Do NOT provide any captions. ONLY present the pseudocode and pseudofunctions used inside the code. Present the pseudofunctions at the beginning and then the pseudocode. Do NOT provide any descriptions inside the code.*

*title: {title}*

*description: {description}*

*steps: {steps}*



### 3.4 Metrics and Evaluation

We observe that using automatic metrics (i) necessitates manual annotation of functions and pseudocodes each time, which significantly hampers the automation of the evaluation process. Moreover, evaluating the function and input<sup>5</sup> separately falls short of flexible and comprehensive evaluation in a protocol manner.

To this end, we propose LLAM-EVAL, an automatic, flexible prompt-based framework to evaluate the quality of LLM responses. This framework requires three elements: two input texts (one serving as the baseline and the other as the target) and an evaluator LLM: Llama-3<sup>6</sup>. This method encompasses predefining a set of scores<sup>7</sup>  $S = \{s_1, s_2, \dots, s_n\}$ , prompting Llama-3 to rate the outcomes of a target LLM with that of GPT-4 in the scale of  $S$ , calculating the probability of each score  $p(s_i)$ , and calculating the final score as following. This method is inspired by G-Eval (Liu et al., 2023).

$$\text{score} = \sum_{i=1}^n s_i p(s_i)$$

Llama-3 is prompted to evaluate according to one criterion at a time. The original prompts targeting summarizing tasks are modified to perform evaluation on SPFT. In this work, we evaluate the pseudocode generated by the target LLM based on six criteria: the four original criteria used in G-Eval (Liu et al., 2023) (Coherence, Consistency, Fluency, and Relevance) and two criteria we propose (Precision, and Coverage), considering the context of SPFT. For example, the definition of Coherence is:

*Coherence (1-5) - the overall quality of all lines in the pseudocode. The target pseudocode should not be a rough overview but should provide a precise description of a baseline pseudocode.*

The definitions of other criteria in prompts can be found at Appendix A.2. To automatically implement chain-of-thoughts (CoT) in the evaluation process, we instructed GPT-4 to create specific evaluation steps for each criterion. GPT is capable of producing these evaluation steps by itself (Liu et al., 2023). GPT-4 was given a

<sup>5</sup>Input refers to the function parameters and arguments.

<sup>6</sup>Llama3-70b

<sup>7</sup> $s_1=1$  and  $s_n=5$  is set in this work.

task and evaluation criteria, then prompted to generate the evaluation steps using a form-filling paradigm. An example prompt containing GPT-4 generated instructions for evaluation can be found at Appendix A.2. We also implemented an automatic feedback loop to regenerate the response up to five or ten times if the output did not contain scores. We evaluated using two baselines: the GPT-generated pseudocode and the original protocol.

This approach is not constrained by the output structure of the target models, eliminates the need for manual annotation efforts during the parsing process as required in reference-based metrics, enables a comprehensive evaluation, and thereby makes *ProtoLLM* significantly more flexible and automatic.

To ensure compatibility, we also use conventional reference-based metrics: Normalized Levenshtein distance ( $\mathcal{L}_{dn}$ ) for function names, BLEU (Papineni et al., 2002), precision, recall, and SciBERTScore (O’Donoghue et al., 2023) for function inputs. SciBERTScore is calculated using the encoded **predicted**  $\mathcal{E}(a_i^{\text{pred}})$  and **baseline** values  $\mathcal{E}(a_i^{\text{BL}})$  using the SciBERT (Beltagy et al., 2019) sentence encoder  $\mathcal{E}$ .

$$\text{SciBERTScore} = \frac{1}{N} \sum_{i=0}^N \frac{\langle \mathcal{E}(a_i^{\text{pred}}), \mathcal{E}(a_i^{\text{BL}}) \rangle}{\|\mathcal{E}(a_i^{\text{pred}})\| \|\mathcal{E}(a_i^{\text{BL}})\|}$$

### 3.5 Evaluator LLM Selection

To select a specific LLM as an evaluator, we propose *self-self comparison task* as a baseline, where an LLM generates a pseudocode<sup>8</sup> for a protocol and then evaluates the score using the same LLM against the generated pseudocode. For example, this means evaluating GPT-4 generated pseudocode against the same pseudocode using GPT-4. Our assumption was that the score should be close to the maximum<sup>9</sup> when the baseline and target pseudocode are the same. Our goal was to select the model with the best results as the evaluator. We evaluated each model based on six criteria in Section 3.4. More details in Appendix A.3.

### 3.6 Evaluating LLMs using LLAM-EVAL

Using LLAM-EVAL, we evaluate across three tasks for each model: (1) GPT-4 generated pseudocode as a baseline with predefined actions given in

<sup>8</sup>Pseudocode with pseudofunctions defined at the beginning to be precise.

<sup>9</sup>maximum score  $s_n = 5$  in this work

prompt, (2) the same task with no predefined actions, (3) the original protocol as a baseline with predefined actions. We evaluate GPT variations (Brown et al., 2020; OpenAI, 2023), Llama, Mixtral, Gemma, Cohere, and Gemini (Google, 2024). Details are in Appendix A.4.

### 3.7 Implementation Details

To ensure a fair evaluation of LLMs, we considered additional factors that may affect performance and present several settings. We consider that LLMs tend to perform better when the actions are presented in the same order as in the protocol. While previous work extracted different actions from each protocol<sup>10</sup>, we predefined the actions which is equivalent to shuffling. Also, while using LLAM-EVAL, we encountered instances where the output was a sentence instead of a score (number). To address this issue, we modified the parameters, dataset, and prompts. Further details are in Appendix A.5.

## 4 Analysis

### 4.1 Evaluator LLM Selection

Llama-3 achieved the highest scores across all six tasks, while there were small differences across models (Table 2). We chose Llama-3 as an evaluator, which is free of cost to date. Note that evaluations for other models not presented in the table were not feasible, as numerical responses were not generated. More details are in Appendix A.3.

### 4.2 Evaluating LLMs on SPFT

Our results show that GPT-4o and Cohere+ is a powerful protocol formulator (Table 3). We found our work compatible to previous work (O’Donoghue et al., 2023).

**Is applying domain knowledge an effective strategy for evaluation?** We applied domain knowledge by predefined the finite set of actions performed in biology labs. To evaluate the efficacy of this method, we compare the responses generated with predefined actions included in the prompts to those generated without them (Table 4). The performance is enhanced for most models, with the exception of the *Recall*. Further research should be conducted to explore these findings.

<sup>10</sup>In previous work, this required shuffling, as LLMs presented the pseudofunctions in the same order as in the protocol.

**Can the original protocol itself serve as a baseline?** Evaluation of LLMs in SPFT in previous work requires manual processes and pseudocode extraction step in SPFT. However, evaluation using the original protocol itself completely eliminates the manual processes of pseudofunction evaluation and the GPT-generated pseudocode extraction step, thereby enhancing flexibility and automation. To this end, we evaluate using the original protocol as a baseline. While scores obtained using this approach is not close to the maximum score (Table 3), we observe that the relative ranking of the models remains relevant to the results of using the pseudocode as a baseline.

**Will LLM as an evaluator prefer responses from itself?** It is reported that LLM as an evaluator prefer responses from itself over human responses in text summarization tasks (Liu et al., 2023). Therefore, a potential concern is that the evaluator may prefer outputs from itself regardless of its quality. While results in Table 2 and 4 address this concern, Table 3 shows that Llama-3 as an evaluator does not prefer its outputs over that of GPT-4. Our results suggest that GPT’s preference for its own responses in previous work (Liu et al., 2023) may be a phenomenon unique to GPT.

### 4.3 The BIOPROT 2.0 Dataset

We introduce BIOPROT 2.0, a dataset with scientific protocols and the corresponding pseudocodes with a larger number of datapoints. Previous work highlights that a dataset with these two components can aid protocol formulation of LLMs (O’Donoghue et al., 2023). The pseudocode extracted from protocols are only composed of pseudofunctions (actions) predefined above the previous step, as each model was prompted to use only the provided functions but to define the arguments on their own. The summary of generated pseudocode are in Table 5. This dataset can be used to formulate scientific protocols to achieve a prompted goal using a toolformer like (Schick et al., 2023) chain-of-thought LLM agent (Wei et al., 2023).

## 5 Conclusion

We introduce *ProtocoLLM*, a flexible and automatic framework designed to evaluate LLMs’ capabilities on Scientific Protocol Formulation Tasks (SPFT).

Models	Original Criteria				New Criteria	
	Coherence	Consistency	Fluency	Relevance	Precision	Coverage
GPT-4o	4.95 ± 0.26	4.98 ± 0.25	4.95 ± 0.27	4.93 ± 0.44	4.97 ± 0.23	4.95 ± 0.08
GPT-4	4.98 ± 0.23	4.99 ± 0.19	4.99 ± 0.19	4.99 ± 0.14	4.99 ± 0.18	4.99 ± 0.17
GPT-3.5	4.96 ± 0.23	4.97 ± 0.21	4.77 ± 0.52	4.96 ± 0.25	4.95 ± 0.30	4.99 ± 0.12
<b>Llama-3</b>	<b>5.00 ± 0.02</b>	<b>5.00 ± 0.00</b>	<b>5.00 ± 0.00</b>	<b>5.00 ± 0.00</b>	<b>5.00 ± 0.00</b>	<b>5.00 ± 0.06</b>

Table 2: *Self-Self Comparison Task Results*: We report the mean and standard deviation of scores over ten runs. Values in bold indicate the highest scores for each criterion. Higher values for all metrics represent better performance. Note that a larger dataset was used for this task. Details in Appendix A.3.

Models	Prompt		Original Criteria				New Criteria		Average
	Ac	Pr	Coherence	Consistency	Fluency	Relevance	Precision	Coverage	
GPT-4o	✓	✗	<b>4.10</b> ± 0.79	<b>3.80</b> ± 0.85	<b>3.86</b> ± 0.67	<b>4.32</b> ± 0.71	<b>4.02</b> ± 0.65	<b>4.26</b> ± 0.73	<b>4.06</b>
	✗	✗	<b>4.28</b> ± 0.50	<b>3.94</b> ± 0.64	<b>4.04</b> ± 0.37	<b>4.45</b> ± 0.54	<b>4.18</b> ± 0.41	<b>4.39</b> ± 0.50	<b>4.21</b>
	✓	✓	<b>4.29</b> ± 0.57	<u>4.73</u> ± 0.50	4.42 ± 0.53	<b>4.75</b> ± 0.48	<b>3.90</b> ± 0.48	<b>4.67</b> ± 0.56	<u>4.46</u>
GPT-4 (Baseline)	✓	✗	5.00 ± 0.00	5.00 ± 0.00	5.00 ± 0.08	5.00 ± 0.00	4.99 ± 0.11	5.00 ± 0.00	5.00
	✗	✗	5.00 ± 0.00	5.00 ± 0.00	5.00 ± 0.04	5.00 ± 0.03	5.00 ± 0.03	5.00 ± 0.00	5.00
	✓	✓	4.32 ± 0.53	4.70 ± 0.58	4.53 ± 0.51	4.75 ± 0.44	3.99 ± 0.29	4.67 ± 0.48	4.49
GPT-3.5	✓	✗	3.61 ± 0.97	3.51 ± 1.02	3.58 ± 0.85	4.11 ± 0.78	3.82 ± 0.73	3.90 ± 0.83	3.75
	✗	✗	3.83 ± 0.82	3.71 ± 0.81	3.76 ± 0.68	4.19 ± 0.64	3.96 ± 0.57	3.97 ± 0.71	3.90
	✓	✓	4.13 ± 0.65	<b>4.76</b> ± 0.49	<u>4.48</u> ± 0.52	<u>4.69</u> ± 0.49	3.79 ± 0.58	<u>4.49</u> ± 0.67	4.39
Llama3-8b	✓	✗	2.25 ± 1.00	1.93 ± 0.99	2.27 ± 0.83	2.39 ± 1.08	2.61 ± 0.96	2.56 ± 1.09	2.33
	✗	✗	2.90 ± 0.89	2.69 ± 0.92	3.02 ± 0.88	3.41 ± 0.81	3.47 ± 0.70	3.19 ± 0.82	3.12
	✓	✓	2.80 ± 1.02	3.00 ± 1.26	3.10 ± 1.00	3.39 ± 1.09	2.93 ± 0.92	3.27 ± 1.06	3.08
Llama3-70b	✓	✗	3.61 ± 0.94	3.14 ± 1.10	3.53 ± 0.82	3.73 ± 0.97	3.72 ± 0.70	3.77 ± 0.79	3.58
	✗	✗	<u>3.98</u> ± 0.64	<u>3.72</u> ± 0.75	3.92 ± 0.49	<u>4.20</u> ± 0.57	<u>4.03</u> ± 0.36	<u>4.09</u> ± 0.53	<u>3.99</u>
	✓	✓	4.02 ± 0.75	4.17 ± 0.98	4.15 ± 0.66	4.37 ± 0.74	3.78 ± 0.59	4.25 ± 0.69	4.12
Mixtral	✓	✗	3.41 ± 1.03	2.90 ± 1.13	3.57 ± 0.83	3.36 ± 1.14	3.68 ± 0.77	3.54 ± 0.93	3.41
	✗	✗	3.95 ± 0.68	3.68 ± 0.79	3.94 ± 0.53	4.18 ± 0.66	4.05 ± 0.43	4.00 ± 0.61	3.97
	✓	✓	4.06 ± 0.69	4.32 ± 0.84	4.28 ± 0.59	4.37 ± 0.71	3.88 ± 0.44	4.31 ± 0.70	4.21
Gemma-7b	✓	✗	3.06 ± 0.97	2.81 ± 1.03	3.47 ± 0.86	3.52 ± 0.93	3.55 ± 0.78	3.19 ± 0.89	3.27
	✗	✗	2.93 ± 0.85	2.66 ± 0.88	3.63 ± 0.76	3.61 ± 0.71	3.66 ± 0.61	3.06 ± 0.80	3.26
	✓	✓	3.81 ± 0.75	4.13 ± 0.83	4.25 ± 0.61	4.26 ± 0.75	3.76 ± 0.60	3.94 ± 0.79	4.02
Cohere+	✓	✗	<u>3.95</u> ± 0.74	<u>3.63</u> ± 0.87	<u>3.87</u> ± 0.60	<u>4.11</u> ± 0.74	<u>3.98</u> ± 0.50	<u>4.07</u> ± 0.63	<u>3.94</u>
	✗	✗	3.97 ± 0.60	3.71 ± 0.73	<u>3.95</u> ± 0.46	4.15 ± 0.56	<u>4.03</u> ± 0.38	4.04 ± 0.50	3.98
	✓	✓	4.44 ± 0.52	4.63 ± 0.61	<b>4.53</b> ± 0.52	4.73 ± 0.47	4.04 ± 0.30	4.66 ± 0.49	<b>4.50</b>
Cohere	✓	✗	3.51 ± 0.91	3.06 ± 1.02	3.56 ± 0.74	3.66 ± 0.87	3.71 ± 0.63	3.70 ± 0.76	3.53
	✗	✗	3.71 ± 0.68	3.44 ± 0.83	3.83 ± 0.56	4.05 ± 0.53	3.94 ± 0.41	3.84 ± 0.56	3.80
	✓	✓	<u>3.98</u> ± 0.63	4.11 ± 0.87	4.14 ± 0.51	4.29 ± 0.63	3.83 ± 0.48	4.24 ± 0.64	4.10
Gemini-1.0	✓	✗	2.77 ± 1.09	2.30 ± 1.08	2.90 ± 0.95	2.80 ± 1.10	3.13 ± 0.92	3.15 ± 1.01	2.84
	✗	✗	3.46 ± 0.93	3.22 ± 1.01	3.59 ± 0.83	3.89 ± 0.77	3.80 ± 0.69	3.66 ± 0.79	3.60
	✓	✓	3.37 ± 0.93	3.68 ± 1.11	3.73 ± 0.80	3.87 ± 0.87	3.42 ± 0.78	3.86 ± 0.84	3.66
Gemini-2.0	✓	✗	3.09 ± 1.05	2.53 ± 1.10	3.75 ± 0.70	2.98 ± 1.08	3.63 ± 0.73	3.43 ± 0.89	3.24
	✗	✗	3.88 ± 0.82	3.61 ± 0.91	4.11 ± 0.60	4.13 ± 0.73	4.14 ± 0.54	3.93 ± 0.73	3.97
	✓	✓	3.80 ± 0.80	3.95 ± 0.97	4.30 ± 0.58	4.18 ± 0.72	3.80 ± 0.49	4.11 ± 0.68	4.02
Gemini-1.5	✓	✗	3.02 ± 1.05	2.48 ± 1.02	3.10 ± 0.93	2.97 ± 1.07	3.32 ± 0.84	3.42 ± 0.93	3.05
	✗	✗	4.12 ± 0.66	3.86 ± 0.72	4.03 ± 0.55	4.33 ± 0.62	4.13 ± 0.50	4.21 ± 0.59	4.11
	✓	✓	3.34 ± 0.95	3.62 ± 1.04	3.76 ± 0.76	3.81 ± 0.86	3.36 ± 0.77	3.80 ± 0.84	3.61

Table 3: *ProtoLLM Evaluation Results* of three tasks for each model: (1) GPT-4 generated pseudocode as a baseline with predefined actions given in prompt, (2) the same task with no predefined actions, (3) the original protocol as a baseline with predefined actions. 'Ac' and 'Pr' represent whether the predefined actions and the original protocol were given for evaluation, respectively. We report the mean, standard deviation, and average of scores over five runs. The best and second best performance besides a baseline (GPT-4) for each criterion and task is bolded and underlined, respectively. The scores range from a minimum of 1 to a maximum of 5. Higher values for all metrics represent better performance.

Models	Actions	Precision	Recall	SciBERT	BLEU	$\mathcal{L}_{dn}$
GPT-4o	✓	0.581 ± 0.390	0.548 ± 0.414	0.783 ± 0.111	0.102 ± 0.189	0.216 ± 0.110
	✗	0.600 ± 0.375	0.620 ± 0.373	<b>0.778</b> ± 0.103	0.118 ± 0.188	0.214 ± 0.106
GPT-4 (baseline)	✓	1.00 ± 0.00	1.00 ± 0.00	1.00 ± 0.00	0.905 ± 0.198	0.055 ± 0.129
	✗	1.00 ± 0.00	1.00 ± 0.00	1.00 ± 0.00	0.911 ± 0.173	0.021 ± 0.043
GPT-3.5	✓	0.817 ± 0.308	0.425 ± 0.404	0.766 ± 0.115	0.102 ± 0.205	<b>0.205</b> ± 0.117
	✗	0.732 ± 0.357	0.572 ± 0.378	0.742 ± 0.099	0.099 ± 0.178	<b>0.200</b> ± 0.106
Llama3-8b	✓	0.763 ± 0.323	0.708 ± 0.411	0.801 ± 0.128	0.135 ± 0.329	0.413 ± 0.351
	✗	0.759 ± 0.322	0.570 ± 0.352	0.744 ± 0.100	0.075 ± 0.174	0.242 ± 0.133
Llama3-70b	✓	0.825 ± 0.319	<b>0.917</b> ± 0.220	<b>0.883</b> ± 0.136	<b>0.563</b> ± 0.464	0.287 ± 0.203
	✗	0.812 ± 0.268	<b>0.769</b> ± 0.260	0.772 ± 0.097	<b>0.161</b> ± 0.210	0.206 ± 0.095
Mixtral	✓	0.855 ± 0.280	0.605 ± 0.393	0.784 ± 0.120	0.135 ± 0.288	0.603 ± 0.366
	✗	0.754 ± 0.291	<u>0.735</u> ± 0.290	0.771 ± 0.093	0.130 ± 0.215	0.499 ± 0.261
Gemma-7b	✓	<u>0.911</u> ± 0.249	0.641 ± 0.406	0.838 ± 0.139	0.205 ± 0.342	0.243 ± 0.130
	✗	<b>0.849</b> ± 0.261	0.651 ± 0.337	<u>0.775</u> ± 0.116	0.092 ± 0.180	0.221 ± 0.096
Cohere+	✓	0.646 ± 0.373	0.548 ± 0.352	0.767 ± 0.110	0.075 ± 0.172	0.363 ± 0.300
	✗	0.600 ± 0.366	0.604 ± 0.368	0.744 ± 0.100	0.095 ± 0.153	0.325 ± 0.265
Cohere	✓	0.645 ± 0.361	0.551 ± 0.380	0.717 ± 0.097	0.077 ± 0.193	0.360 ± 0.247
	✗	<u>0.767</u> ± 0.295	0.630 ± 0.314	0.750 ± 0.099	0.091 ± 0.165	<u>0.204</u> ± 0.105
Gemini-1.0	✓	0.852 ± 0.319	0.867 ± 0.313	<u>0.875</u> ± 0.133	<u>0.444</u> ± 0.497	0.410 ± 0.699
	✗	0.758 ± 0.319	0.584 ± 0.360	0.765 ± 0.111	0.112 ± 0.211	0.247 ± 0.182
Gemini-2.0	✓	<b>0.942</b> ± 0.147	0.878 ± 0.288	0.843 ± 0.165	0.342 ± 0.415	0.381 ± 0.254
	✗	0.736 ± 0.350	0.651 ± 0.339	0.758 ± 0.104	0.128 ± 0.197	0.308 ± 0.268
Gemini-1.5	✓	0.889 ± 0.258	<u>0.896</u> ± 0.202	0.814 ± 0.163	0.355 ± 0.461	0.371 ± 0.217
	✗	0.628 ± 0.377	0.682 ± 0.367	0.773 ± 0.101	<u>0.135</u> ± 0.205	0.214 ± 0.116

Table 4: **Evaluation Results Using Reference-Based Metrics.** Comparison with and without predefined actions given in prompts. We report mean and standard deviation of scores over five runs. The best and second best performance for each criterion is bolded and underlined, respectively. Except for  $\mathcal{L}_{dn}$ , higher values for all metrics represent better performance.

Statistic	Value ( $m \pm \sigma$ )
# of protocols	300
Tokens / protocol	812.3 ± 469.9
# of steps	14.81 ± 10.74
Tokens / step	54.28 ± 42.41
Tokens / description	139.0 ± 135.7
Tokens / generated pseudocode	623.8 ± 223.2
# of lines / generated pseudocode	83.06 ± 28.89
# of pseudofunctions / edited pseudocode	10.28 ± 6.582

Table 5: **Statistics of BIOPROT 2.0.** ‘Edited Pseudocode’ refers to the pseudocode that was reformatted, while preserving its content, to obtain the scores presented in Table 4.

Cohere to be particularly effective in formulating scientific protocols. Additionally, we present BIOPROT 2.0, a dataset containing biology protocols and corresponding pseudocodes, which supports LLMs in the formulation and evaluation of SPFT. Our work is extensible to the assessment of LLMs on SPFT across various domains and other fields that require protocol generation for specific goals.

474  
475  
476  
477  
478  
479  
480  
481

This framework prompts the target model and GPT-4 to extract pseudocode from biology protocols using only predefined lab actions, then evaluates the target model’s output using LLAM-EVAL, with the GPT-4 generated pseudocode as a baseline and Llama-3 as the evaluator. Our prompt-based evaluation method, LLAM-EVAL, provides significant flexibility in terms of evaluation models, materials, criteria, and is free of cost. We assess various models, including GPT variants, Llama, Mixtral, Gemma, Cohere, and Gemini, and find GPT and

463  
464  
465  
466  
467  
468  
469  
470  
471  
472  
473



## 6 Limitations

We recognize several limitations. The predefined actions may not encompass all actions performed in a biology labs. The definitions of predefined actions may be incomplete. To precisely define an action, it is necessary to define not only the function but also the function arguments. The number of protocols in BIOPROT 2.0 may be insufficient for evaluation purposes. The performance of *Pro- tocoLLM* may decline outside of biology. Addressing this requires redefining domain-specific actions and exploring other LLMs for diverse fields. Future work should investigate these cross-disciplinary implications. LLMs are continuously evolving due to regular updates. The LLMs used for evaluation in this work might become unavailable in the future. Upgraded versions of LLMs may result in performance degradation and metrics may differ from those obtained using previous models. Due to selecting Llama-3 as the evaluator, our results may be susceptible to its biases and hallucinations. The outcomes when evaluated with models other than Llama-3 are unknown. Future work should investigate the outcomes using different LLMs as an evaluator. Using an API of LLMs as an evaluator, such GPT, is often not free of charge and can be costly. We used GPT-4 generated responses as a baseline; however, it may not accurately represent the ground truth. Future work should explore the implications of employing alternative resources (e.g., manually annotated pseudocodes, responses generated by other models) as the baseline. We observed basic actions classified as NoAction in minor cases. It has been reported that GPT prefers outputs from LLMs, which also produced our evaluation materials including all ground truth and target pseudocodes. This can potentially influence the scores. The four criteria mentioned in G-Eval may not sufficiently fulfill the role of evaluating protocols where real-world validation is crucial. Also, applying these criteria originally designed for summarization tasks may be inappropriate for evaluating SPFT. Even if the protocol pseudocode is successfully synthesized, real-world experiments may fail depending on the person performing the protocol or the condition of the physical equipment, especially in cases that are more complex than stem cell culture or require delicate manual work and experience.

## Ethical Considerations

The use of manually verified protocols in LLMs is strictly prohibited for generating false protocols on platforms like STAR Protocols (Cell Press) and Nature Protocols. Numerous sites also prohibit the use of these protocols in conjunction with any form of AI tool. Our framework can be applied to the protocols of these sites. Although we have endeavored to exclude protocols that can create dangerous substances, there remains the potential for generating protocols that inadvertently produce hazardous products or byproducts.

## Acknowledgements

We would like to thank Karim Md.Adnan for assisting us with the action defining process.

## References

- Walid Ahmad, Elana Simon, Seyone Chithrananda, Gabriel Grand, and Bharath Ramsundar. 2022. [Chemberta-2: Towards chemical foundation models](#).
- Anthropic. 2023. [The claude 3 model family: Opus, sonnet, haiku](#). Technical report, Anthropic. Accessed: 2024-05-16.
- Satanjeev Banerjee and Alon Lavie. 2005. [METEOR: An automatic metric for MT evaluation with improved correlation with human judgments](#). In *Proceedings of the ACL Workshop on Intrinsic and Extrinsic Evaluation Measures for Machine Translation and/or Summarization*, pages 65–72, Ann Arbor, Michigan. Association for Computational Linguistics.
- Iz Beltagy, Kyle Lo, and Arman Cohan. 2019. [SciBERT: A pretrained language model for scientific text](#). In *Proceedings of the 2019 Conference on Empirical Methods in Natural Language Processing and the 9th International Joint Conference on Natural Language Processing (EMNLP-IJCNLP)*, pages 3615–3620, Hong Kong, China. Association for Computational Linguistics.
- Andres M Bran, Sam Cox, Oliver Schilter, Carlo Baldasari, Andrew D White, and Philippe Schwaller. 2023. [Chemcrow: Augmenting large-language models with chemistry tools](#).
- Tom B. Brown, Benjamin Mann, Nick Ryder, Melanie Subbiah, Jared Kaplan, Prafulla Dhariwal, Arvind Neelakantan, Pranav Shyam, Girish Sastry, Amanda Askell, Sandhini Agarwal, Ariel Herbert-Voss, Gretchen Krueger, Tom Henighan, Rewon Child, Aditya Ramesh, Daniel M. Ziegler, Jeffrey Wu, Clemens Winter, Christopher Hesse, Mark Chen, Eric Sigler, Mateusz Litwin, Scott Gray, Benjamin Chess, Jack Clark, Christopher Berner, Sam McCandlish,

582	Alec Radford, Ilya Sutskever, and Dario Amodei.	Schmidt, Ian Foster, Andrew D. White, and Ben	638
583	2020. <a href="#">Language models are few-shot learners.</a>	Blaiszik. 2023. <a href="#">14 examples of how llms can transform materials science and chemistry: a reflection on a large language model hackathon.</a> <i>Digital Discovery</i> , 2:1233–1250.	639
584	Seyone Chithrananda, Gabriel Grand, and Bharath		640
585	Ramsundar. 2020. <a href="#">Chemberta: Large-scale self-supervised pretraining for molecular property prediction.</a>		641
586			642
587		Shuo Jiang, Daniel Evans-Yamamoto, Dennis Bersenev, Sucheendra K. Palaniappan, and Ayako Yachie-Kinoshita. 2024. <a href="#">Protocolcode: Leveraging large language models (llms) for automated generation of machine-readable pcr protocols from scientific publications.</a> <i>SLAS Technology</i> , 29:100134.	643
588	Esin Durmus, He He, and Mona Diab. 2020. <a href="#">FEQA: A question answering evaluation framework for faithfulness assessment in abstractive summarization.</a> In <i>Proceedings of the 58th Annual Meeting of the Association for Computational Linguistics</i> , pages 5055–5070, Online. Association for Computational Linguistics.		644
589			645
590			646
591			647
592			648
593		Seungone Kim, Jamin Shin, Yejin Cho, Joel Jang, Shayne Longpre, Hwaran Lee, Sangdoon Yun, Seongjin Shin, Sungdong Kim, James Thorne, and Minjoon Seo. 2023. <a href="#">Prometheus: Inducing fine-grained evaluation capability in language models.</a>	649
594			650
595	Yuchen Fan, Yantao Liu, Zijun Yao, Jifan Yu, Lei Hou, and Juanzi Li. 2024. <a href="#">Evaluating generative language models in information extraction as subjective question correction.</a>		651
596			652
597			653
598		Matt Kusner, Yu Sun, Nicholas Kolkin, and Kilian Weinberger. 2015. <a href="#">From word embeddings to document distances.</a> In <i>Proceedings of the 32nd International Conference on Machine Learning</i> , volume 37 of <i>Proceedings of Machine Learning Research</i> , pages 957–966, Lille, France. PMLR.	654
599	Jinlan Fu, See-Kiong Ng, Zhengbao Jiang, and Pengfei Liu. 2023. <a href="#">Gptscore: Evaluate as you desire.</a>		655
600			656
601	Gemini Team Google. 2024. <a href="#">Gemini: A family of highly capable multimodal models.</a>		657
602			658
603	Yu Gu, Robert Tinn, Hao Cheng, Michael Lucas, Naoto Usuyama, Xiaodong Liu, Tristan Naumann, Jianfeng Gao, and Hoifung Poon. 2020. <a href="#">Domain-specific language model pretraining for biomedical natural language processing.</a>	R. Kwok. 2010. <a href="#">Five hard truths for synthetic biology.</a> <i>Nature</i> , 463:288–290.	660
604			661
605		Yanis Labrak, Mickael Rouvier, and Richard Dufour. 2024. <a href="#">A zero-shot and few-shot study of instruction-finetuned large language models applied to clinical and biomedical tasks.</a>	662
606			663
607			664
608	Taicheng Guo, Kehan Guo, Bozhao Nan, Zhenwen Liang, Zhichun Guo, Nitesh V. Chawla, Olaf Wiest, and Xiangliang Zhang. 2023. <a href="#">What can large language models do in chemistry? a comprehensive benchmark on eight tasks.</a>		665
609		Jinhyuk Lee, Wonjin Yoon, Sungdong Kim, Donghyeon Kim, Sunkyu Kim, Chan Ho So, and Jaewoo Kang. 2019. <a href="#">BioBERT: a pre-trained biomedical language representation model for biomedical text mining.</a> <i>Bioinformatics</i> , 36(4):1234–1240.	666
610			667
611			668
612			669
613	Tanishq Gupta, Mohd Zaki, N. M. Anoop Krishnan, and Mausam. 2021. <a href="#">Matscibert: A materials domain language model for text mining and information extraction.</a>	Chin-Yew Lin. 2004. <a href="#">ROUGE: A package for automatic evaluation of summaries.</a> In <i>Text Summarization Branches Out</i> , pages 74–81, Barcelona, Spain. Association for Computational Linguistics.	671
614			672
615			673
616			674
617	T Inagaki, Akari Kato, Koichi Takahashi, Haruka Ozaki, and Genki N. Kanda. 2023. <a href="#">Llms can generate robotic scripts from goal-oriented instructions in biological laboratory automation.</a>	Yang Liu, Dan Iter, Yichong Xu, Shuohang Wang, Ruochen Xu, and Chenguang Zhu. 2023. <a href="#">G-eval: Nlg evaluation using gpt-4 with better human alignment.</a>	675
618			676
619			677
620			678
621	Kevin Maik Jablonka, Qianxiang Ai, Alexander Al-Feghali, Shruti Badhwar, Joshua D. Bocarsly, Andres M. Bran, Stefan Bringuier, L. Catherine Brinson, Kamal Choudhary, Defne Circi, Sam Cox, Wibe A. de Jong, Matthew L. Evans, Nicolas Gastellu, Jerome Genzling, María Victoria Gil, Ankur K. Gupta, Zhi Hong, Alishba Imran, Sabine Kruschwitz, Anne Labarre, Jakub Lála, Tao Liu, Steven Ma, Sauradeep Majumdar, Garrett W. Merz, Nicolas Moitessier, Elias Moubarak, Beatriz Mouriño, Brenden Pelkie, Michael Pieler, Mayk Caldas Ramos, Bojana Ranković, Samuel G. Rodrigues, Jacob N. Sanders, Philippe Schwaller, Marcus Schwarting, Jiale Shi, Berend Smit, Ben E. Smith, Joren Van Herck, Christoph Völker, Logan Ward, Sean Warren, Benjamin Weiser, Sylvester Zhang, Xiaoqi Zhang, Ghezal Ahmad Zia, Aristana Scourtas, K. J.	Renqian Luo, Liai Sun, Yingce Xia, Tao Qin, Sheng Zhang, Hoifung Poon, and Tie-Yan Liu. 2022. <a href="#">Biogpt: generative pre-trained transformer for biomedical text generation and mining.</a> <i>Briefings in Bioinformatics</i> , 23(6).	679
622			680
623			681
624			682
625			683
626		Potsawee Manakul, Adian Liusie, and Mark J. F. Gales. 2023. <a href="#">Selfcheckgpt: Zero-resource black-box hallucination detection for generative large language models.</a>	684
627			685
628			686
629			687
630		Odhran O’Donoghue, Aleksandar Shtedritski, John Ginger, Ralph Abboud, Ali Ghareeb, and Samuel Rodrigues. 2023. <a href="#">BioPlanner: Automatic evaluation of LLMs on protocol planning in biology.</a> In <i>Proceedings of the 2023 Conference on Empirical Methods</i>	688
631			689
632			690
633			691
634			692

693	<i>in Natural Language Processing</i> , pages 2676–2694,	Andrew D. White, Glen M. Hocky, Heta A. Gandhi,	747
694	Singapore. Association for Computational Linguistics.	Mehrad Ansari, Sam Cox, Geemi P. Wellawatte, Subarna Sasmal, Ziyue Yang, Kangxin Liu, Yuvraj Singh, and Willmor J. Peña Ccoa. 2023. <i>Assessment of chemistry knowledge in large language models that generate code</i> . <i>Digital Discovery</i> , 2:368–376.	748
695			749
696	OpenAI. 2023. <a href="#">Gpt-4 technical report</a> .		750
697	Kishore Papineni, Salim Roukos, Todd Ward, and Weijing Zhu. 2002. <a href="#">Bleu: a method for automatic evaluation of machine translation</a> . In <i>Proceedings of the 40th Annual Meeting on Association for Computational Linguistics</i> , ACL '02, page 311–318, USA. Association for Computational Linguistics.		751
698			752
699		Tianyi Zhang, Varsha Kishore, Felix Wu, Kilian Q. Weinberger, and Yoav Artzi. 2020. <a href="#">Bertscore: Evaluating text generation with bert</a> .	753
700			754
701			755
702		Wei Zhao, Maxime Peyrard, Fei Liu, Yang Gao, Christian M. Meyer, and Steffen Eger. 2019. <a href="#">Moverscore: Text generation evaluating with contextualized embeddings and earth mover distance</a> .	756
703	Yifan Peng, Shankai Yan, and Zhiyong Lu. 2019. Transfer learning in biomedical natural language processing: An evaluation of bert and elmo on ten benchmarking datasets. In <i>Proceedings of the 2019 Workshop on Biomedical Natural Language Processing (BioNLP 2019)</i> , pages 58–65.		757
704			758
705			759
706		Ruiyang Zhou, Lu Chen, and Kai Yu. 2024. <a href="#">Is LLM a reliable reviewer? a comprehensive evaluation of LLM on automatic paper reviewing tasks</a> . In <i>Proceedings of the 2024 Joint International Conference on Computational Linguistics, Language Resources and Evaluation (LREC-COLING 2024)</i> , pages 9340–9351, Torino, Italia. ELRA and ICCL.	760
707			761
708			762
709	Timo Schick, Jane Dwivedi-Yu, Roberto Dessì, Roberta Raileanu, Maria Lomeli, Luke Zettlemoyer, Nicola Cancedda, and Thomas Scialom. 2023. <a href="#">Toolformer: Language models can teach themselves to use tools</a> .		763
710			764
711			765
712			766
713	Thibault Sellam, Dipanjan Das, and Ankur P. Parikh. 2020. <a href="#">Bleurt: Learning robust metrics for text generation</a> .		
714			
715			
716	P. Shetty, A.C. Rajan, C. Kuenneth, et al. 2023. <a href="#">A general-purpose material property data extraction pipeline from large polymer corpora using natural language processing</a> . <i>npj Computational Materials</i> , 9:52.		
717			
718			
719			
720			
721	Hoo-Chang Shin, Yang Zhang, Evelina Bakhturina, Raul Puri, Mostofa Patwary, Mohammad Shoeybi, and Raghav Mani. 2020. <a href="#">Biomegatron: Larger biomedical domain language model</a> .		
722			
723			
724			
725	Leonid Teytelman, Alexei Stoliartchouk, Lori Kindler, and Bonnie L. Hurwitz. 2016. <a href="#">Protocols.io: Virtual communities for protocol development and discussion</a> . <i>PLOS Biology</i> , 14(8):1–6.		
726			
727			
728			
729	Rian Touchent, Laurent Romary, and Eric de la Clergerie. 2024. <a href="#">Camembert-bio: Leveraging continual pre-training for cost-effective models on french biomedical data</a> .		
730			
731			
732			
733	Alain C. Vaucher, Federico Zipoli, Jonas Geluykens, et al. 2020. <a href="#">Automated extraction of chemical synthesis actions from experimental procedures</a> . <i>Nature Communications</i> , 11:3601.		
734			
735			
736			
737	Alex Wang, Kyunghyun Cho, and Mike Lewis. 2020. <a href="#">Asking and answering questions to evaluate the factual consistency of summaries</a> . In <i>Proceedings of the 58th Annual Meeting of the Association for Computational Linguistics</i> , pages 5008–5020, Online. Association for Computational Linguistics.		
738			
739			
740			
741			
742			
743	Jason Wei, Xuezhi Wang, Dale Schuurmans, Maarten Bosma, Brian Ichter, Fei Xia, Ed Chi, Quoc Le, and Denny Zhou. 2023. <a href="#">Chain-of-thought prompting elicits reasoning in large language models</a> .		
744			
745			
746			



## Appendix

### A BIOPROT 2.0

#### A.1 Data Curation

We used *protocols.io* (Teytelman et al., 2016) API for data collection. Protocols of  $1 \leq score \leq 5$  and  $3 \leq steps$  are collected. The collected data was in a *.json* format, every data point with slight differences in keys. Some protocols were present in the git repository but could not be found when retrieved using the API, and vice versa<sup>11</sup>. Also, even if the file ID in the git repository and the protocol ID retrieved using the API are the same, the dictionary key *number\_of\_steps* may differ<sup>12</sup>. Keywords<sup>13</sup> extracted from the *keywords.txt* file and the descriptions were converted to lowercase temporarily for comparison and scoring. As of May 2024, we collected a total of approximately 15k mirrored public protocols from *protocols.io*'s GitHub before refinement. Protocols were excluded if dictionary key *steps* is empty. Protocols were manually verified by experts in biology. The protocols were removed if they were multiple duplicated files for an identical protocol<sup>14</sup>. For the same title, we score the latest version of the protocol.

#### A.2 Metrics and Evaluation

##### Definitions of Evaluation Criteria

- **Consistency:** Consistency (1-5) - the factual alignment between the source and the target pseudocode. A factually consistent pseudocode contains only statements that are entailed by the source pseudocode. Annotators

<sup>11</sup>The protocol with ID 3737 exists in *protocol.io* but doesn't exist in git repository.

<sup>12</sup>The *number\_of\_steps* for the protocol with ID 10489 is 3 in the git repository but 0 when retrieved using the API.

<sup>13</sup>The keywords are: Biology, Cell, DNA, Protein, Stem Cell, Molecular Biology, Molecular, Gene, Virus, E. coli, cDNA, Agarose, Agarose Gel, in vitro, PCR, NGS, Ethanol, Illumina, Cell Theory, Evolution, Genetics, Homeostasis, Cell Membrane, Mitochondria, Nucleus, Ribosomes, DNA Replication, Mutation, Chromosomes, Gene Expression, Natural Selection, Speciation, Adaptation, Phylogenetics, Ecosystems, Biodiversity, Conservation, Bacteria, Viruses, Fungi, Pathogens, Proteins, Enzymes, Metabolism, Photosynthesis, Gel Electrophoresis, Cloning, CRISPR-Cas9, Neurons, Brain, Synapses, Neurotransmitters, Antibodies, Vaccines, Immune Response, Autoimmunity, Embryogenesis, Stem Cells, Morphogenesis, Regeneration, Pollination, Growth Hormones, Tropisms, Coral Reefs, Oceanic Zones, Marine Conservation, Aquatic Ecosystems, Endangered Species, Habitat Destruction, Conservation Strategies, Rewilding, Genetic Engineering, Bioreactors, Bioinformatics, and Synthetic Biology.

<sup>14</sup>such as protocol ID: 9216

were also asked to penalize summaries that contained hallucinated facts.

- **Fluency:** Fluency (1-5): the quality of the pseudocode in terms of grammar, spelling, punctuation, word choice, and structure.
- **Relevance:** Relevance (1-5) - selection of important information from the source pseudocode. The target pseudocode should include only important information from the source document. Annotators were instructed to penalize summaries which contained redundancies and excess information.
- **Precision:** Precision (1-5) - the exactness and accuracy of the expressions and terminology used in the pseudocode. The target pseudocode should avoid vague or ambiguous terms and should use specific and appropriate terminology that accurately reflects the intended operations and logic.
- **Coverage:** Coverage (1-5) - the extent to which the target pseudocode addresses all aspects of the source pseudocode. The target pseudocode should comprehensively represent all the necessary steps, operations, and details present in the source pseudocode without omitting any critical information.

Note that above are criteria used for evaluation when GPT-generated pseudocode was a baseline. This was slightly modified when evaluating based on original protocol.

**Example LLAM-EVAL Prompt** Below is a prompt evaluating the generated pseudocode from a target LLM based on the criteria Coherence using the GPT-generated pseudocode as the ground truth. The GPT-generated pseudocode for each protocol is placed inside `{{Ground_truth_pseudocode}}`, and the target model-generated pseudocode is placed inside `{{Target_pseudocode}}`.

*You will be given a source pseudocode as a ground truth. You will then be given a target pseudocode which is generated from an identical source of protocol.*

*Your task is to rate the target pseudocode on one metric. Please make sure you read and understand these instructions carefully. Please keep this document open while reviewing, and refer to it as needed.*

*Evaluation Criteria: Coherence (1-5) - the overall quality of all lines in the pseudocode. The target*



*pseudocode should not be a rough overview but should provide a precise description of the ground truth pseudocode.*

*Evaluation Steps:*

*1. Read the Ground Truth Pseudocode: Carefully read and understand the source pseudocode provided as the ground truth. Ensure you comprehend the logic, flow, and details of the algorithm or protocol described.*

*2. Read the Target Pseudocode: Thoroughly read the target pseudocode that needs to be evaluated. Pay attention to the details, structure, and clarity of the pseudocode.*

*3. Compare Against Ground Truth: Compare each line and section of the target pseudocode with the corresponding parts of the ground truth pseudocode. Ensure that all critical steps, variables, and logic present in the ground truth are accurately reflected in the target pseudocode.*

*4. Assess Coherence: Evaluate the overall quality of the target pseudocode based on how well it translates the ground truth. Consider the following aspects: Clarity: Is the pseudocode easy to understand? Completeness: Does it cover all the steps and details present in the ground truth? Precision: Are the descriptions and instructions in the pseudocode precise and unambiguous? Consistency: Are there any contradictions or logical inconsistencies?*

*5. Assign a Coherence Rating (1-5):*

*1 (Poor): The target pseudocode is incomplete, confusing, and lacks most details from the ground truth. 2 (Fair): The target pseudocode is partially complete but has significant gaps and is often unclear. 3 (Good): The target pseudocode covers most details from the ground truth but has some minor inconsistencies or lacks clarity in parts. 4 (Very Good): The target pseudocode is mostly complete and clear, with very few minor issues. 5 (Excellent): The target pseudocode is complete, clear, precise, and fully coherent with the ground truth.*

*Source Pseudocode:*

*{{Ground\_truth\_pseudocode}}*

*Target Pseudocode:*

*{{Target\_pseudocode}}*

*Evaluation Form (scores ONLY):*

*- Coherence:*

### A.3 Evaluator LLM Selection

Models without numerical responses include: Llama3-8b, Llama3-70b, Mixtral, and Gemma.

### A.4 Evaluating LLMs on SPFT

#### Versions of LLMs

Model Name	Call Strings
GPT-4o	gpt-4o
GPT-4	gpt-4
GPT-3.5	gpt-3.5-turbo-1106
Llama3-8b	llama3-8b-8192
Llama3-70b	llama3-70b-8192
Mixtral	mixtral-8x7b-32768
Gemma-7b	gemma-7b-it
Cohere+	command-r-plus
Cohere	command-r
Gemini-1.0	gemini-1.0-pro-001
Gemini-1.5	gemini-1.5-pro-001
Gemini-2.0	gemini-1.0-pro-002

Table 6: **Versions of LLMs.** Exact API call strings for corresponding models.

### A.5 Implementation Details

Except for n and seed, parameters were set to their default values. We used approximately \$1000 for GPT API calls, \$20 for Gemini, and other models were free of cost.

**Counting Tokens** We counted the tokens of the concatenated string of the title, original description, and steps, separated by "\n\n". The reason for this approach is to match the token count with that of the previous work.

**Inconsistencies LLAM-EVAL Outputs** To address this issue, we attempted the following methods: (1) Modified `max_token = 5` to `max_token = 1`: The scores became integers, but the model still generated sentences in addition to scores. (2) Use different versions of the model: Other model variations, such as `gpt-3.5-turbo-1106`, did not enhance the results.