

PSCsystem User Study

Three Tasks & Questionnaires

This user experiment aims to validate whether large language models can enhance the intelligence of multimodal frameworks in protein structure tasks.

Why are we conducting this experiment?

We aim to understand whether the new PSCsystem Pipeline, empowered by large language models, can offer a more intelligent and efficient solution for multimodal protein-structure-related tasks than two classic Pipelines for



Experiment Overview

We will conduct **two tasks** related to protein structure model/task matching and implementation. One task will be a baseline pipeline for model/task matching. Another task will be the PSCsystem model/task matching. The third task will be auto-training the matched model with biodata or asking the model to generate direct results.

Before starting each task, we will provide user guides for the Baseline Matching Pipeline and PSCsystem Pipeline, along with biological examples for data inputting.

Help us to know you!

Background Survey

As the very first step of the experiment, you will complete a user background survey to help us better understand your background and requirements.

Please enter the user background questionnaire by following the link:
<https://forms.gle/3J8UoRuwiXc9RVx7>

Questionnaire Overview

Background Survey: <https://forms.gle/5CpbrrrVtFriZ112A>

System Study: Please express your insights into the framework's usability in the SUS questionnaire and provide your feedback on the workload of this experiment in the TLX questionnaire. Please access them through the following links:

- For **Baseline**:

SUS: <https://forms.gle/mpyDWZXhrh3JTnCA6>

TLX: <https://forms.gle/PHvk7dcnm5miqrnp6>

- For **PSCsystem Pipeline**:

SUS: <https://forms.gle/zw9KeGfYzR4xHiw98>

TLX: <https://forms.gle/jCZTmLkEUbNKsXC96>

3. Enter a Path and then enter the number, and click enter. You will complete one set of testing.

```
Please enter your command (输入您的命令): ESM1v for inversefolding
Please enter the path to the protein's pdb file: jansn/sjansn.pdb
Please enter the number of samples: 4
Running Model name: ESM1v model for results...
python main.py [Model name: ESM1v, sample, jansn/sjansn.pdb, 4]
*****
Searching summary (搜索总结):
Detected Model: Model name: ESM1v
Detected Task: Task name: inversefolding
Model Match: Yes
Task Match: Yes
Command to execute: python main.py [Model name: ESM1v, sample, jansn/sjansn.pdb, 4]
Total processing time: 13.95 seconds
*****
```

4. Repeat three times. Then, copy the script of testing and submit it to the **SUS baseline matching questionnaire**.

Task 2

Explore PSCsystem Pipeline

Instructions for Use

After downloading the compressed file, please run the PSCsystem executable. This will start the code testing. The exploration basically focuses on two parts: model and task matching and sample data processing:

1. Model/Task matching: Enter natural language command in Chinese or English, and click return to process. The command should include the purpose or the task at least. Irrelevant commands, such as “anything is possible,” will not work for model and task matching:

```
Set 1 of 4
Please enter your command (输入您的命令): anything is possible
Manager Suggestion (AI建议): The phrase "anything is possible" can refer to any of the models and tasks mentioned above,
as all of them can be used to predict and analyze proteins and protein complexes.
Processing Time (处理时间): 1.0248980522155762 s
process or not: Y/N
Y
There is no specific model or task mentioned in the given command.

Time for extracting model and task name: 1.4770143032073975s
```

2. Enter natural language commands that contain data or data directories in Chinese or English, and click return to process. If the AI interpretation does not match what you meant, feel free to add clarification or correct your previous command and return.

```
Please enter your data or data path(输入数据或文件路径): Given AA sequence SJJABJSNDIINSNNDBDB
Detected:This input represents direct data (like text) without punctuations. It does not resemble a file or directory pa
th.
Time:0.8496992588043213
Enter clarification:
process or not: Y/N
Y

*****Recorded the input*****
"Certainly! Here is part of the result. Token1: 0.425026, Token2: 3.189499, Token3: 2.111793, Token4: 3.782089. For the
complete result, please refer to the generated output file."
Time: 1.3603460788726807s
```

3. After you see ******all task complete******, you can copy everything and submit to the **SUS PSCsystem questionnaire**.

```
****ALL TASKS COMPLETE****
Enter anything to end all sessions
|
```

Examples for Task 3 Data Input

1. "Use the model to predict the mutant score for a protein, the path to pdb file 'protein.pdb' and fasta file: 'protein.fasta'."
2. "I would like to use the model for an inverse folding task. Here is the path to the pdb file: protein.pdb'. Additionally, sample it twice."
3. "Use the model to predict the mutant score for a protein. Here is the sequence: [sequence], and the mutant csv file is 'result.csv'."
4. "Use the model to perform the task. Here is the sequence: [sequence]. Additionally, the corresponding mutant file is located at 'F7YBW8/F7YBW8.csv', and the msa file path is 'F7YBW8/F7YBW8.a3m'."
5. "Test the EC dataset, using the attr weight to test it."
6. "Use the model to predict mutant score on GFP protein dataset. Here are the pdb file, msa file, and mutant tsv file: 'Dataset/GFP.pdb', 'Dataset/GFP.a3m', 'Dataset/GFP.tsv'"
- 7.....

We thank you for your participation and feedback!
Your input is valuable to our research!

For further comments or concerns, please contact fh2450@tc.columbia.edu