

## **Evaluation**

#### **Criteria:**

- 1. Procedure
- 2. Personalization
- 3. Proactivity

## **Criteria Standard:**

The standard of score 0~5 for each metric

## Query and Model output in uniform format

User Profile and Toolutilizing Preference

## **Analysis format**

Key-point for analysis

# GPT-40

Evaluating



### **Key-point for Procedure:**

- 1. Completeness; 2. Avoid Unnecessary Action;
- 3. Call the tool accurate; 4. Summary the query clearly and comprehensive

## **Key-point for Personalization:**

- 1. Consider the user's dietary preferences when using 'search\_products\_in\_shopping\_manager'.
- 2. Consider user's occupation and budgets for shopping.

## **Key-point for Proactivity:**

- 1. Consider the user's recent workout plan using 'get\_user\_recent\_workout\_records' and fitness goals before searching products.
- 2. Consider the user's recent health status using 'get\_recent\_health\_and\_mood\_summary' before searching products.

•••••

```
"Procedure": {
     "Keypoints for Procedure": {......},
     "Final Assessment": {"anlysis": ".....", "score": 5}
"Personalization": {
     "Keypoints for Personalization": {......},
     "Final Assessment": { "analysis": ".....", "score": 5}
"Proactivity": {
     "Keypoints for Proactivity": {
          "Consider the user's recent workout plan .....": {
               "analysis": "The assistant does not reference .....",
               "score": 0
                                   0 if not satisfied
                                   1 if part satisfied
                                   2 if satisfied
     "Final Assessment": {
          "Analysis": "The assistant demonstrates minimal
proactivity, missing several opportunities to .....",
          "score": 2
```