

介绍Smart Prompt框架



汇报人：庄表伟



时间：2024.11



CONTENTS

目录

01

如何理解这一轮AI浪潮？

02

SmartPrompt的设计思路

03

SmartPrompt功能示例

04

SmartPrompt的后续发展计划

05

结语

如何理解这一轮AI浪潮？

01



什么是LLM?

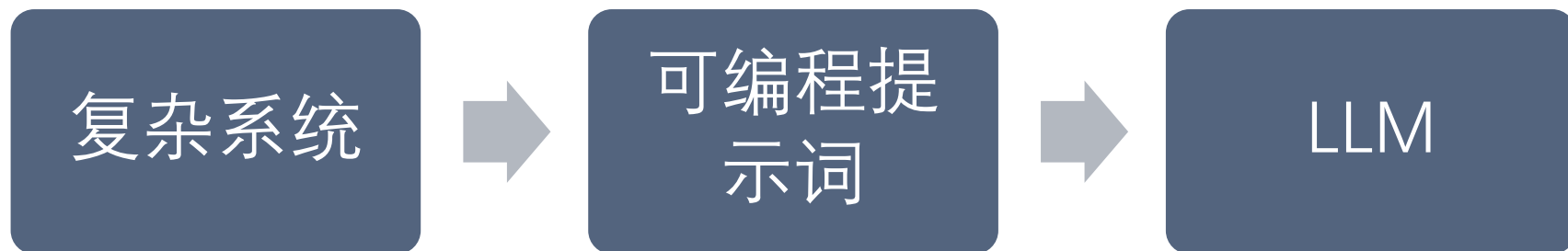
- **定义:**
 - LLM是通过深度学习训练的大型神经网络模型，专门用于理解和生成自然语言。
- **工作原理:**
 - 通过大规模的文本数据学习语言结构、语法和语义。
 - 具备处理多种语言任务的能力，如翻译、文章生成、问题回答等。
- **应用场景:**
 - AI对话系统
 - 自动翻译
 - 内容创作
 - 自然语言处理 (NLP)

什么是提示词？

- **定义：**
 - 提示词是输入给LLM的文本，用于引导模型生成相关输出或执行特定任务。
- **功能与重要性：**
 - 清晰、具体的提示词能提高模型的回答质量和准确性。
 - 模糊、不明确的提示词可能导致不准确的输出。
 - 不同的模型（开源的/闭源的，小的/大的），针对不同的任务，不同的提示词，可能有不同的效果
- **示例：**
 - 提示词1：“请解释什么是机器学习。”
 - 提示词2：“翻译以下句子为法语：‘Hello, how are you?’”
 - 提示词3：“告诉我这个出错信息是什么意思，我要怎么解决？”

什么是可编程的提示词？

- **定义：**本质上是一种传统架构与AI架构集成起来的**架构模式**
 - 并非所有的问题，都适合用LLM来解决
 - 在选择调用何种LLM时，如何才能灵活切换？



SmartPrompt的设计思路

02



Ruby DSL

- 什么是DSL?
 - 领域特定语言（DSL）是一种专门为特定问题域设计的编程语言或语法。
 - 与通用编程语言（如Ruby、Python等）不同，DSL通常具有简化的语法，更贴近业务领域的需求。
- Ruby中的DSL特点：
 - 简洁性与可读性：Ruby的语法非常简洁、优雅，适合快速开发领域特定的语言。
 - 高效的抽象：通过Ruby的元编程功能，可以创建简洁的DSL，隐藏复杂性。
 - 自然语言风格：Ruby的语法使得DSL通常看起来像自然语言，易于理解和使用。
- 常见应用：
 - 配置文件：DSL常用于编写配置文件，提升可读性和简洁性。例如，Ruby on Rails的路由配置。
 - 测试框架：RSpec是Ruby中的一个DSL，用于编写易于理解的测试用例。
 - 查询语言：ActiveRecord的查询接口提供了一个Ruby风格的DSL，用于构建数据库查询。

```
# 示例: RSpec测试DSL
describe "Calculator" do
  it "adds two numbers" do
    expect(1 + 1).to eq(2)
  end
end
```

```
# 使用DSL定义任务
scheduler = TaskScheduler.new

scheduler.task("task_1") do
  puts "Executing task 1..."
end

scheduler.task("task_2") do
  puts "Executing task 2..."
end

scheduler.task("task_3") do
  puts "Executing task 3..."
end

# 运行所有任务
scheduler.run
```


用ERB作为提示词模板

- **什么是ERB?**

- ERB (Embedded Ruby) 是Ruby语言中的一种模板引擎，它允许将Ruby代码嵌入到HTML或其他文本格式中。
- 通过ERB模板，开发者可以将动态内容嵌入到静态页面中，生成动态HTML页面或其他格式的输出。

- **ERB 的工作原理:**

- 模板语法：在ERB模板中，Ruby代码通过特定的标签（`<% %>`）嵌入到文本中。
- `<%= %>`：执行Ruby代码并将结果插入到输出中。
- `<% %>`：执行Ruby代码但不将结果插入到输出中。
- ERB解析器会将模板中的Ruby代码解析并执行，生成最终的HTML或文本输出。

- **常见应用场景:**

- Web开发：在Rails等Web框架中，ERB广泛用于生成动态HTML页面，特别适合与数据库或用户输入交互。
- 邮件模板：生成动态电子邮件内容时，使用ERB模板可以嵌入动态数据（如用户名称、订单信息等）。
- 报告生成：根据数据生成动态报告，ERB模板可以根据输入的数据结构动态生成报告内容。

模板示例：

Translate the following text from
`<%= source_language %>` to `<%= target_language %>`:

`<%= text %>`

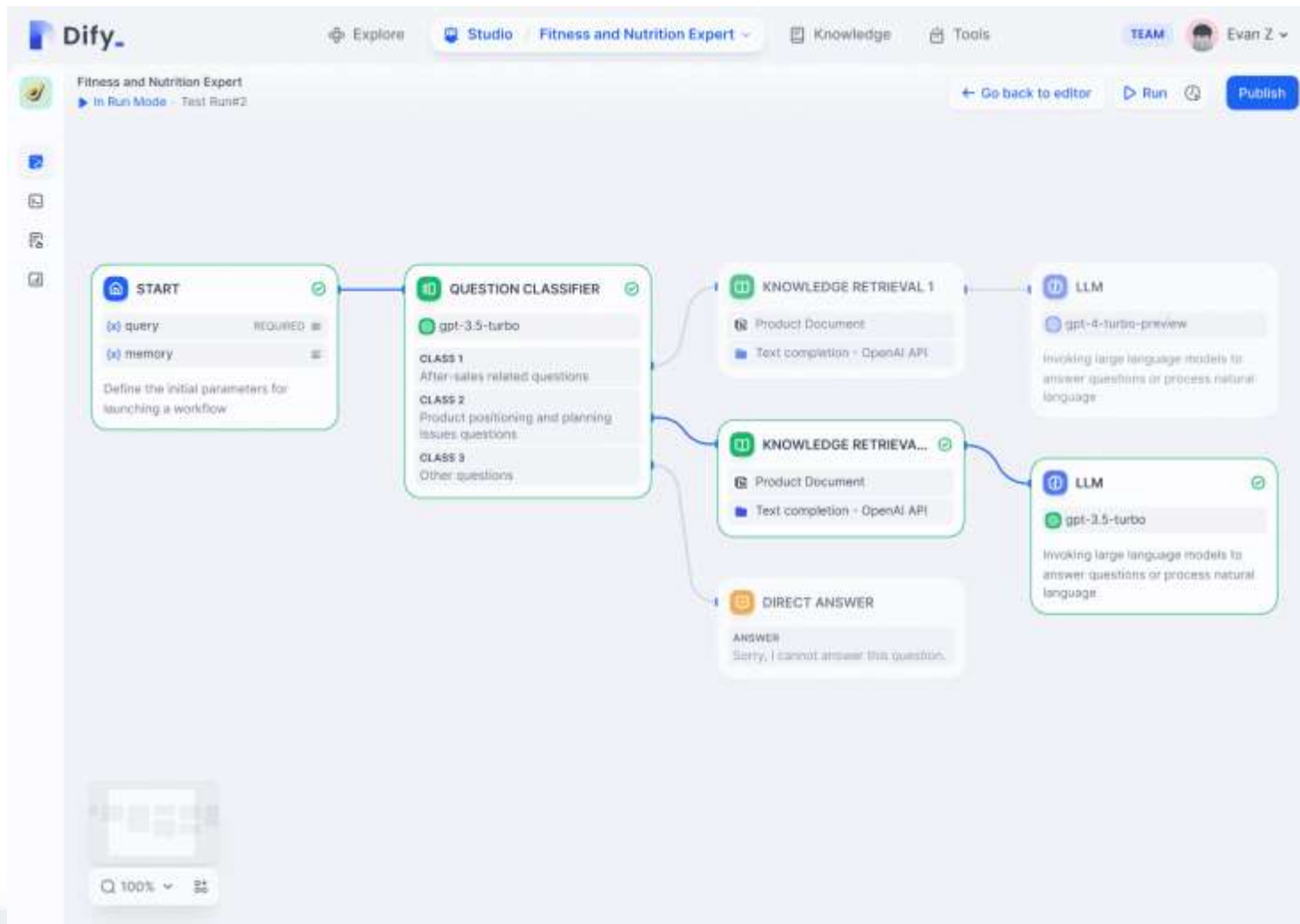
The translation should be accurate, maintain the original meaning, and be appropriate for the cultural context of the target language.

多模型切换

```
workers > get_code.rb
1  SmartPrompt.define_worker :get_code do
2    use "siliconflow"
3    sys_msg "You are a helpful programmer."
4    model "Qwen/Qwen2.5-Coder-7B-Instruct"
5    prompt :generate_code, {
6      name: params[:name],
7      description: params[:description],
8      input: params[:input],
9      output: params[:output]
10   }
11   code_text = safe_send_msg
12   if code_text.include?("Failed to call LLM after")
13     code_text
14   else
15     model "meta-llama/Meta-Llama-3.1-8B-Instruct"
16     prompt :get_code, {code_text: code_text}
17     safe_send_msg
18   end
19 end
```

选择不同的模型

Worker + Workflow



低代码的工作流看起来好看，其实还不如自己写代码来得快。

SmartPrompt功能示例

03



应用的目录结构

- config
 - llm_config.yml
- log
 - log.txt
- templates
 - *.erb
- workers
 - *.rb
- hello.rb
- code.rb
- translat.rb
- embedding.rb
- rag.rb
-

```
# gem install smart_prompt
```

```
config > ! llm_config.yml
1  adapters:
2    openai: OpenAIAdapter
3    ollama: OllamaAdapter
4  logger_file: ./log/log.txt
5  llms:
6    siliconflow:
7      adapter: openai
8      url: https://api.siliconflow.cn/v1/
9      api_key: ENV["APIKey"]
10     default_model: Qwen/Qwen2.5-7B-Instruct
11  llamacpp:
12     adapter: openai
13     url: http://localhost:8080/
14  ollama:
15     adapter: ollama
16     url: http://localhost:11434/
17     default_model: qwen2.5
18  default_llm: siliconflow
19  worker_path: "./workers"
20  template_path: "./templates"
```


翻译

```
translate_txt.rb > ...
1 require "../smart_prompt/lib/smart_prompt"
2 engine = SmartPrompt::Engine.new("../config/llm_config.yml")
3 Words = {}
4
5 def translate_text(engine, text)
6   checked = true
7   count = 0
8   result = ""
9   while checked == true && count < 10
10    puts "try #{count}"
11    result = engine.call_worker(:categorized_translation, {
12      text: text,
13      source_language: "English",
14      target_language: "Chinese"})
15    checked =
16      (result.include?("(") && result.include?("}")) ||
17      (result.include?("```") || result.include?("```json")) ||
18      result.include?("**Output**") ||
19      result.include?("**输出**") ||
20      (result.split("\n").size > text.split("\n").size)
21    count += 1
22  end
23  result
24 end
25
26 book = File.read("../OpenLife.txt")
27 new_book = File.new("../OpenLife_new.txt", "w+")
28
29 book.split("\n").each do |line|
30   if line.strip.empty?
31     result = line
32   elsif Words.has_key?(line)
33     result = Words[line]
34   else
35     result = translate_text(engine, line)
36     Words[line] = result
37   end
38   new_book.puts result
39 end
40 new_book.close
```

```
SmartPrompt.define_worker :categorized_translation do
  use "siliconflow"
  model "Qwen/Qwen2.5-Coder-7B-Instruct"
  prompt :cat_trans, {text: params[:text], source_language: params[:source_language]}
  json = send_msg
  model "Qwen/Qwen2.5-72B-Instruct-128K"
  if word?(params[:text])
    prompt :senior_translator, {
      json: json,
      text: params[:text],
      source_language: params[:source_language],
      target_language: params[:target_language]
    }
    result = send_msg
  else
    prompt :senior_translator2, {
      json: json,
      text: params[:text],
      source_language: params[:source_language],
      target_language: params[:target_language]
    }
    result = send_msg
  end
end
```

翻译

templates > cat_trans.erb

```
1 # Command:
2 **Input**: <%= source_language %>
3 **Output**: Analysis
4 1. **Category**: Based on the content of a given text, determine to which of t
5     - Abbreviations and acronyms
6     - Dates and times (No translation required)
7     - Numerical value (No translation required)
8     - Monetary units
9     - Units of measurement
10    - Chapter headings
11    - Proper nouns
12    - Literature citations
13    - Fixed expressions, phrases, idioms and expressions
14    - Titles of literary and artistic works
15    - Legal and regulatory provisions
16    - Technical terms
17    - File name (No translation required)
18    - Dialogue, inflections, slang and colloquial expressions
19    - Historical and cultural contexts
20    - Formatted texts, formulas and symbols
21    - URI/URL (No translation required)
22    - Other complex texts
23 2. **Instances**: Under this category, enter actual examples from the text.
24
25 # Example:
26
27 **Input**: NASA's mission was launched on March 1, 2023, at 10:30 PM.
28 **Output**(JSON format):
29 {
30   "category": ["Abbreviations and acronyms", "Dates and times"],
31   "instances": {
32     "Abbreviations and acronyms": ["NASA (name of organization)"],
33     "Date & times": ["March 1, 2023", "10:30 PM"]
34   }
35 }
36
37 **Input**: <%= text %>
38 **Output**(JSON format):
```

templates > senior_translator.erb

```
1 # Analysis:
2 <%= json %>
3
4 # Command:
5 1. Output only translated content
6 2. Please follow the input file format strictly and do not add anything else.
7 3. Based on the results of the above analysis, translate the following <%= source_language %>, into <%= target_language %>.
8
9 # Example:
10 **Input**: , John Wiley & Sons.
11 **Output**: , 约翰·威利父子.
12
13 # Task:
14 **Input**: <%= text %>
15 **Output**(Plain text):
```

templates > senior_translator2.erb

```
1 # Analysis:
2 <%= json %>
3 # Command:
4 0. The input just one word
5 1. Output only translated content
6 2. Please follow the input file format strictly and do not add anything else.
7 3. Based on the results of the above analysis, translate the following <%= source_language %>, into <%= target_language %>.
8 # Example:
9 **Input**: world
10 **Output**: 世界
11 # Task:
12 **Input**: <%= text %>
13 **Output**(Plain text):
```


抓取新闻

```
workers >  get_news.rb
1  SmartPrompt.define_worker :get_news do
2    url = params[:text]
3    html = call_worker(:download_page, {url: url})
4    text = call_worker(:html_to_text, {html: html})
5    use "siliconflow"
6    sys_msg "You're a journalist familiar with open source-related news coverage."
7    model "Qwen/Qwen2.5-7B-Instruct"
8    prompt :analyzing_news_content, {news: text}
9    news_json = safe_send_msg
10   f = File.open("news.json", "w+")
11   f.puts news_json
12   f.close
13   prompt :generate_sql2, {news: news_json}
14   sql = safe_send_msg
15   sql
16 end
```

templates >  analyzing_news_content.erb

```
1  Based on the content of the page provided below, the content of N news items
2  was analyzed and presented in the following format:
3
4  ```json
5  [
6    {
7      original_title: "",
8      chinese_title: "",
9      date: "",
10     summary: "",
11     categories: ["", ""]
12   },
13 ]
14  ```
15
16  Input news:
17  <%= news %>
18
19
20
```

生成代码

```
code.rb
1 require "../smart_prompt/lib/smart_prompt"
2 engine = SmartPrompt::Engine.new("../config/llm_config.yml")
3 result = engine.call_worker(:get_code, {
4   name: "calculate_triangle_area",
5   description: "calculates the area of a triangle",
6   input: "the base and height of the triangle",
7   output: "the area as a float"
8 })
9
10 if result.include?("Failed to call LLM after")
11   puts result
12 else
13   code_str = result + "\n" + "puts calculate_triangle_area(8, 10)"
14   eval(code_str)
15 end
```

```
tes > <> generate_code.erb
Write a Ruby function called "<%= name %>" that <%= description %>.
The function should accept <%= input %> as input and return <%= output %>.
Just return the source code to me, no other explanations needed.
```

```
workers > get_code.rb
1 SmartPrompt.define_worker :get_code do
2   use "siliconflow"
3   sys_msg "You are a helpful programmer."
4   model "Qwen/Qwen2.5-Coder-7B-Instruct"
5   prompt :generate_code, {
6     name: params[:name],
7     description: params[:description],
8     input: params[:input],
9     output: params[:output]
10  }
11  code_text = safe_send_msg
12  if code_text.include?("Failed to call LLM after")
13    code_text
14  else
15    model "meta-llama/Meta-Llama-3.1-8B-Instruct"
16    prompt :get_code, {code_text: code_text}
17    safe_send_msg
18  end
19 end
20
```

```
tes > <> get_code.erb
Remove all non-code text and return only the code. Not even characters like ``.

<%= code_text %>
```

```
def calculate_triangle_area(base, height)
  (base * height) / 2.0
end
puts calculate_triangle_area(8, 10)
40.0
```

SmartPrompt的后续发展计划

04



支持多模态的大模型

模型广场

类型

对话

生图

嵌入

重排序

语音

视频

价格

只看免费

可用赠费

系列 / 厂商

腾讯混元

Qwen2.5

Qwen2

书生万象

书生浦语

NVIDIA

FLUX.1

DeepSeek

智谱 AI

智源研究院

网易有道

电信

Fish Audio

FunAudioLLM

LLama3.1

LLama3

零一万物

Stability AI

Google

Mistral AI

字节跳动

腾讯 ARC

InstantX

隐藏筛选器

请输入模型名称

Qwen/QwQ-32B-Preview

Qwen2.5 | ¥1.26 / M Tokens

QwQ-32B-Preview是Qwen 最新的实验性研究模型。专注于提升AI推理能力。通过探索语言混合、递归推理等复杂机制。主要优势包括强大的推理分析能力、数学和编程能力。...

对话 32B 32K

Qwen/Qwen2.5-Coder-32B-Instruct

Qwen2.5 | ¥1.26 / M Tokens

Qwen2.5-Coder-32B-Instruct 是基于 Qwen2.5 开发的代码特定大语言模型。该模型通过 5.5 万亿 tokens 的训练。在代码生成、代码推理和代码修复方面都取得了显著提升。它是当...

对话 FIM Coder 32B 32K

Tencent/Hunyuan-A52B-Instruct

腾讯混元 | ¥21 / M Tokens

混元大模型 (Hunyuan-Large) 是业界最大的开源 Transformer 架构 MoE 模型。拥有 3890 亿总参数量和 520 亿激活参数量。该模型采用了高质量合成数据训练。KV 缓存压缩。专...

对话 MoE 32K

fishaudio/fish-speech-1.4

Fish Audio | ¥105 / M UTF-8 bytes

Fish Speech V1.4 是一个先进的文本转语音 (TTS) 模型。在 70 万小时的多语言音频数据集上训练而成。该模型支持 8 种语言。包括约 30 万小时的英语和中文数据。以及约 2 万...

语音 多语言

stabilityai/stable-diffusion-3-5-large

Stability AI | 免费

Stable Diffusion 3.5 Large 是 Stable Diffusion 系列中最强大的基础模型。拥有 80 亿参数。该模型具有卓越的图像质量和更好的提示遵循能力。特别适合在 1 兆像素分辨率下进行...

生图 8B

stabilityai/stable-diffusion-3-5-large-turbo

Stability AI | ¥0.0032 / M px / Steps

Stable Diffusion 3.5 Large Turbo 是 Stable Diffusion 3.5 Large 的高速版本。该模型能够在仅 4 步内生成高质量图像。并具有出色的提示遵循能力。使其比 Stable Diffusion 3.5 Large ...

生图 8B

nvidia/Llama-3.1-Nemotron-70B-Instruct

NVIDIA | ¥4.13 / M Tokens

Llama-3.1-Nemotron-70B-Instruct 是由 NVIDIA 定制的大型语言模型。旨在提高 LLM 生成的响应对用户查询的帮助程度。该模型在 Arena Hard、AlpacaEval 2 LC 和 GPT-4-Turbo MT-...

对话 70B 32K

Qwen/Qwen2-VL-72B-Instruct

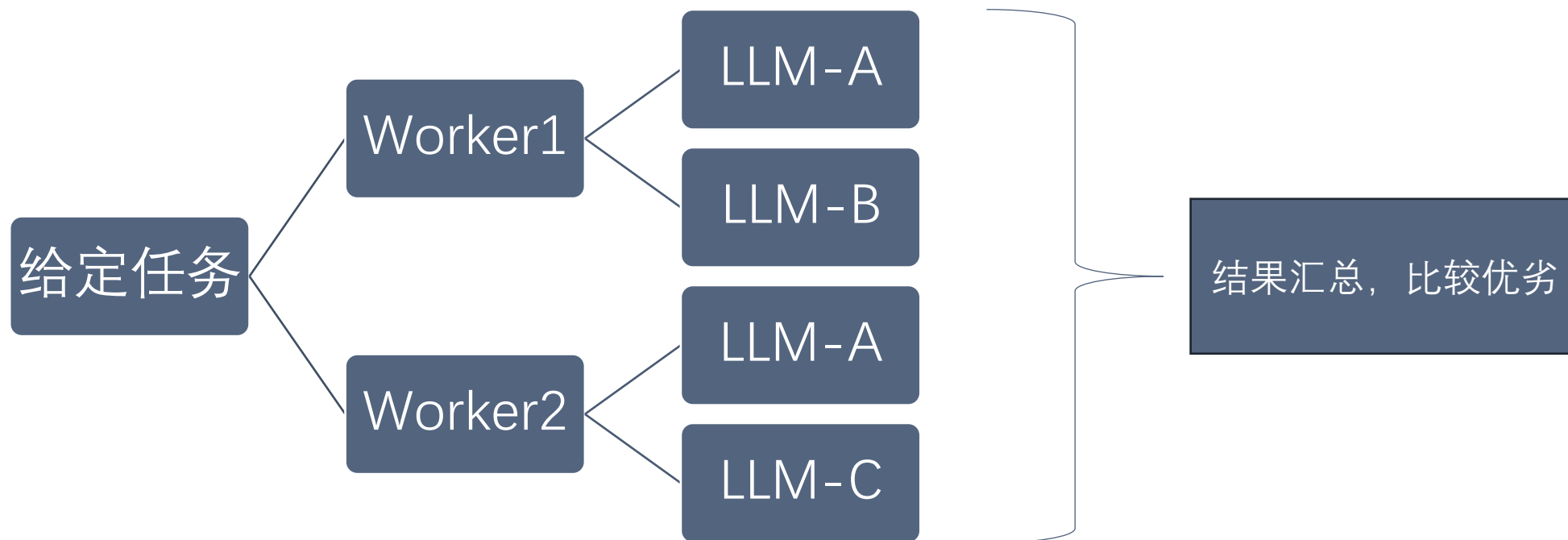
Qwen2 | ¥4.13 / M Tokens

Qwen2-VL 是 Qwen-VL 模型的最新迭代版本。在视觉理解基准测试中达到了最先进的性能。包括 MathVista、DocVQA、RealWorldQA 和 MTVQA 等。Qwen2-VL 能够理解超过 2...

对话 视觉 72B 32K

zhuangbiaowei
135*****931

多个Worker、多个模型并发



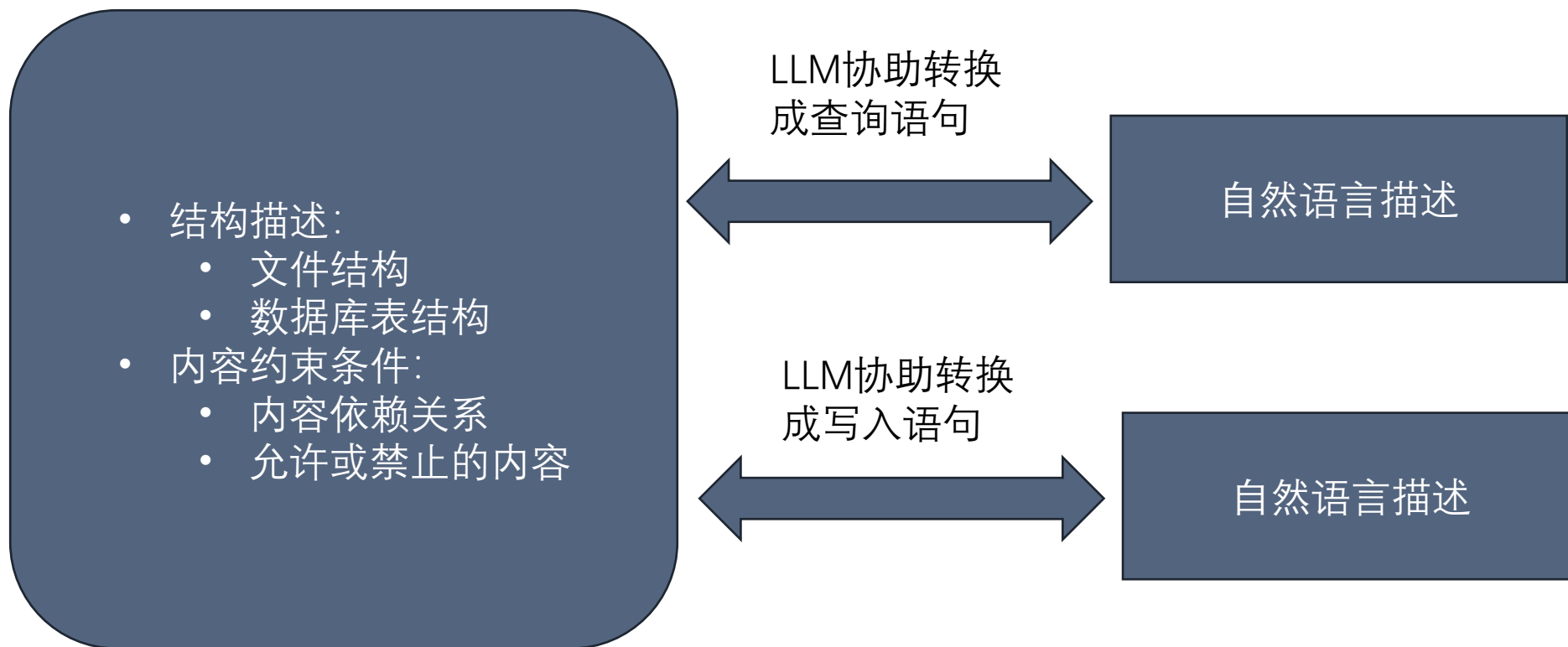
对于Worker的输出结果进行调优

现在的worker里，写死了模型名称，也有明确调用的模板。

如果能够定义：“理想的输出结果”，然后针对可选的模板，反复优化模板的提示词，直到找到最佳组合。

效果就会更加理想。

系统级Worker：文件管理、数据库管理等



结语

05



目前，这还是一个非常早期的项目，因为很多想要实现的功能，还没有找到自己觉得满意的写法。
估计在上述设想实现之后，再正式开源。

https://github.com/zhuangbiaowei/smart_prompt

<https://github.com/zhuangbiaowei/SmartPromptDemo>

谢谢大家



汇报人：庄表伟



时间：2024.11

