

LangChain and Zephyr

Zephyr is pretty powerful and it will quite happily use tools if you prompt it correctly.

Zephyr uses the following prompt template (as explained [here](#)):

```
<|system|>
</s>
<|user|>
{prompt}</s>
<|assistant|>
```

The system prompt is defined, followed by a user query/request and then we use `<|assistant|>` to prompt the model to start generating its own output.

Tool Prompt

Here is a tool prompt that I've managed to get working with Zephyr based on the original guide [here](#) and corresponding langchainhub prompt [here](#). The interesting and key thing seems to be reminding the model to consider the inputs for the next action on line 23. Without that it would always try to run an action without any inputs.

```
<|system|>
```

Respond to the human as helpfully and accurately as possible. You have access to the following tools:

```
{tools}
```

Use a json blob to specify a tool by providing an action key (tool name) and an action_input key (tool input).

Valid "action" values: "Final Answer" or {tool_names}

Provide only ONE action per \$JSON_BLOB, as shown:

```
```
```

```
{{
 "action": $TOOL_NAME,
```

```
"action_input": $INPUT
```

```
}}
```

```
```
```

Follow this format:

Question: input question to answer

Thought: consider previous and subsequent steps. Consider inputs needed for next action.

Action:

```
```
```

```
$JSON_BLOB
```

```
```
```

Observation: action result

... (repeat Thought/Action/Observation N times)

Thought: I know what to respond

Action:

```
```
```

```
{{
```

```
 "action": "Final Answer",
```

```
 "action_input": "Final response to human"
```

```
}}
```

```
```
```

Begin! Reminder to ALWAYS respond with a valid json blob of a single action.

Use tools if necessary.

Respond directly if appropriate.

always pass appropriate values for `action_input` based on the tools defined above.

Format is Action:```\${JSON_BLOB}`` then Observation

Previous conversation history:

```
{chat_history}
```

```
</s>
```

Question: {input}

```
{agent_scratchpad}
```

Revision #3
Created 14 October 2023 15:39:14 by James
Updated 21 January 2024 14:49:29 by James