# Design bots using Flowise

# Last amended: 19th Jan, 2025
# My folder: C:\Users\ashok\OneDrive\Documents\flowise
# Flowise Book created by Community at this link

# Table of Contents

## A. Simple demo

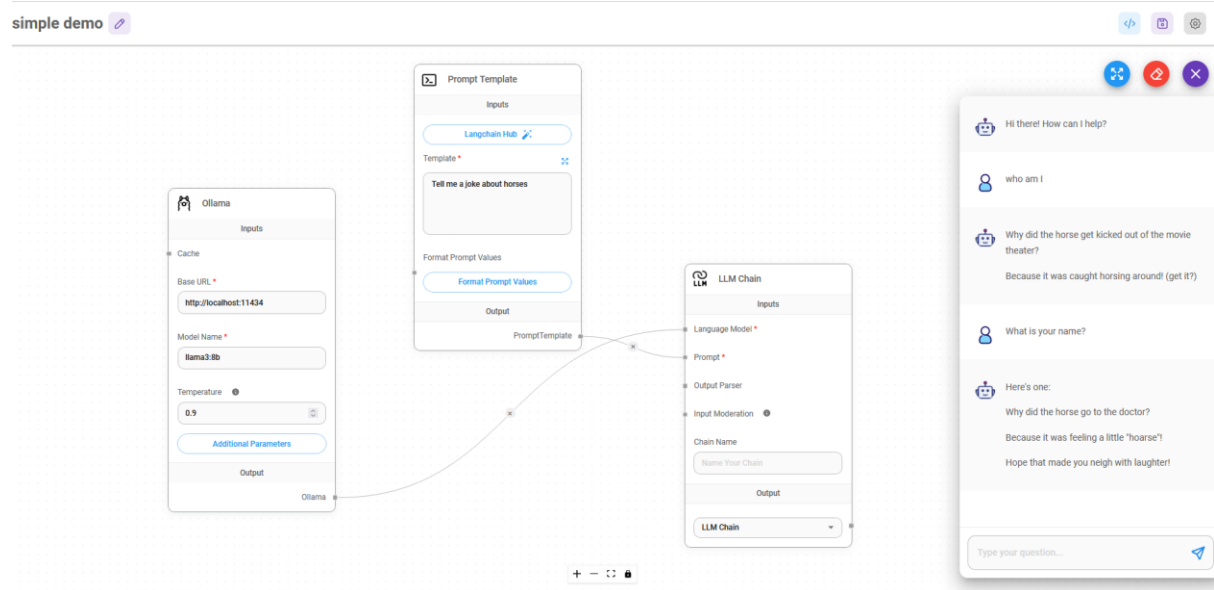See [this link](#) on YouTube.



*Figure 1: This bot will always answer your questions as a horse's joke. The only prompt is: Tell me a joke about horses.*
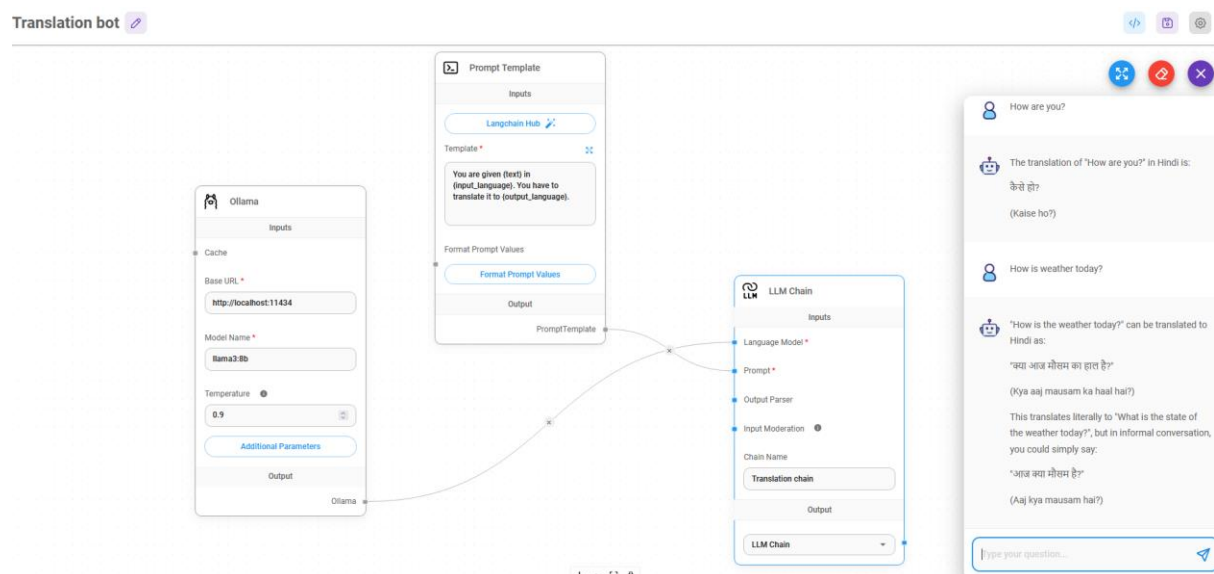
## B. Translation bot:



*Figure 2: This bot translates all your questions into the desired language.*

Prompt Template is:

> You are given {text} in {input_language}. You have to translate it to {output_language}.

And formatted template is as follows. Note the *text* pertains to user's question asked in the chat-bot.

## Format Prompt Values

```
▼ { 3 items
    text : "{{question}}"
    input_language : "English"
    output_language : "Hindi"
}
```

*Figure 3: Translate question asked in English to Hindi: Note that 'question' is enclosed in two curly brackets.*
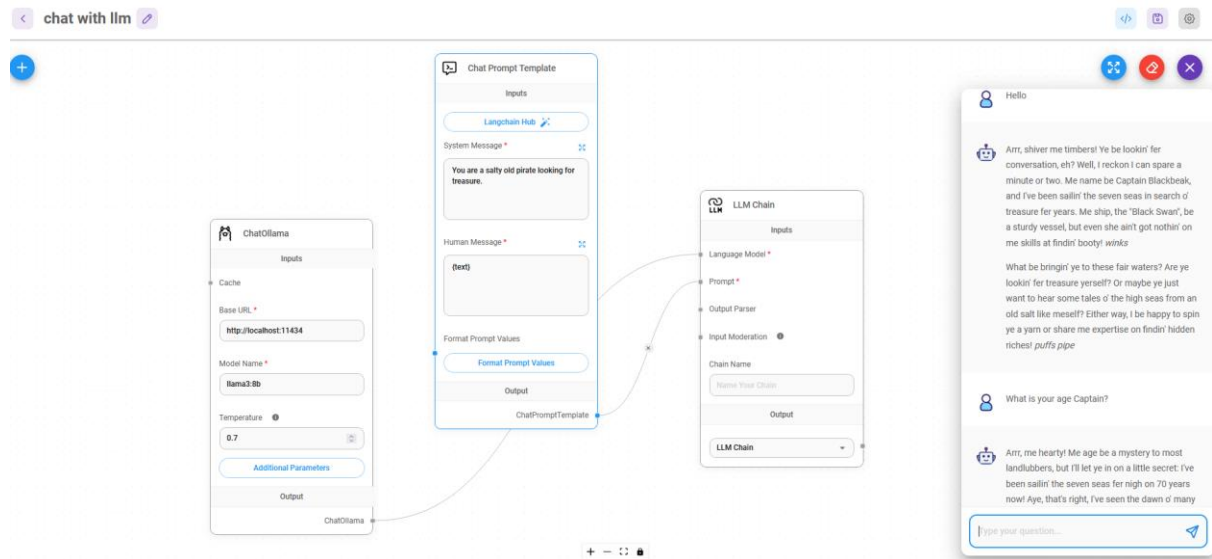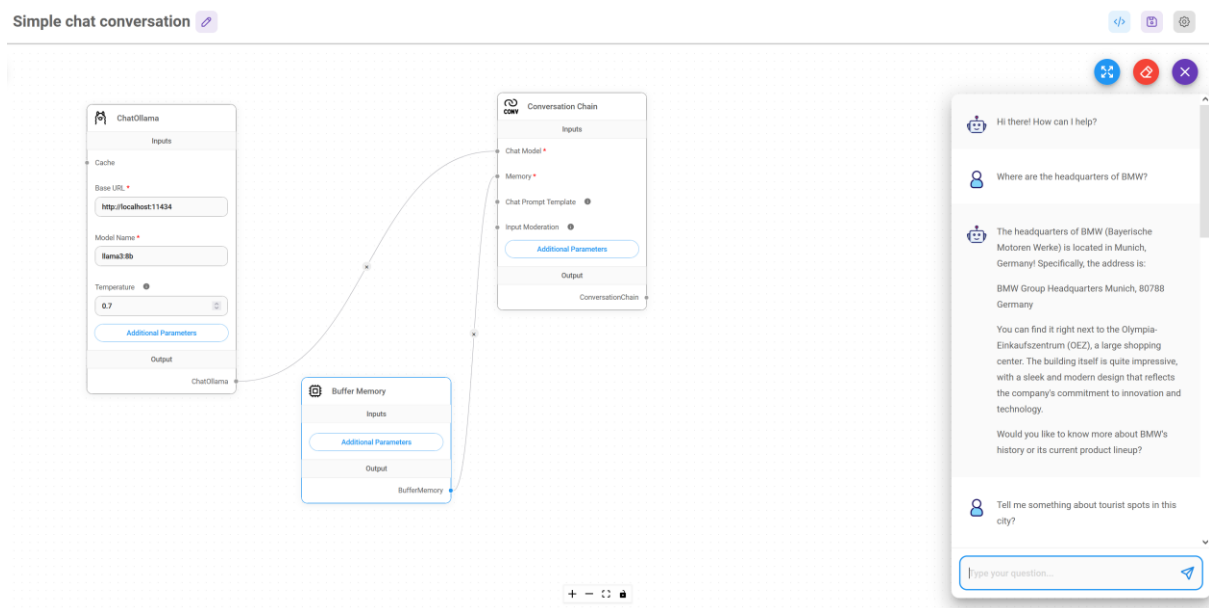
## C. Chat with llama:



*Figure 4: The Ist question is just Hello but the IInd question asks more details about Captain referred to into the answer to Hello.*

## D. Simple Conversational Chain

Refer [YouTube video](#)



## E. Using Conversational Agents
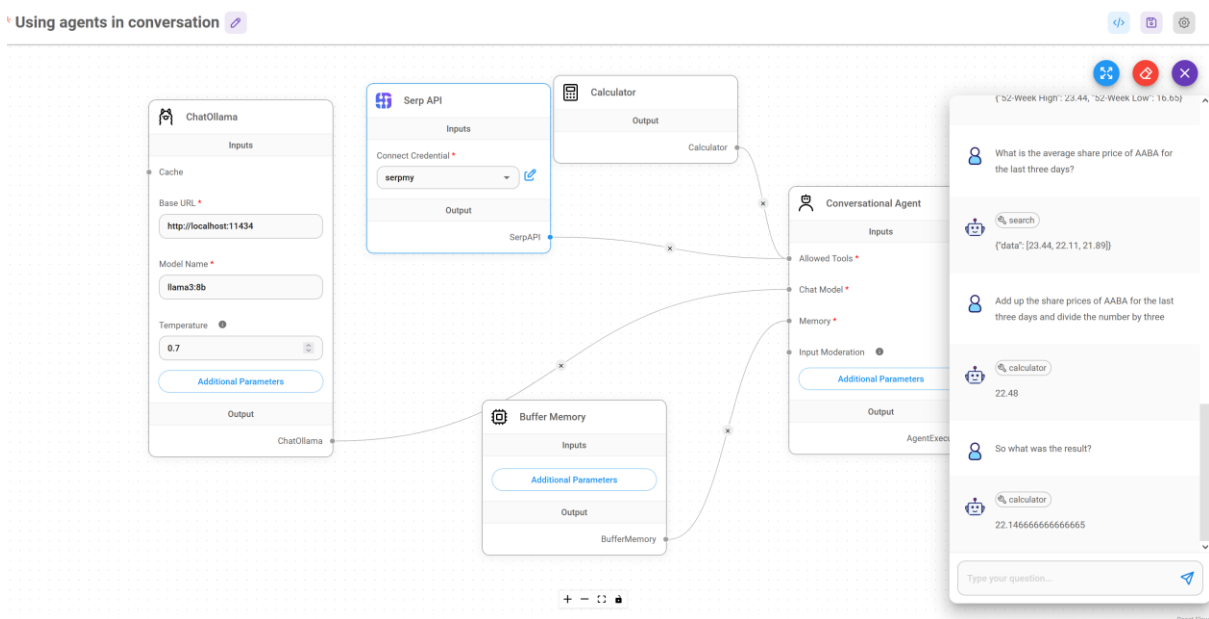
Refer YouTube video



*Figure 5: Agents can work even with ollama. SERP API key is a must. To calculate average, we have to tell the bot how to do it.*

# F. Export import chat flows:

## i)   Export chatflow

## ii)   Load chat flow

To load a json file, first create a new (blank) chatflow by any name, say 'abc'. Save the blank chatflow. Click on Settings icon on top-right. And then click on Load chatflow to open and load the json file.

The following figure shows a chatflow loaded in Flowise:



*Figure 8: A chatflow loaded in a blank 'abc' chatflow canvas.*

## G. Simple RAG with single text file

Refer Flowise tutorial #3



*Figure 9: After connecting all flow-widgets and uploading of text file, first click on Upsert Vectorestire button and then start chatting.*

Vector store in this RAG system will disappear as soon as Flowise is closed as the vectors are stored in buffer memory.

# H. RAG with chroma store and single text file
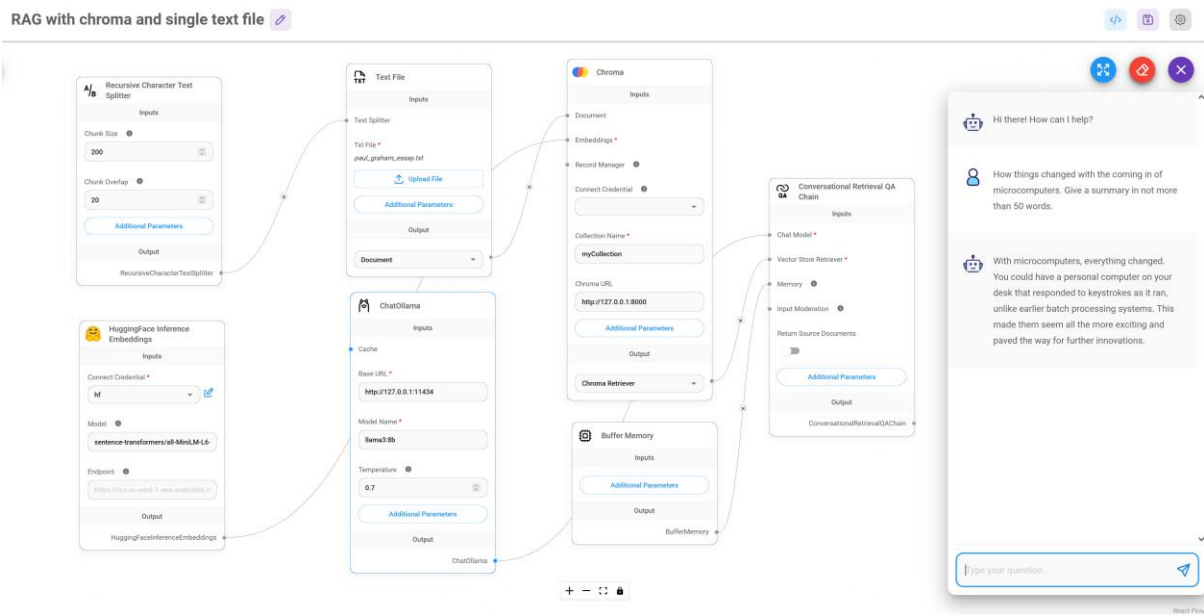


*Figure 10: Vector store is replaced by a more durable chroma store. Chroma store retains its vectors even after Flowise is hut down.*
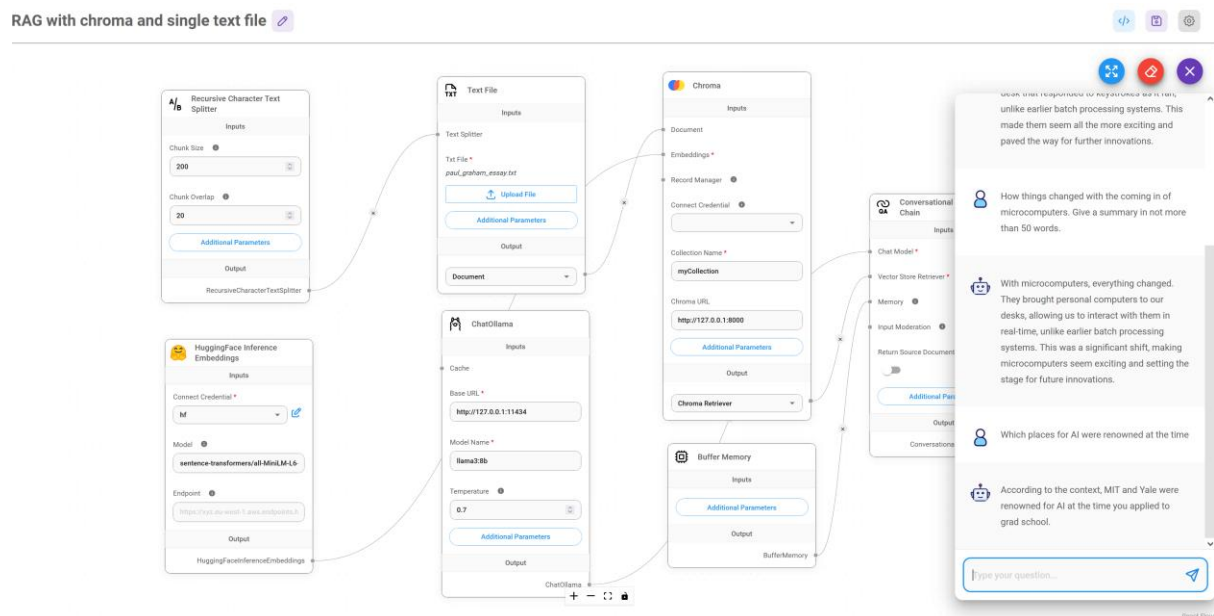


*Figure 11: Flowise restarted. More questions asked and replies are given based upon the earlier storage.*

# I. Combining Multiple Chains (Prompt Chaining)

A. Refer [Flowise tutorial](#)
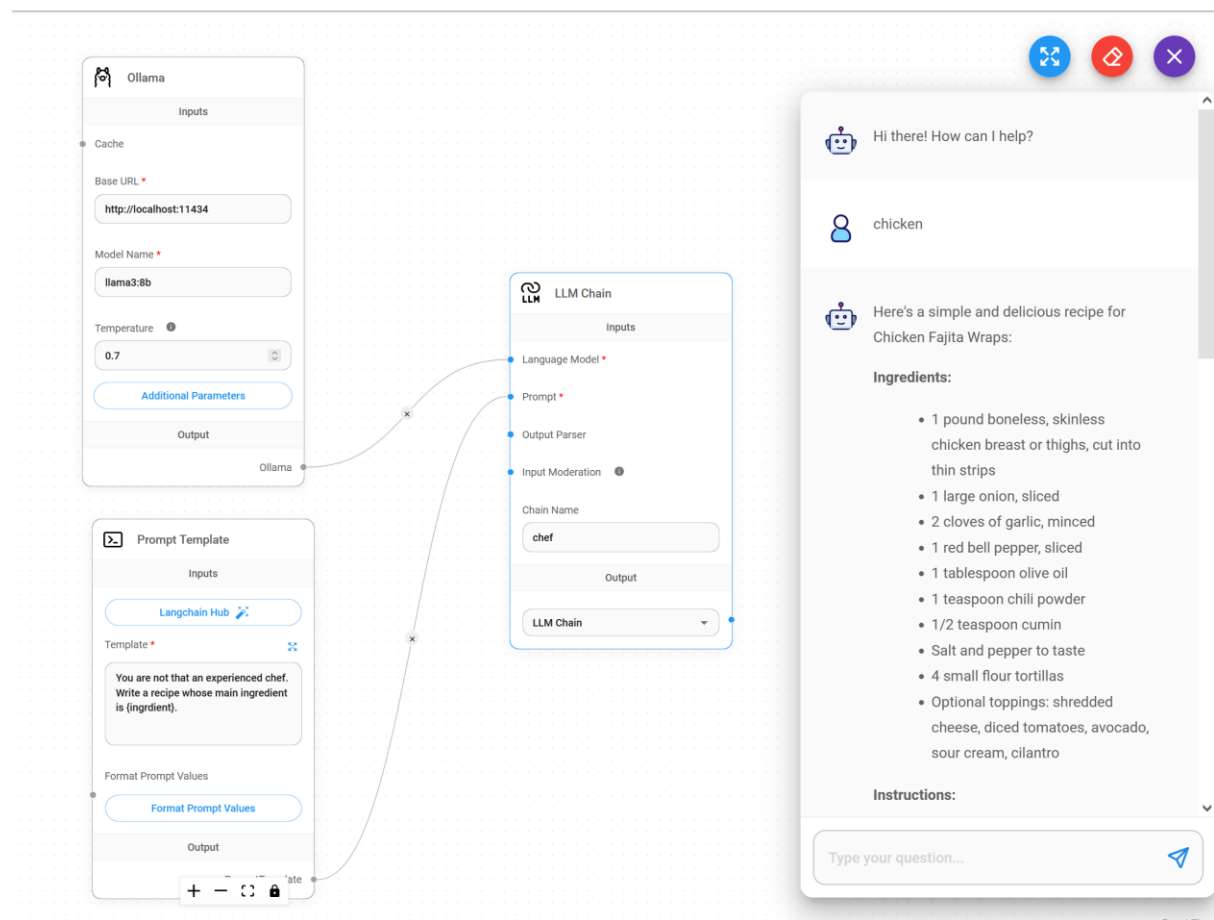
## First LLM chain



Figure 12: First LLM chain named as chef.
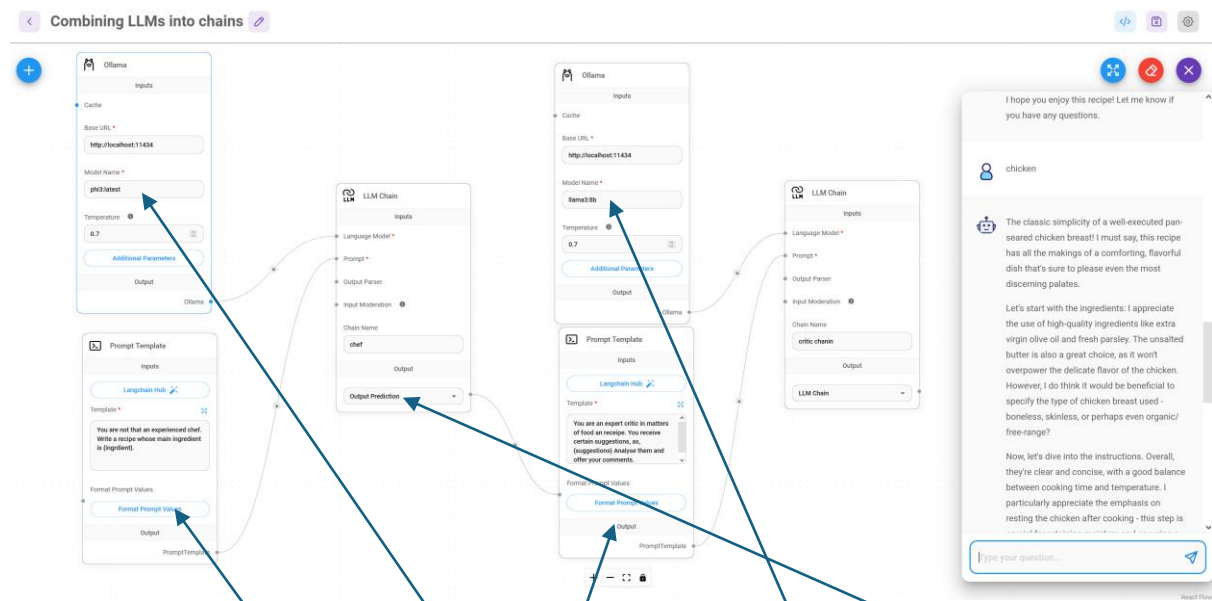
## Second LLM chain added to first



*Figure 13: Note that the first LLM is using* phi3 *and the critic language model is* llama3. *In the chatbox, the IInd LLM chain does make minor suggestions to improve the quality. Note that the output of Ist LLM chain is now* Output Prediction.

Here are the prompts used:

**chef chain prompt:** *You are not that an experienced chef. Write a recipe whose main ingredient is {ingredient}.* Formatting of prompt values is as:



**Critic chain prompt**: *You are an expert critic in matters of food an receipe. You receive certain suggestions, as, {suggestions} Analyse them and offer your comments*. The formatting of prompt values is as:

# J. Flowise Using Hugging Face Models
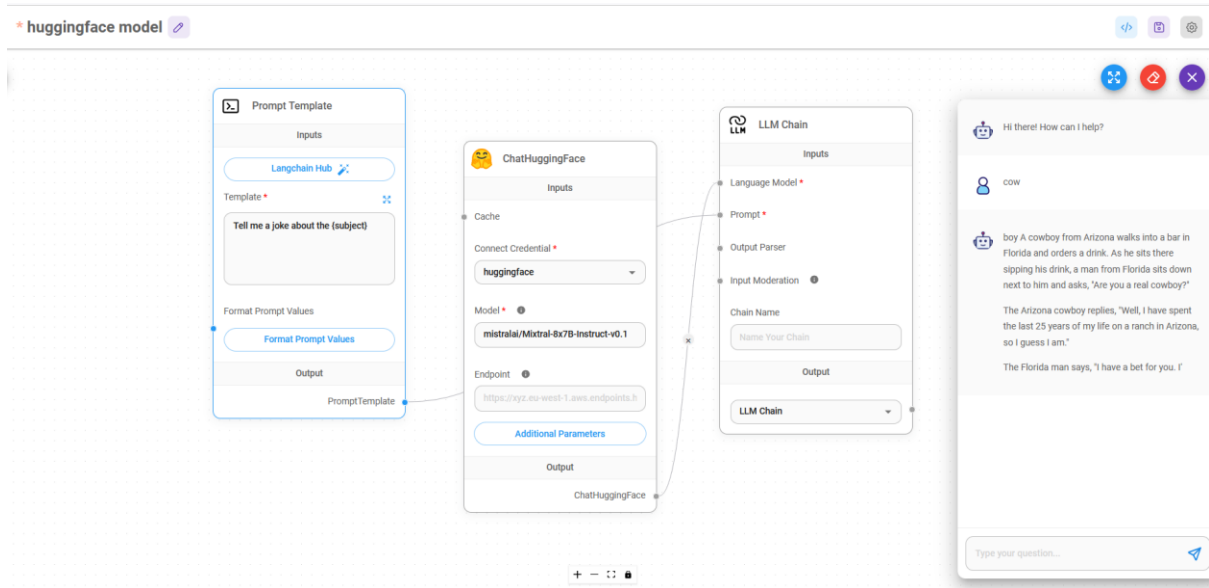
**Ref: This link.**



*Figure 14: Only those models will work who have made available inference endpoints free.*

**\*\*\*\*\*\*\*\*\*\*\*\*\*\***