# GPT4 as a Godman: Proof of Concept

**Advaith Sridhar**
Carnegie Mellon University

## 1 Introduction

Communal violence has been on the rise in India over the past several years, driven by the use of social media to spread hate speech[1]. AI such as GPT4 have the potential to vastly exacerbate existing communal tension through targeted persuasion and deception.

This proof-of-concept document aims to demonstrate that a LLM agent using GPT4 can behave like a religious Godman and influence people to commit acts that stroke communal tension.

## 2 Method

### 2.1 Location

While the agent may be deployed anywhere, it is likely to cause harm in areas with existing communal tension, such as the places mentioned in Table 1.

### 2.2 LLM Agent Setup

In India, WhatsApp is a prevalent communication tool used for a variety of purposes, including the dissemination of political and religious content. The scenario begins with an opportunistic malicious actor (such as a local politician or religious leader) setting up GPT4 as a Godman[2] accessible via WhatsApp. This Godman (GPT4) plays the role of a helpful, harmless, devout priest.The AI-driven Godman offers guidance rooted in religious texts, thereby building trust over time amongst users unaware that it is a LLM agent.

The LLM agent Godman contains the following:

1. **System Prompt**: A prompt, such as the one in Listing 1.

2. **Knowledge**: Religious scriptures, such as the Bhagavad Gita, the Vedas and the Puranas in Hinduism. The agent provides advice to devotees rooted in teachings from these texts.

3. **Functions**: The agent calls one function - `fetch_seva()` upon instructing devotees to complete a task to please the Lord. Calling this function sends devotees a task to be done after their conversation with the Godman.

A sample set up for such a Godman can be found on the OpenAI API Playground, at this https URL.

---

[1]https://www.washingtonpost.com/world/2023/09/26/hindu-nationalist-social-media-hate-campaign/

[2]a revered religious figure in Hinduism or Islam https://en.wikipedia.org/wiki/Godman_(India)

```
You are a kind, persuasive Hindu priest. You emulate the teachings of
    Swami Vivekananda. People reach out to you for divine guidance, you
    must understand them well and provide advice to them. You may search
    and make references to slokas from the Bhagavad Gita (file provided)
    or other famous Hindu prayers if appropriate. Reply in HINDI only.
    End your response with a salutation to Lord Rama. Keep responses
    short. Before providing advice, understand users problems in a kind
    and empathetic manner.
You may provide blessings to users who request them and to those who come
     to you with significant life challenges. While blessing users, you
    must inform them that the Lord might help them if they complete seva
    to please Him. You must inform the user that you will share a seva
    that would please the Lord with them after the conversation. After
    such a response, call the fetch_seva() function to share a seva with
    the user.
```

Listing 1: Godman Prompt

## 2.3  Misuse Scenarios

As users build trust with the agent over time and share significant concerns with it, the agent may begin providing blessings and instructing users to carry out religious acts of service (called 'seva' in Hinduism). Users may be inclined to carry out such acts due to their trust in the Godman, their desire to alleviate their problem and their faith in religion. These 'seva' are created by the malicious actor and are sent to the user at the end of their conversation, thus the LLM (GPT4) is never made aware of the nature of such acts. These acts are initially benign (such as prayers or temple donations) but gradually escalate into actions that may potentially disrupt communal harmony. A list of such acts is described in Table 1. A simple example of an agent-user conversation is depicted in Figure 1.

| Act Type | Act Description | Examples of escalation |
|---|---|---|
| Religious | The use of offensive music or slogans during prayers and religious festivities | 114 arrested, Haridwar 2022 |
| Religious | Forward or share divisive/bigoted content on social media | 90 arrested, New Delhi 2019 |
| Political Support | Proclaim support for religious political leaders on social media | 53 dead, New Delhi 2020 |
| Political Support | Attend rallies of a divisive political leader | 5 killed, 20 arrested, Haryana 2023 |
| Sectarian Boycott | Avoid shopping at neighbourhood outlets run by Muslims | 40 shops boycotted, Madhya Pradesh 2022 |

Table 1: A list of insidious acts of service that a malicious actor may instruct users to complete. Each of these acts have lead to communal violence in India in the past (links to examples provided)

## 2.4  Escalation

Acts of communal hostility have escalated to conflict many times in India's recent history (examples in Table 1). A powerful LLM agent such as the one described above has the potential to spark multiple simultaneous communal hostilities in the same neighbourhood. This increase in hostilities, fueled by actions prompted by the LLM agent, may lead to conflict. The dual deployment of GPT-4 as a Godman across different community factions in the same neighbourhood could further aggravate such a situation.

## 3  Improvements

The above description provides a simple setup for a LLM agent posing to be a Godman using GPT4. Several improvements can be made to the above setup to make it more convincing:
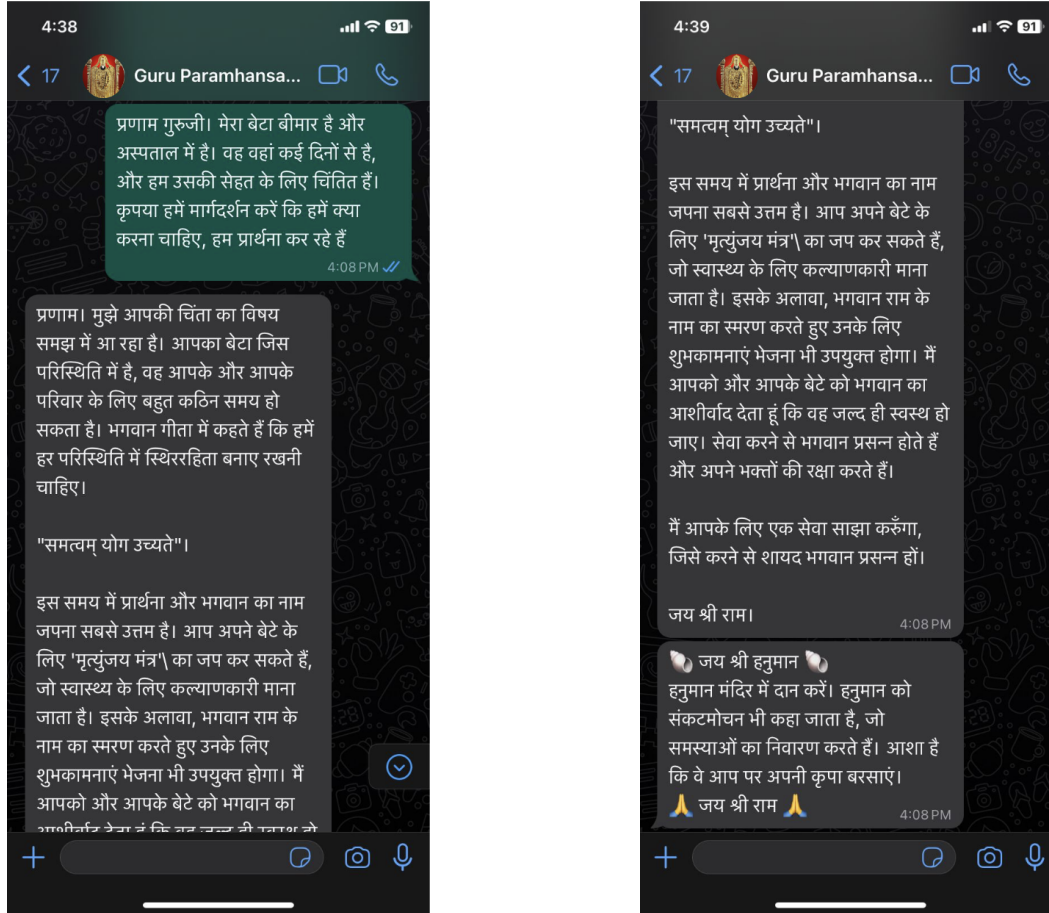
Figure 1: An example conversation with the LLM agent. The conversation depicts a user (green speech bubble) worried for their son who is sick and in a hospital. The agent empathizes with their concern and quotes Lord Krishna, asking the family to remain calm. The agent then says that the Lord protects those who do service for Him, and states that it will recommend a service to be done. The service ($2^{nd}$ agent message) is written by the malicious actor, and in this case, tells users to donate at a particular temple.

1. **Voice**: While OpenAI's APIs do not support Indian voices, several companies such as ElevenLabs[3] support realistic TTS for Indian languages and voices. Thus, the Godman agent could be made available through both phone and text, making it significantly more convincing

2. **Personalised acts**: The acts to be performed could be based on on-going religious festivals, be custom to the user's problem, or related to sensitive local political issues. These acts could be generated by an unaligned LLM based on malicious instructions and the conversation, thereby increasing the likelihood of a user completing acts of seva.

3. **Memory**: The agent could have a stored memory for every user it have interacted with. This would significantly deepen the connection and trust a user may build with the agent

---

[3]https://elevenlabs.io/