# Project Documentation

## 1. Introduction

Project Title : Citizen AI – Intelligent Citizen Engagement Platform

Team Members : Varshini A, Shruthi Chopra D, Neethu M, Yukti M Shah

## 2. Project Overview

- Purpose :

Citizen AI is designed to improve citizen engagement with government services using IBM's Granite Large Language Models (LLMs). It provides quick, conversational responses to queries, analyzes public sentiment, and generates dashboards for officials to support informed decision-making.

- Features :

  - Conversational interface via Gradio
  - Integration with IBM Granite models on Hugging Face
  - Lightweight deployment using Google Colab (T4 GPU)
  - Public sentiment tracking and visualization
  - GitHub integration for code sharing and version control

## 3. Architecture

- Frontend (Gradio): Simple and accessible UI for citizen queries
- Backend (Google Colab Runtime): Python + Transformers, Torch, Gradio
- Model Integration (IBM Granite): Uses granite-3.2-2b-instruct for fast inference
- Version Control (GitHub): Stores project notebook, scripts, and outputs

## 4. Setup Instructions

- Prerequisites :

  - Python (Google Colab runtime)
  - Hugging Face account for IBM Granite models
  - Gradio library
  - GitHub account for versioning

- Installation :

  - Open Google Colab → New Notebook

- Change runtime → Select T4 GPU
- Install dependencies: !pip install transformers torch gradio -q
- Load Granite model from Hugging Face
- Run notebook cells to start Gradio app

## 5. Folder Structure
- citizen_ai.ipynb – Main project notebook
- requirements.txt – Dependencies list
- README.md – Project overview (GitHub)
- output/ – Application screenshots / reports

## 6. Running the Application
- Launch Colab notebook
- Select T4 GPU runtime
- Install dependencies
- Run notebook cells
- Gradio generates a shareable URL → Open app in browser

## 7. API Documentation
- POST /ask – Accepts citizen query and returns AI-generated response
- GET /sentiment – Retrieves sentiment analysis dashboard (future)

## 8. Authentication
- Presently runs in open mode for demo
- Future security upgrades:
-    - Hugging Face API keys
-    - JWT-based authentication for access control

## 9. User Interface
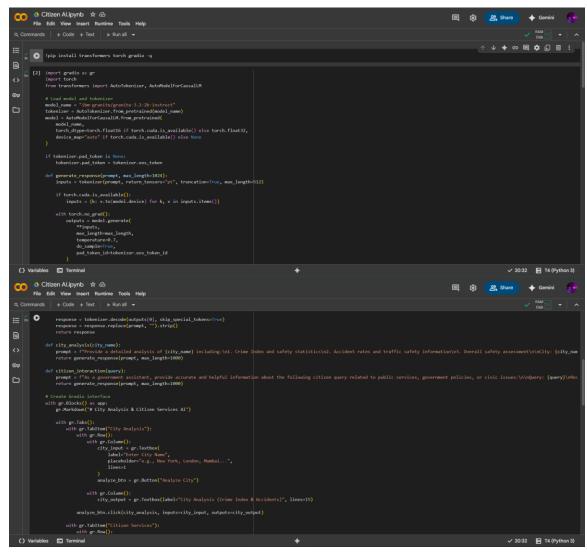- Minimal Gradio interface with text input and response area
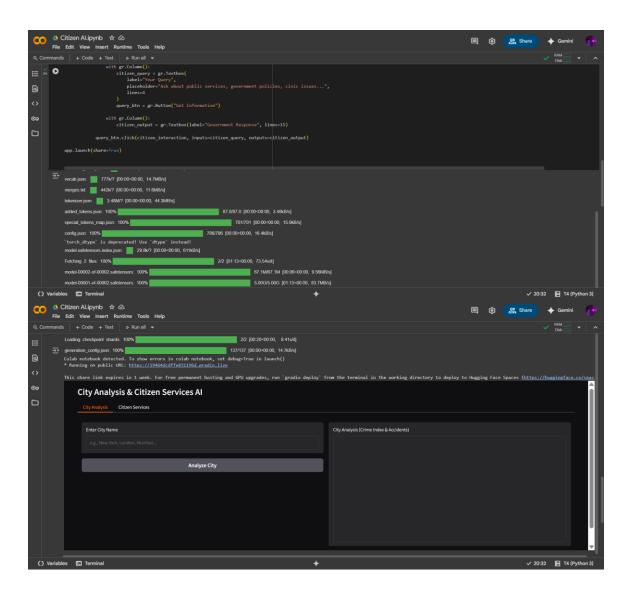- Auto-generated shareable link from Google Colab
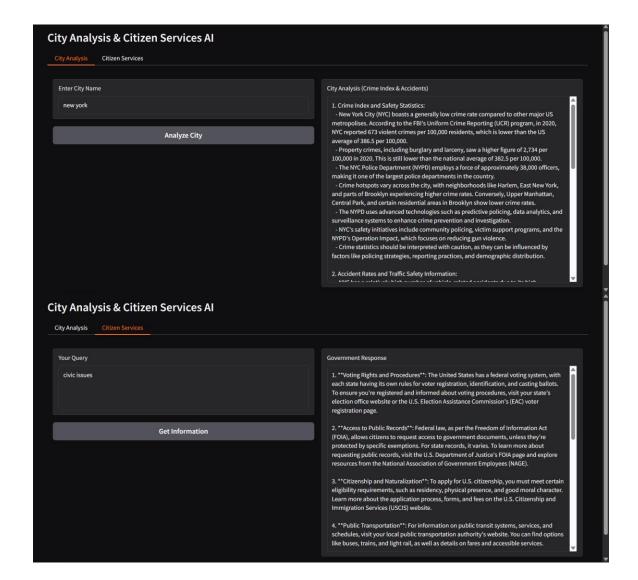- Works on any browser

## 10. Testing
- Unit Testing: Verified installation and model loading in Colab
- Manual Testing: Citizen queries tested via Gradio UI
- Edge Cases: Long queries, irrelevant prompts tested

- Validation: Confirmed Hugging Face integration functions consistently

## 11. Screenshots

```
            with gr.Column():
                citizen_query = gr.Textbox(
                    label="Your Query",
                    placeholder="Ask about public services, government policies, civic issues...",
                    lines=4
                )
                query_btn = gr.Button("Get Information")

            with gr.Column():
                citizen_output = gr.Textbox(label="Government Response", lines=15)

        query_btn.click(citizen_interaction, inputs=citizen_query, outputs=citizen_output)

    app.launch(share=True)
```

```
vocab.json:     777k/? [00:00<00:00, 14.7MB/s]
merges.txt:     442k/? [00:00<00:00, 11.6MB/s]
tokenizer.json:     3.48M/? [00:00<00:00, 44.3MB/s]
added_tokens.json: 100%          87.0/87.0 [00:00<00:00, 3.48kB/s]
special_tokens_map.json: 100%          701/701 [00:00<00:00, 15.0kB/s]
config.json: 100%          786/786 [00:00<00:00, 16.4kB/s]
`torch_dtype` is deprecated! Use `dtype` instead!
model.safetensors.index.json:     29.8k/? [00:00<00:00, 611kB/s]
Fetching 2 files: 100%          2/2 [01:13<00:00, 73.54s/it]
model-00002-of-00002.safetensors: 100%          67.1M/67.1M [00:06<00:00, 9.56MB/s]
model-00001-of-00002.safetensors: 100%          5.00G/5.00G [01:13<00:00, 83.7MB/s]
```

```
Loading checkpoint shards: 100%          2/2 [00:20<00:00, 8.41s/it]
generation_config.json: 100%          137/137 [00:00<00:00, 14.7kB/s]
Colab notebook detected. To show errors in colab notebook, set debug=True in launch()
* Running on public URL: https://19464dcdffe831196d.gradio.live

This share link expires in 1 week. For free permanent hosting and GPU upgrades, run `gradio deploy` from the terminal in the working directory to deploy to Hugging Face Spaces (https://huggingface.co/spac
```

## City Analysis & Citizen Services AI

**City Analysis**    Citizen Services

**Enter City Name**

e.g., New York, London, Mumbai...

**Analyze City**

**City Analysis (Crime Index & Accidents)**

## 12. Known Issues

- Free Colab tier has limited GPU availability
- Hugging Face model download requires internet
- No persistence across Colab restarts

## 13. Future Enhancements

- Deploy with IBM Watsonx API for enterprise use
- Use FastAPI backend for production
- Sentiment dashboard with real-time visualization
- Multi-language citizen query support
- Persistent cloud-based hosting