

```
!pip install transformers torch gradio -q
```

```
# run this project file in google collab by changing run type to T4 GPU
```

```
!pip install transformers torch gradio -q
```

```
import gradio as gr
```

```
import torch
```

```
from transformers import AutoTokenizer, AutoModelForCausalLM
```

```
# Load model and tokenizer
```

```
model_name = "ibm-granite/granite-3.2-2b-instruct"
```

```
tokenizer = AutoTokenizer.from_pretrained(model_name)
```

```
model = AutoModelForCausalLM.from_pretrained(  
    model_name,  
    torch_dtype=torch.float16 if torch.cuda.is_available() else torch.float32,  
    device_map="auto" if torch.cuda.is_available() else None  
)
```

```
if tokenizer.pad_token is None:
```

```
    tokenizer.pad_token = tokenizer.eos_token
```

```
def generate_response(prompt, max_length=1024):
```

```
    inputs = tokenizer(prompt, return_tensors="pt", truncation=True, max_length=512)
```

```
    if torch.cuda.is_available():
```

```
        inputs = {k: v.to(model.device) for k, v in inputs.items()}
```

```
    with torch.no_grad():
```

```
        outputs = model.generate(  
            **inputs,  
            max_length=max_length,  
            temperature=0.7,  
            do_sample=True,  
            pad_token_id=tokenizer.eos_token_id  
        )
```

```
    response = tokenizer.decode(outputs[0], skip_special_tokens=True)
```

```
    response = response.replace(prompt, "").strip()
```

```
    return response
```

```
def city_analysis(city_name):
```

```
    prompt = f"Provide a detailed analysis of {city_name} including:\n1. Crime Index and safety
```

```
    return generate_response(prompt, max_length=1000)
```

```
def citizen_interaction(query):
```

```
    prompt = f"As a government assistant, provide accurate and helpful information about the fo
```

```
    return generate_response(prompt, max_length=1000)
```

```
# Create Gradio interface
```

```
with gr.Blocks() as app:
```

```
    gr.Markdown("# City Analysis & Citizen Services AI")
```

```
    with gr.Tabs():
```

```
        with gr.TabItem("City Analysis"):
```

```
            with gr.Row():
```

```
                with gr.Column():
```

```
                    city_input = gr.Textbox(  
                        label="Enter City Name",  
                        placeholder="e.g., New York, London, Mumbai...",  
                        lines=1
```

```

        )
        analyze_btn = gr.Button("Analyze City")

    with gr.Column():
        city_output = gr.Textbox(label="City Analysis (Crime Index & Accidents)", 1

    analyze_btn.click(city_analysis, inputs=city_input, outputs=city_output)

with gr.TabItem("Citizen Services"):
    with gr.Row():
        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)

```

Fetching 2 files: 100%

2/2 [01:14<00:00, 74.20s/it]

model-00001-of-00002.safetensors: 100%

5.00G/5.00G [01:13<00:00, 66.3MB/s]

Loading checkpoint shards: 100%

2/2 [00:25<00:00, 10.40s/it]

generation\_config.json: 100%

137/137 [00:00<00:00, 12.0kB/s]

Colab notebook detected. To show errors in colab notebook, set debug=True in launch()

\* Running on public URL: <https://dd0bb62acecb74bf02.gradio.live>

This share link expires in 1 week. For free permanent hosting and GPU upgrades, run `gradio dep



**No interface is running right now**