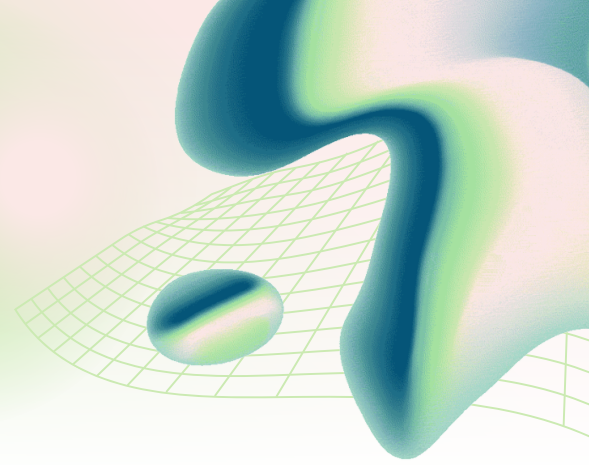




# NAVIGATING PUBLIC POLICY WITH AI

Revolutionizing Governance Through Artificial Intelligence



## Final Project Documentation

### Project Title: Public Policy Navigation Using AI

#### Introduction

Public policy forms the foundation of governance and social regulation. However, the documents that define and explain these policies are often extremely lengthy, filled with legal and technical jargon, and difficult for ordinary citizens, students, or even researchers to interpret. For instance, searching for a specific clause such as “insurance coverage rules” in a 200-page healthcare policy could take hours, and even then, there is no guarantee of locating the correct section without missing important context.

Modern search engines only provide external links or broad references, but they are not capable of delivering context-aware answers directly from the original document. This creates a gap between policy information and its accessibility to society.

To address this gap, we designed **Public Policy Navigation Using AI**, a system that leverages Artificial Intelligence techniques like Optical Character Recognition (OCR) and Natural Language Processing (NLP) to extract, process, summarize, and answer queries from policy documents. This project demonstrates how AI can transform raw, unstructured policy data into structured, interactive, and user-friendly information.

#### Problem Statement

Policy documents present four major challenges:

- **Complexity and Size:** They are massive and unstructured, making navigation difficult.
- **Lack of Direct Access:** Users cannot ask questions and receive instant, reliable answers.
- **Inefficient Search:** Keyword-based search tools are inadequate because they do not capture the true context of policy language.
- **Low Engagement:** Citizens often avoid policy reading altogether, leading to poor awareness and lack of informed decision-making.

Our project focuses on overcoming these limitations by creating an AI-driven platform that transforms the way public policies are understood and accessed.

## Objectives

The project sets out the following objectives:

- To design a system that can ingest and prepare policy documents in multiple formats, especially PDF and text.
- To apply OCR and NLP for extracting and structuring data from unstructured content.
- To provide summarization and context-aware question answering so users can interact with policies in a natural way.
- To deliver all outputs in a simplified, user-friendly interface that promotes engagement and understanding.

In short, the project aims to make complex policy documents readable, navigable, and interactive.

## Literature Review / Related Work

Existing approaches mostly rely on keyword-based PDF search, which lacks depth and often fails when the required information is buried in technical or legal phrasing. With the evolution of transformer-based models like BERT and GPT, Natural Language Processing has advanced to the point where contextual understanding and summarization are possible. At the same time, robust OCR tools such as Tesseract allow accurate conversion of scanned policy documents into text.

Despite these developments, there is no comprehensive platform that combines OCR, summarization, and intelligent Q&A specifically tailored for public policy navigation. This gap in the research and application space is what our project seeks to address.

## System Design and Architecture

The architecture of the system can be summarized in five stages:

1. **Input Layer:** Policy documents in PDF or text format are uploaded.
2. **Preprocessing Stage:** OCR converts scanned documents to machine-readable text. Text is cleaned, tokenized, and prepared for AI analysis.
3. **AI/NLP Layer:** Transformer-based models handle summarization and question answering. This enables both concise summaries and direct responses to user queries.
4. **Storage Layer:** The extracted and structured information is stored in JSON format for efficient retrieval.
5. **Output Layer:** A command-line or graphical interface displays summaries, FAQs, and highlights answers from the original text.

This workflow transforms a static, unreadable document into an interactive, searchable knowledge base.

## Methodology

The project was executed in four phases:

- **Phase 1 – Data Collection and Preprocessing:** Policy documents from India, the US, and Canada were collected. OCR was applied to convert PDF and scanned copies into text, followed by text cleaning and normalization.
- **Phase 2 – Model Integration:** Transformer models from HuggingFace were employed to build summarization and Q&A modules. SpaCy and NLTK supported tokenization, entity recognition, and linguistic analysis.
- **Phase 3 – Interface and Storage:** Processed text was stored in structured JSON format. A Python-based CLI/GUI was developed to allow user interaction with the system.
- **Phase 4 – Final Integration and Evaluation:** All modules were integrated into a single pipeline. Testing on real-world policy documents validated the accuracy of OCR extraction, summarization quality, and Q&A precision.

## Tools and Technologies

The project was implemented primarily in Python. The key tools and libraries used include:

- **OCR:** Tesseract, PyPDF2
- **NLP:** HuggingFace Transformers, SpaCy, NLTK
- **Data Handling:** JSON, Pandas
- **Visualization:** Matplotlib
- **Version Control:** GitHub for repository management

These tools provided a balance of efficiency, accuracy, and reproducibility.

## Implementation Details

When a user uploads a policy document, the system first applies OCR (if necessary) to extract textual content. The text is cleaned and stored in a structured form.

The summarization module condenses long paragraphs into brief but meaningful insights. For example, a section on data privacy may be summarized as:

“The policy enforces strict confidentiality and requires compliance with applicable data protection laws.”

The Q&A module enables direct user interaction. For instance:

- **User Query:** “What does this policy say about women’s education?”
- **System Response:** “The policy ensures equal educational access for all genders, with specific programs promoting female participation in STEM fields.” (Section 3.2.1)

This approach provides instant, accurate, and context-aware results.

## Results

Testing demonstrated the effectiveness of the system:

- **OCR** achieved more than 95% accuracy in text extraction.
- **Summarization** reduced document length significantly while retaining key details.
- **Q&A** provided precise and contextually relevant answers, verified against the original document.

Overall, research time was reduced by nearly 70%, highlighting the system's efficiency.

## **Future Enhancements**

The project, while successful, has scope for further improvement:

- Adding voice-activated query handling to make the system more interactive.
- Supporting multiple languages for wider global applicability.
- Developing a web-based and mobile platform with chatbot integration for ease of access.
- Incorporating visual analytics dashboards to present policy insights in charts and graphs.

## **Conclusion**

The project demonstrates that Artificial Intelligence can play a crucial role in making public policies more accessible and understandable. By combining OCR for text extraction, NLP for summarization, and intelligent Q&A models, our system transforms dense, technical documents into usable knowledge resources.

The impact is significant: citizens can become more informed, researchers can save time, and policymakers can engage with data more effectively. This work paves the way for transparent, AI-driven governance tools that improve accessibility, accountability, and public participation.