**SnT 2025**
CTBT: SCIENCE AND TECHNOLOGY CONFERENCE

8 SEPTEMBER
ONLINE DAY
9 TO 12 SEPTEMBER
AT HOFBURG PALACE, VIENNA & ONLINE

P4.3-883

# Enhancing Information Access to CTBTO Staff through Generative AI: A Smart Chat Application for Simplified Access to Complex Directives and Protocols

## Marko Bosancic

Comprehensive Nuclear-Test-Ban Treaty Organization (CTBTO)

**CTBTO** PREPARATORY COMMISSION
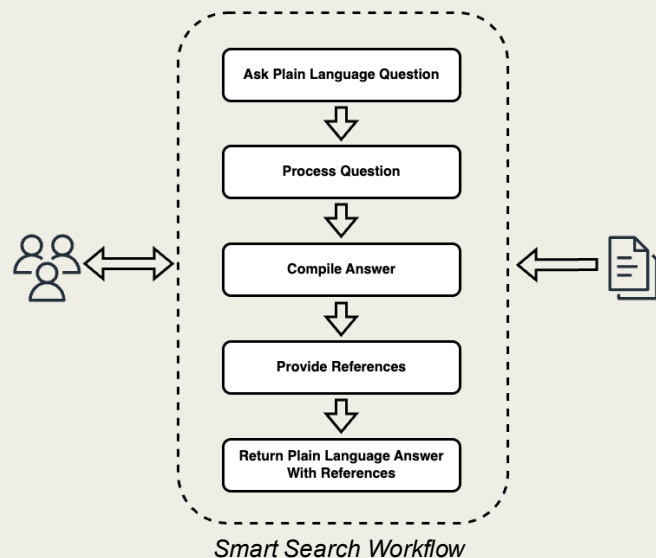
### INTRODUCTION AND MAIN RESULTS

Presenting a Generative AI based application that enhances organization staff access to complex directives and protocols. Providing natural language conversation queries with accurate, reference backed responses and iterative refinement, the system reduces the time and effort of manual search and interpretation of documentation, thereby improving operational efficiency and knowledge accessibility

# SnT 2025
CTBT: SCIENCE AND TECHNOLOGY CONFERENCE

8 SEPTEMBER
ONLINE DAY
9 TO 12 SEPTEMBER
AT HOFBURG PALACE, VIENNA & ONLINE

# Enhancing Information Access to CTBTO Staff through Generative AI: A Smart Chat Application for Simplified Access to Complex Directives and Protocols

Marko Bosancic
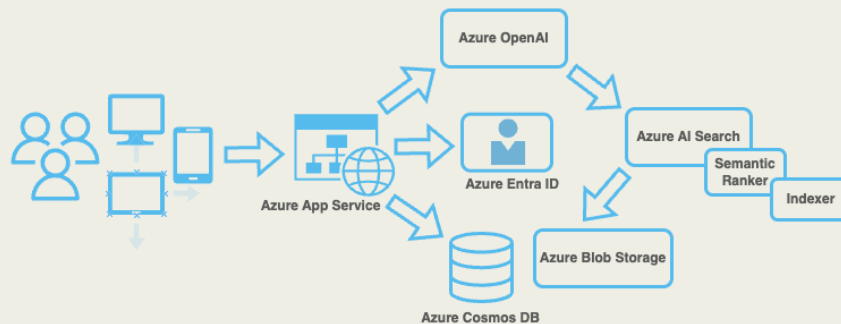Comprehensive Nuclear-Test-Ban Treaty Organization (CTBTO)

## Introduction

Successful market adaptation of large language models (LLMs) and the rapidly expanding capabilities of generative AI has demonstrated how such technologies can be applied to various organizational challenges with growing demand for assistive software that supports staff in navigating information heavy processes and minimizing cognitive burden. CTBTO approach was to utilize applications based on Generative AI aimed at enhancing access to complex directives and protocols by enabling natural language conversation-like queries, delivering accurate and reference-backed responses and allowing iterative refinement of results.



*Smart Search Workflow*

## System Architecture and Implementation



*Retrieval Augmented Generation (RAG) architecture*
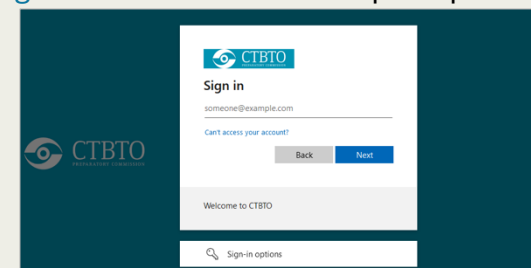
Generative AI chat application was implemented using Microsoft Azure cloud services to provide rapid development of scalable, secure and high-performance solutions for CTBTO staff. The system integrates multiple Azure resources to handle document storage, content search and AI processing orchestrated according to a market proven Retrieval Augmented Generation (RAG) design pattern, together with authentication and conversation history management.

- **Azure AI Search**: Indexing the directives, protocols and other related documentation to provide best ranked search results to the LLM
- **Azure OpenAI**: GPT-4 LLM is deployed for natural language reasoning capabilities to generate responses to user prompts based on provided results
- **Azure Blob Storage**: Central repository for targe-domain directive and protocol documents, accessible by search indexers, to constrain generative AI responses

- **Azure Cosmos DB**: Stores conversation histories, enabling iterative refinement of user queries and continuity across sessions
- **Azure App Service**: Hosts the user-facing custom developed application, providing a responsive UI for all devices and secure interface for staff queries
- **Azure Entra ID**: Ensures secure, role-based organizational 2F Authentication and access to the application for selected or all organization staff, depending on the business domain sensitivity the application covers

## Authentication and Authorization

CTBTO utilizes Azure Entra ID services for authentication and authorization across different systems. Federated identity for seamless Authentication and Role-Based Access Control (RBAC) for managing authorization through Azure App Services container provided endpoints /.auth/login/aad and /.auth/me for principal data.



*Azure Entra ID 2FA*

# Enhancing Information Access to CTBTO Staff through Generative AI: A Smart Chat Application for Simplified Access to Complex Directives and Protocols
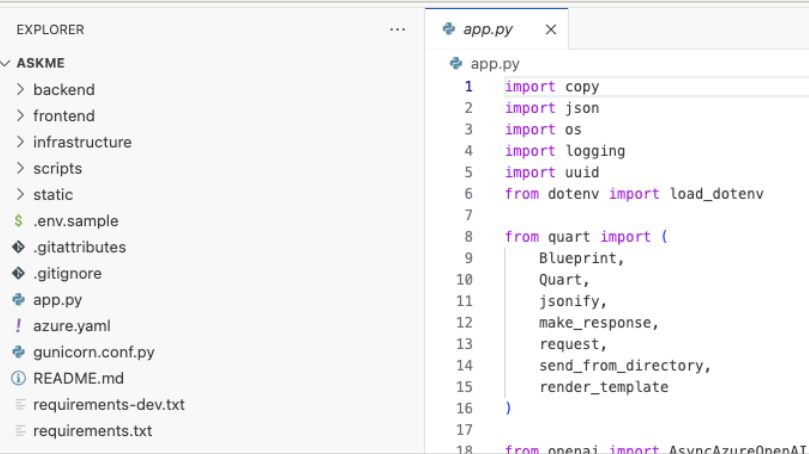
Marko Bosancic
Comprehensive Nuclear-Test-Ban Treaty Organization (CTBTO)
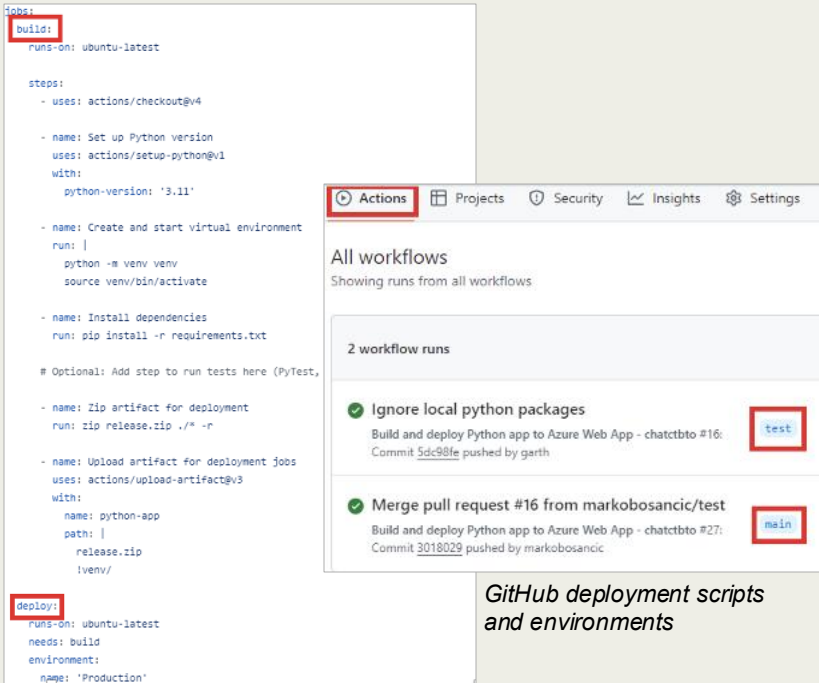
P4.3-883

## Software development and deployment

To enable intuitive and effortless utilization of AI capabilities on targeted domain documentation, custom software was developed using Node.js, TypeScript, React, and FluentUI components on the Frontend and lightweight Quart Python framework with multiple Azure libraries on the Backend.
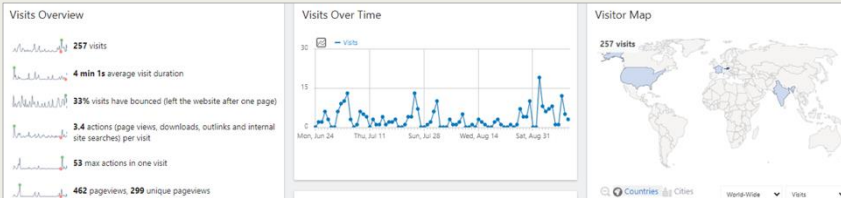
GitHub repository is used for software versioning with GithHub Actions for seamless Azure App Service deployments compose out of main PremiumV3 (P1v3 - 2 vCPU and 8GB of memory) and multiple "Slot" instances for different development environments and stages.



```python
# app.py
import copy
import json
import os
import logging
import uuid
from dotenv import load_dotenv

from quart import (
    Blueprint,
    Quart,
    jsonify,
    make_response,
    request,
    send_from_directory,
    render_template
)

from openai import AsyncAzureOpenAI
```



GitHub deployment scripts and environments

| Language | files | blank | comment | code |
|---|---|---|---|---|
| JSON | 9 | 0 | 0 | 5448 |
| TypeScript | 24 | 258 | 38 | 3119 |
| Python | 16 | 571 | 304 | 3113 |
| CSS | 7 | 119 | 1 | 836 |
| YAML | 8 | 30 | 2 | 190 |
| Markdown | 1 | 66 | 0 | 188 |
| SVG | 2 | 0 | 0 | 33 |
| HTML | 2 | 1 | 0 | 27 |
| Text | 2 | 0 | 0 | 24 |
| SUM: | 71 | 1045 | 345 | 12978 |

*Lightweight codebase*

## Analytics and costs

Implemented Matomo analytics showed early traction with staff members shortly after release with development and post-release Azure costs remaining negligible.



*2-Months Analytics after the initial release*



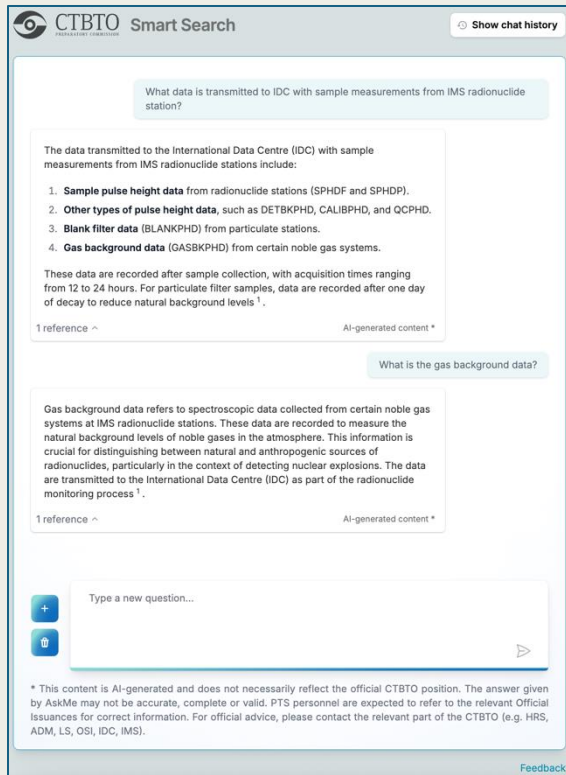*Azure costs for the development and post-release period*

## Results

After an initial investment of time and effort to develop a generic and highly customizable smart chat solution, multiple instances of the application were rapidly deployed to serve various business domains within the organization. Access to each instance is tailored based on the sensitivity of the information, ranging from restricted access for select staff to full access for all personnel. Instances and POC-s for different organizational contexts cover analysis of the International Labor Organization documents for the Legal department, operational Manuals for the On-Site Inspection teams, clarification documents for results presented on the Secure Web Portal and simplifying Administrative Directives for CTBTO Staff.

# Enhancing Information Access to CTBTO Staff through Generative AI:
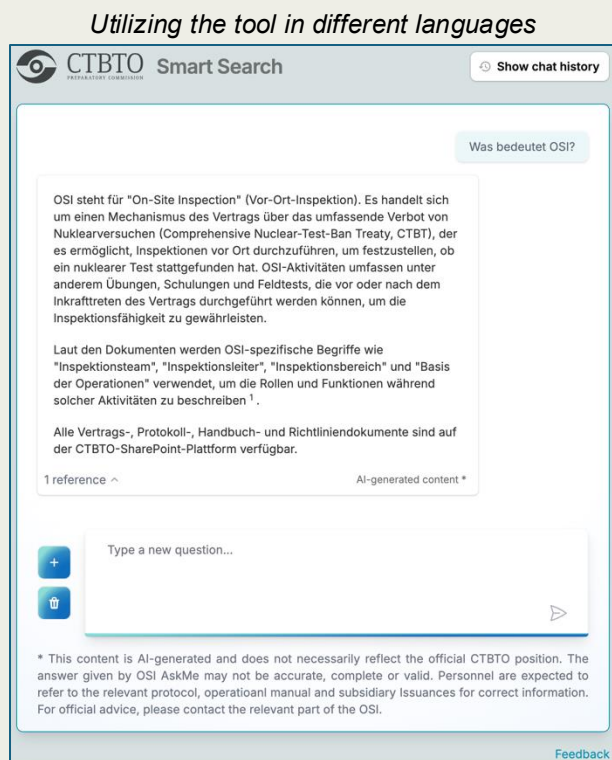## A Smart Chat Application for Simplified Access to Complex Directives and Protocols

P4.3-883

Marko Bosancic
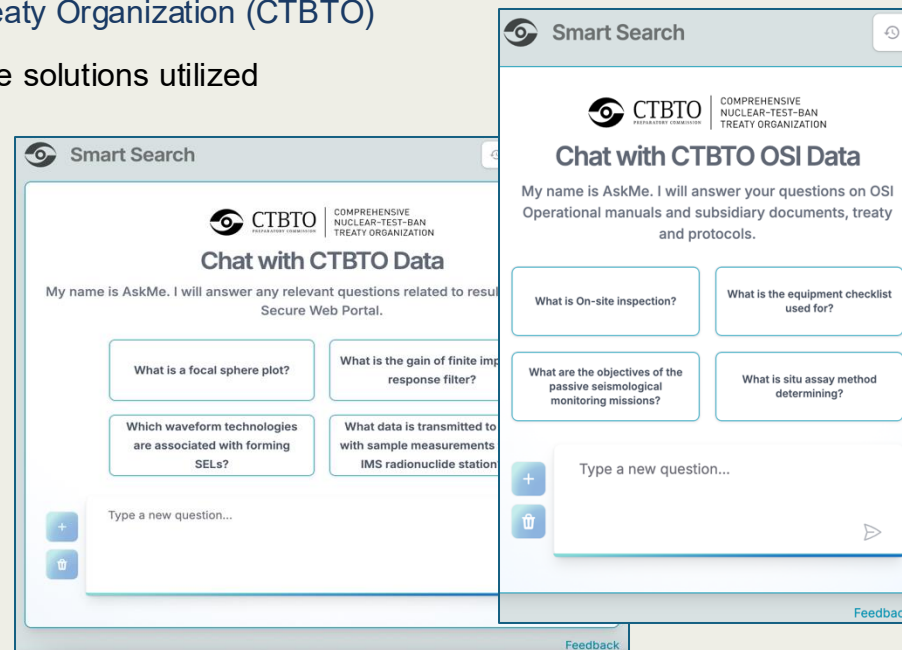Comprehensive Nuclear-Test-Ban Treaty Organization (CTBTO)

The following graphics illustrate the design, functionality and practical application of the solutions utilized in the CTBTO organization.



*Simplifying complex topics*

*Utilizing the tool in different languages*







*Web client Interface on different devices*

## Conclusion

In conclusion, the AI based systems significantly reduces the time and effort associated with manual search and interpretation of complex or domain specific documentation. Removing repetitive workloads, fostering consistency in interpretation, enables staff to focus more effectively on their core responsibilities, benefiting not only those who seek answers but also staff responsible for interpreting and providing relevant information. In doing so, such systems strengthens both operational efficiency and knowledge accessibility across the organization, operated at negligible costs with highly controllable and predictable expenses.

Presented solutions are sustainable for long-term organizational adoption and more importantly, the potential for further expansion and contribution to organizational efficiency is endless.