

# **PROJECT REPORT**

## **ARTIFICIAL INTELLIGENCE AND EXPERT SYSTEM**

**CT-361**



### **GROUP MEMBERS:**

**CT-22056 YUSRA HAQUE**

**CT-22062 FARIA SHAHID**

**CT-22066 HAMDIA NOUMAN**

**CT-22069 NADEEYAH JAMIL**

# AI URDU TUTOR

## TABLE OF CONTENTS

No	Topic	Pg No
1.1	Problem Statement	3
1.2	Project Overview	3
2	Technologies & Libraries used	3
3	Features	4
4	How it Works	5
5	Code	5
6	Conclusion	14

## **1.Introduction:**

### **1.1 PROBLEM STATEMENT:**

In Pakistan and other Urdu-speaking regions, many children face challenges in accessing engaging and interactive learning resources in their native language. Traditional educational methods often lack personalization, interactivity, and accessibility—especially for young learners. Additionally, there is a scarcity of AI-powered educational tools specifically designed for kids in Urdu, which can make language learning and comprehension more difficult for early learners. This gap hinders not only linguistic development but also the adoption of modern learning methods in local contexts.

### **1.2 PROJECT OVERVIEW:**

The "AI Urdu Tutor for Kids" is an intelligent, interactive educational tool designed to help children learn and understand the Urdu language through an engaging conversational interface. The system leverages artificial intelligence to provide real-time answers to children's questions using a curated knowledge base that includes an Urdu dictionary and a collection of age-appropriate stories. Children can ask the AI questions in Urdu, and it will respond with accurate and easy-to-understand answers, making language learning more interactive and enjoyable. The project aims to promote Urdu literacy in a fun and accessible way, especially for early-age learners, by combining modern AI technology with culturally relevant content.

## **2.Technologies & Libraries Used:**

<b>Tool / Library</b>	<b>Purpose</b>
<b>streamlit</b>	To build the user interface
<b>speech_recognition</b>	Convert Urdu audio to text
<b>gtts</b>	Convert Urdu text to speech
<b>fitz (PyMuPDF)</b>	Read and extract text from PDF
<b>SentenceTransformer</b>	Convert sentences to embeddings (numbers)
<b>faiss</b>	Fast vector-based search
<b>langchain</b>	To handle prompts, memory, and chains
<b>ChatGroq + LLaMA 3</b>	The actual LLM generating the response
<b>audio_recorder_streamlit</b>	For recording audio in Streamlit
<b>tempfile, os, BytesIO</b>	File and memory handling

### **3. Features :**

<b>RAG (Retrieval-Augmented Generation)</b>	Combines knowledge from a PDF with AI's intelligence to answer accurately.
<b>Urdu PDF Support</b>	Upload and read Urdu PDF files.
<b>Voice Input</b>	User can speak in Urdu to ask questions.
<b>Voice Output (TTS)</b>	Assistant speaks back the answers.
<b>Memory</b>	Remembers previous questions and answers (context-aware).
<b>Caching</b>	If the same question is asked again, it gives a faster answer using stored results.
<b>Control Buttons</b>	Options to clear memory, toggle memory visibility, etc.

### **4..How It Works:**

1. User uploads an Urdu PDF.
2. Bot processes the PDF and creates a searchable vector store.
3. User asks a question (voice or text).
4. Bot converts voice to Urdu text incase of voice input.
5. Bot retrieves related chunks from the PDF using FAISS.
6. Bot generates an answer using LLaMA 3 (Groq) with context + memory.
7. Bot converts the Urdu answer to audio and shows + speaks it.

## 5.Code:

```
4 import streamlit as st
5 import speech_recognition as sr
6 from gtts import gTTS
7 from io import BytesIO
8 import fitz # PyMuPDF
9 import os
10 import numpy as np
11 import faiss
12 from sentence_transformers import SentenceTransformer
13 from langchain_groq import ChatGroq
14 from langchain.text_splitter import RecursiveCharacterTextSplitter
15 import torch
16 from audio_recorder_streamlit import audio_recorder
17 import tempfile
18 import time
19 from langchain.memory import ConversationBufferMemory
20 from langchain.chains import ConversationChain
21 from langchain.prompts import PromptTemplate
22 from langchain.chains import LLMChain
23
```

### 1. User Interface (UI)

- **streamlit:**  
Builds the interactive web app interface with chat input, buttons, and audio playback.
- **audio\_recorder\_streamlit:**  
Adds a user-friendly audio recording button directly to the Streamlit app.

### 2. Voice Interaction

- **speech\_recognition:**  
Converts spoken Urdu input into text using Google's Speech API.
- **gTTS (Google Text-to-Speech):**  
Converts Urdu text responses to audio for playback.
- **BytesIO:**  
Handles audio input/output entirely in memory (no file writing needed).

- **tempfile:**  
Manages temporary files, especially for handling audio during text-to-speech processing.

### 3. Document Processing

- **fitz (PyMuPDF):**  
Extracts text from uploaded PDF files, including support for right-to-left Urdu.
- **RecursiveCharacterTextSplitter:**  
Splits large documents into smaller, manageable chunks for retrieval.

### 4. Natural Language Processing & Retrieval

- **SentenceTransformer:**  
Converts each text chunk into vector embeddings for similarity search.
- **numpy:**  
Formats embeddings into compatible structures for FAISS.
- **faiss:**  
Performs efficient vector-based retrieval to find the most relevant document chunks for a user's question.

### 5. Language Model Integration (LLM)

- **ChatGroq via LangChain:**  
Connects to Groq's LLaMA3 model to generate answers based on context and conversation history.
- **PromptTemplate:**  
Defines the structure of inputs sent to the language model (history + context + question).
- **LLMChain:**  
Links the prompt, memory, and model together to execute the query pipeline.

- **ConversationChain:**  
Manages the overall conversational flow with memory and model response handling.

## 6. Memory & Conversation Handling

- **ConversationBufferMemory:**  
Stores the full chat history, enabling context-aware and follow-up questions.

## 7. Utility Libraries

- **torch:**  
Supports PyTorch-based models, useful for GPU checks or model compatibility.
- **time:**  
Handles delays, tracks performance, or controls animation effects (like loading spinners).

```

25 @st.cache_resource 1 usage
26 def load_models():
27     device = "cuda" if torch.cuda.is_available() else "cpu"
28     EMBED_MODEL = SentenceTransformer(model_name_or_path="paraphrase-multilingual-MiniLM-L12-v2", device=device)
29     GROQ_API_KEY = "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"
30     os.environ['GROQ_API_KEY'] = GROQ_API_KEY
31     groq_model = ChatGroq(model_name="llama3-70b-8192", temperature=0.3, max_tokens=512)
32     urdu_splitter = RecursiveCharacterTextSplitter(
33         separators=[".", "\n\n", "\n", ".", " "],
34         chunk_size=400,
35         chunk_overlap=100,
36         length_function=len
37     )
38     return EMBED_MODEL, groq_model, urdu_splitter
39 
```

`st.cache_resource` helps to remember (or cache) expensive operations like loading big machine learning models so they don't run again and again every time the user interacts with the app. Without caching, the app would reload everything every time a button is clicked. With `@st.cache_resource`, the app only loads it once, and reuses it

The first time `load_model()` runs, it loads the model and stores it in memory. After that, any time you call `load_model()`, it uses the already-loaded model.



This function `load_models()` is designed to prepare and load everything your chatbot needs to work efficiently. First, it checks if a GPU is available (using CUDA) and chooses the best device (GPU or CPU) for faster performance. Then it loads a multilingual sentence embedding model called `paraphrase-multilingual-MiniLM-L12-v2`, which helps the app convert user input or PDF content into vector format for smart searching. After that, it sets up an API key for using Groq's powerful language model (LLaMA 3) to generate intelligent replies. It creates a `ChatGroq` object with specific settings like how creative the response should be (temperature) and how long the replies can be (`max_tokens`). Finally, it defines how to split Urdu text into smaller parts using the `RecursiveCharacterTextSplitter`, so long texts (like PDFs) can be broken down and processed effectively. In the end, it returns all three tools — the embedding model, the chatbot model, and the text splitter — ready to use in your app.

```
26
27 def create_faiss_vectorstore(text, EMBED_MODEL):
28     chunks = urdu_splitter.split_text(text[:100000])
29     embeddings = EMBED_MODEL.encode(chunks, show_progress_bar=True)
30     embeddings = np.array(embeddings).astype('float32')
31     dimension = embeddings.shape[1]
32     index = faiss.IndexFlatL2(dimension)
33     index.add(embeddings)
34     return {"chunks": chunks, "index": index}
35
```

The `create_faiss_vectorstore` function is designed to process a large Urdu text and prepare it for efficient semantic search using FAISS. It first splits the input text into smaller chunks using an Urdu-specific text splitter. These chunks are then converted into vector embeddings using a provided embedding model (`EMBED_MODEL`). Since FAISS requires embeddings to be in `float32` format, the function ensures this conversion. It then initializes a FAISS index using L2 (Euclidean) distance and adds the generated embeddings to this index. Finally, the function returns a dictionary containing both the list of text chunks and the FAISS index, which can later be used for searching relevant information from the original text.

```
38
39 def init_memory(): 1 usage
40     return ConversationBufferMemory(
41         memory_key="history",
42         input_key="input",
43         return_messages=True
44     )
45
```

This function sets up conversation memory using `ConversationBufferMemory`, which helps the chatbot remember previous user inputs and responses. This makes the conversation more natural and coherent like how a human would remember what was said earlier.

```

47 # Session State Setup
48 if "messages" not in st.session_state:
49     st.session_state.messages = []
50 if "processing" not in st.session_state:
51     st.session_state.processing = False
52 if "vectorstore" not in st.session_state:
53     st.session_state.vectorstore = None
54 if "last_processed" not in st.session_state:
55     st.session_state.last_processed = None
56 if "audio_cooldown" not in st.session_state:
57     st.session_state.audio_cooldown = 0
58 if "qa_cache" not in st.session_state:
59     st.session_state.qa_cache = {}
60 if "conversation_memory" not in st.session_state:
61     st.session_state.conversation_memory = init_memory()
62 if "show_memory" not in st.session_state:
63     st.session_state.show_memory = False
64

```

Streamlit resets everything on each user interaction, which means that without storing information in `st.session_state`, all previous data would be lost after each action. To prevent this and maintain the chatbot's functionality, the code initializes several session variables. The `messages` variable keeps track of the conversation history, while `processing` acts as a flag to indicate if the chatbot is currently working on a response, which can be helpful for showing loading animations. The `vectorstore` is used to store vector embeddings, such as those generated from PDF text, enabling quick and efficient semantic searches. The `last_processed` variable stores the name or timestamp of the most recently processed file to avoid redundant operations. `audio_cooldown` ensures that the audio recording feature is not triggered too frequently, improving usability. The `qa_cache` holds previously asked questions and their answers to avoid repeating computations and to improve performance. The `conversation_memory` stores past inputs and responses using the `init_memory()` function so that the chatbot can maintain context across the conversation. Lastly, `show_memory` is a toggle that allows developers to show or hide memory-related information in the user interface. Additionally, the line `EMBED_MODEL, groq_model, urdu_splitter = load_models()` loads the core components of the app: `EMBED_MODEL` converts text into numerical vectors, `groq_model` generates responses using Groq's powerful language model, and `urdu_splitter` breaks long Urdu text into manageable chunks for better processing and retrieval.

```

69 with st.sidebar:
70     st.title("Settings")
71     uploaded_file = st.file_uploader("Upload Urdu PDF", type=["pdf"])
72     if uploaded_file is not None and st.session_state.vectorstore is None:
73         with st.spinner("Processing PDF..."):
74             with BytesIO(uploaded_file.getvalue()) as pdf_data:
75                 doc = fitz.open(stream=pdf_data, filetype="pdf")
76                 raw_text = "".join([page.get_text() for page in doc])
77                 doc.close()
78                 st.session_state.vectorstore = create_faiss_vectorstore(raw_text, EMBED_MODEL)
79                 st.success("PDF loaded successfully!")
80
81     if st.button("Clear Conversation Memory"):
82         st.session_state.conversation_memory.clear()
83         st.success("Conversation memory cleared!")
84     st.sidebar.checkbox("Show Memory", key="show_memory")
85

```

This sidebar section titled "Settings" allows users to upload Urdu PDF files, which are processed only if they haven't been handled before. When a PDF is uploaded, the app reads and extracts its text, converts it into vector embeddings using the embedding model, and stores these vectors for fast searching. A loading spinner informs the user during processing, and a success message confirms when the PDF is ready. Additionally, there is a button to clear the chatbot's conversation memory, allowing the chat history to be reset, with a confirmation message shown after clearing. Finally, a checkbox lets users toggle whether the conversation memory details are displayed within the app interface.

```

89 def urdu_rag_query(question: str) -> str: 2 usages
90     # Return cached answer if available
91     if question in st.session_state.qa_cache:
92         return st.session_state.qa_cache[question]
93
94     # Ensure vectorstore is loaded
95     vectorstore = st.session_state.vectorstore
96     # Retrieve context
97     q_emb = EMBED_MODEL.encode([question])
98     q_emb = np.array(q_emb).astype('float32')
99     distances, indices = vectorstore["index"].search(q_emb, 15)
100     context = "\n\n".join([vectorstore["chunks"][i] for i in indices[0]])
101
102     # Initialize ConversationChain if needed
103     if "conversation_chain" not in st.session_state:
104         prompt = PromptTemplate(
105             input_variables=["history", "input", "context"],
106             template="""
107             آپ کا کردار: ایک مایہ ناز اردو استاد
108             ہدایات:

```

```

109     # 1. پچھلی گفتگو اور دیا گیا سیاق و سباق استعمال کریں۔
110     # 2. اگر سوال مختصر اور معلوماتی ہو، تو مختصر جواب دیں۔
111     # 3. اگر سوال میں وضاحت یا خلاصہ مانگا گیا ہو، تو تفصیلی جواب دیں۔
112     # 4. جہاں ضروری ہو وہاں مشکل اردو الفاظ کی وضاحت کریں۔
113     # 5. جواب کے شروع میں 'جواب:' نہ لکھیں۔
114
115     پچھلی گفتگو:
116     {history}
117
118     سیاق:
119     {context}
120
121     سوال:
122     {input}
123
124     جواب:
125     """
126
127     st.session_state.conversation_chain = LLMChain(

```

```

127         st.session_state.conversation_chain = LLMChain(
128             llm=groq_model,
129             prompt=prompt,
130             memory=st.session_state.conversation_memory,
131             verbose=False
132         )
133     # Build history string from memory
134     history = "\n".join([
135         f"<انسان>: {m.content}" if m.type == "human" else f"AI: {m.content}"
136         for m in st.session_state.conversation_memory.chat_memory.messages
137     ])
138     resp = st.session_state.conversation_chain.predict(
139         input=question,
140         context=context
141     )
142     st.session_state.qa_cache[question] = resp
143     return resp
144

```

This function `urdu_rag_query` is designed to answer user questions in Urdu by combining semantic search with a powerful language model, creating a smooth and context-aware chatbot experience.

1. **Caching for Efficiency:** It first checks if the question was already asked and answered before, using a cache stored in `st.session_state.qa_cache`. If yes, it returns the cached answer instantly, saving time and resources.
2. **Vector Search for Context Retrieval:** If the question is new, it converts the question into a vector embedding using the embedding model. It then searches a pre-built vector index (stored in `vectorstore`) to find the most relevant text chunks from the uploaded Urdu PDF or other sources. These chunks provide the context needed to answer the question accurately.
3. **Dynamic Prompt Setup:** If the conversation chain (which ties the language model, memory, and prompts together) hasn't been initialized yet, it sets up a custom prompt template in Urdu. This prompt instructs the language model to behave like an expert

Urdu teacher, use the conversation history and retrieved context, and give clear, detailed, or concise answers depending on the question.

4. **Memory for Conversation Flow:** The function builds a history of the conversation from previous messages, allowing the chatbot to maintain context and respond naturally, not just in isolation.
5. **Generating Response:** The language model (Groq-powered) is called with the question, the retrieved context, and the conversation history to generate a coherent and relevant reply in Urdu.
6. **Caching the Answer:** The new response is saved in the cache to speed up future identical queries.

```
# Text-to-speech helper
def text_to_speech(text: str, lang: str = "ur") -> BytesIO: 2 usages
    tts = gTTS(text=text, lang=lang, slow=False)
    b = BytesIO(); tts.write_to_fp(b); b.seek(0)
    return b
```

The `text_to_speech` function converts written text into spoken audio using the Google Text-to-Speech (gTTS) library, with Urdu as the default language. It takes the input text and generates speech at a normal speed, then saves the audio not as a file on disk but directly into an in-memory buffer called BytesIO. This approach allows the audio to be stored temporarily in the computer's memory, making it quick and easy to access without needing to read or write files on disk. After writing the audio data to this buffer, the function resets the reading position to the start so the audio can be played or processed immediately. This function is useful in applications like chatbots or web apps where you want to provide instant spoken responses without any delay caused by file handling.

```

152 st.title("🇵🇰 Urdu Tutor")
153 for msg in st.session_state.messages:
154     with st.chat_message(msg["role"]):
155         st.markdown(
156             f'<div style="text-align: right; direction: rtl;">{msg["content"]}</div>',
157             unsafe_allow_html=True
158         )
159         if msg.get("audio"): st.audio(msg["audio"], format="audio/mp3")
160 if st.session_state.show_memory:
161     st.subheader("Conversation Memory")
162     st.code(st.session_state.conversation_memory.buffer_as_str, language="text")
163 t1, t2 = st.columns([5,1])
164 with t1:
165     user_input = st.chat_input("Type your question...", key="input")
166 with t2:
167     audio_bytes = None
168     if time.time() - st.session_state.audio_cooldown > 3:
169         audio_bytes = audio_recorder(text="", pause_threshold=2.5, key="recorder")
170     else:
171         st.warning("Please wait...")

```

This Streamlit code builds an Urdu chat tutor interface showing past messages aligned for Urdu reading. It supports text input and audio recording, but limits audio recording frequency to avoid overload. Users can also toggle viewing the full conversation memory.

```

175 if audio_bytes and len(audio_bytes)>0 and not st.session_state.processing and current - st.session_state.audio_cooldown>3:
176     st.session_state.processing=True; st.session_state.audio_cooldown=current
177     with st.spinner("Transcribing..."):
178         try:
179             with tempfile.NamedTemporaryFile(suffix=".wav", delete=False) as tmp:
180                 tmp.write(audio_bytes); path=tmp.name
181                 r=sr.Recognizer();
182                 with sr.AudioFile(path) as source:
183                     audio=r.record(source); text=r.recognize_google(audio, language="ur-PK")
184                 os.remove(path)
185                 if text and text!=st.session_state.last_processed:
186                     st.session_state.messages.append({"role":"user","content":text})
187                     st.session_state.last_processed=text
188                     if st.session_state.vectorstore:
189                         ans=urdu_rag_query(text)
190                         audio=text_to_speech(ans)
191                         st.session_state.messages.append({"role":"assistant","content":ans,"audio":audio})
192                         st.rerun()
193                     else: st.warning("Upload a PDF first.")
194         except Exception as e:
195             st.error(f"Error: {e}")
196     finally: st.session_state.processing=False

```

This code handles voice input by converting recorded Urdu audio into text, then uses that text to get an answer from the chatbot if a PDF is loaded. It adds both the text and audio reply to the chat and prevents multiple rapid inputs using timers and flags, showing errors if anything goes wrong.

```

198 # Process text input
199 if user_input and not st.session_state.processing and user_input!=st.session_state.last_processed:
200     st.session_state.processing=True
201     st.session_state.messages.append({"role":"user","content":user_input})
202     st.session_state.last_processed=user_input
203     if st.session_state.vectorstore:
204         with st.spinner("Generating answer..."):
205             ans=urdu_rag_query(user_input)
206             audio=text_to_speech(ans)
207             st.session_state.messages.append({"role":"assistant","content":ans,"audio":audio})
208     else:
209         st.warning("Upload a PDF first.")
210         st.session_state.processing=False
211         st.rerun()
212

```

This code processes the user's typed question. If the input is new and the system isn't already processing, it stores the question, checks if a PDF is uploaded, and then generates an answer using the chatbot. The answer is converted to audio, both are saved in the chat, and the app refreshes to display the response.

## 6. Conclusion:

This Urdu RAG Chatbot integrates multiple advanced components to provide an interactive, context-aware tutoring experience in Urdu. It efficiently processes uploaded Urdu PDFs by splitting the text and creating a FAISS vector store for fast semantic retrieval. The chatbot uses an embedding model combined with a powerful LLM (Groq) to generate contextually relevant answers while maintaining conversation memory for a natural, coherent dialogue flow. The system supports both text and voice input, utilizing speech recognition for Urdu speech-to-text and text-to-speech for audio responses, enhancing accessibility and user engagement. Caching mechanisms optimize repeated queries for faster response times. The Streamlit UI offers seamless interaction, including conversation memory control and PDF uploads, making the chatbot a robust tool for Urdu language learning and information retrieval.

