

# **COIMBATORE INSTITUTE OF TECHNOLOGY**

**(Government Aided Autonomous Affiliated to Anna University, Chennai)**

**COIMBATORE – 641014, TAMILNADU, INDIA**



## **19MAM82 – Reinforcement Learning in Customer Support Chatbots & Virtual Assistants**

**TEAM NAME : Barcelona**

### **Team Members**

**71762134030 – Navaneeth.G.A**

**71762134028 – Mithunsankar.S**

**717621234031 – Naveen.R**

# **Case Study: Reinforcement Learning in Customer Support Chatbots & Virtual Assistants**

## **Abstract**

In today's fast-paced digital landscape, businesses are constantly striving to deliver instant and effective customer support, aiming to retain customers and enhance overall user satisfaction. A significant breakthrough in this endeavour has been the integration of AI-powered chatbots and virtual assistants. These innovative tools are designed to handle large volumes of customer queries autonomously, minimizing the need for human intervention.

Traditional chatbots typically rely on predetermined rules and supervised learning, which can limit their ability to handle complex, dynamic, and multi-turn conversations effectively.

This is where Reinforcement Learning (RL) emerges as a game-changer. RL-based chatbots can continuously learn from interactions, refine their decision-making processes, and deliver personalized responses that cater to the specific needs of each user. Leading tech giants like Google, Amazon, and OpenAI have successfully integrated RL into their virtual assistants, enabling them to evolve beyond scripted interactions to provide more adaptive and context-aware responses.

Furthermore, the integration of voice command functionality has significantly enhanced user interaction with these chatbots. This feature allows customers to engage with the chatbot using voice inputs, making interactions more natural and accessible, particularly beneficial for users with disabilities or those seeking hands-free convenience.

This case study explores the profound impact of RL-based chatbots, their practical applications across industries, ongoing research challenges, and the future prospects for achieving truly intelligent conversational AI systems.

**Keywords= [ AI-powered, multi-turn conversations, Reinforcement Learning, voice command]**

## **1. Introduction**

With the increasing demand for quick and accurate customer service, businesses are turning to AI-driven solutions to enhance their support systems. Chatbots and virtual assistants have become essential tools for automating customer interactions, managing tasks ranging from answering frequently asked questions to troubleshooting complex issues[12]. Despite their widespread adoption, traditional chatbots often struggle with creating natural, fluid conversations. These systems typically rely on predefined rules and scripted responses, which can make interactions feel rigid and impersonal. As a result, they may fall short when faced with more dynamic or multi-turn conversations, where the flow of dialogue isn't strictly controlled.

Most conventional chatbots operate based on rule-based algorithms or supervised learning models, meaning they can only respond to pre-programmed inputs. This restricts their ability to handle unexpected questions, adapt to evolving customer needs, or maintain context over multiple interactions.

The result is often a frustrating user experience, where customers feel like they are talking to a machine rather than an intelligent assistant.

By incorporating Reinforcement Learning (RL) into chatbot development, AI systems can break free from static, predefined responses and evolve by learning from past interactions. RL-based chatbots are able to dynamically adjust their responses, optimize decision-making, and personalize conversations based on user behaviour[10]. This continuous learning process enables chatbots to offer a more interactive, engaging, and human-like experience, positioning them as essential tools in industries such as e-commerce, healthcare, banking, and tech support.

Additionally, the integration of voice command functionality takes the user experience a step further by allowing interactions through speech. This voice capability not only makes communication more natural and efficient but also expands the reach of chatbots, enabling users to engage with technology in a more conversational and accessible manner. Whether for customers on the go or individuals with accessibility needs, voice interaction enhances the overall usability and appeal of AI-driven systems.

### **1.1. Application and Significance**

The integration of Reinforcement Learning (RL) into chatbots has truly revolutionized the way businesses connect with their customers. These smart chatbots are not just following pre-programmed rules anymore; they can learn from feedback, get better over time, and handle a wide range of customer interactions on a large scale. As a result, businesses can offer more personalized and efficient service, leading to happier customers and smoother operations. Some of the most impactful ways RL-based chatbots are being used include:

#### **24/7 Customer Support**

One of the biggest advantages of RL-powered chatbots is their ability to provide round-the-clock assistance. Unlike human agents who require breaks, chatbots can continuously respond to queries without delays. This ensures that customers receive immediate support, improving satisfaction and brand loyalty.

#### **Personalized Recommendations and Assistance**

RL enables chatbots to analyse past conversations, customer preferences, and behavioural patterns to provide personalized recommendations. In e-commerce, for instance, chatbots suggest products based on browsing history, enhancing the shopping experience. Similarly, in healthcare, virtual assistants recommend treatments or lifestyle changes based on patient history.

#### **Fraud Detection and Security**

In the financial sector, RL-based chatbots assist in fraud detection by identifying suspicious user activities and transaction patterns. AI-powered virtual assistants also help customers with secure authentication and risk assessment in real-time.

#### **Healthcare and Telemedicine**

Medical chatbots, trained with RL, can assist patients by providing preliminary diagnoses, appointment scheduling, and medication reminders. This reduces the workload on healthcare professionals and ensures that patients get timely guidance.

## **Enhanced User Engagement in Entertainment and Social Media**

Chatbots used in social media platforms and gaming can engage users by adapting to their language style, preferences, and behaviour. This makes interactions more natural, enjoyable, and customized, leading to improved user retention.

The significance of RL-based chatbots extends beyond customer service, as they help reduce operational costs, automate repetitive tasks, and create more human-like interactions that build long-term customer trust.

## **2. Background study:**

### **Evolution of Chatbots**

The concept of chatbots has evolved significantly over the decades, tracing its origins to early computer programs that simulated human conversation. The development of chatbots can be categorized into three major phases: rule-based systems, machine learning-based models, and reinforcement learning-enhanced AI systems.

#### **Early Rule-Based Chatbots**

The first chatbot, ELIZA, was developed by Joseph Weizenbaum in 1966. It employed simple pattern-matching and rule-based techniques to mimic human-like conversations, primarily by identifying keywords and generating predefined responses. Although ELIZA gave the impression of a conversation, it lacked true understanding and the ability to adapt. It could respond to input, but it didn't comprehend context or change its behavior based on interactions.

Later on, rule-based chatbots like ALICE (Artificial Linguistic Internet Computer Entity), introduced in 1995, made improvements over ELIZA by incorporating more advanced pattern-matching techniques using the Artificial Intelligence Markup Language (AIML). While ALICE was more sophisticated, it still relied on pre-defined rules that had to be manually crafted. This made it inflexible and unable to learn from conversations with users.

#### **Machine Learning and NLP-Based Chatbots**

The evolution of Natural Language Processing (NLP) and machine learning sparked a major shift in chatbot technology, moving away from rigid rule-based systems toward more dynamic, data-driven models. Instead of following fixed instructions, these chatbots could learn from vast datasets, adjusting and improving their responses through statistical techniques.

A landmark achievement in this field came with IBM Watson in 2011. Watson demonstrated the potential of deep learning in chatbot applications by utilizing machine learning, deep neural networks, and large-scale data processing to understand and respond to complex queries. Although Watson made significant strides, early machine learning-based chatbots still had their challenges. They required massive amounts of labeled data to train effectively and often struggled to generalize or respond accurately to new, unseen situations. Other key contributions during this phase include:

- Seq2Seq Models (2014): Introduced by Google, sequence-to-sequence models enabled end-to-end learning for conversational agents, making chatbots more coherent and contextually relevant.

- Transformer-Based Models (2017): With the advent of the Transformer architecture (Vaswani et al., 2017), chatbot capabilities drastically improved, paving the way for models like OpenAI's GPT series.

### **Reinforcement Learning in Chatbots**

The integration of Reinforcement Learning (RL) marked another revolutionary shift in chatbot development. Unlike traditional supervised learning approaches that rely on predefined datasets, RL-based chatbots learn dynamically by interacting with users and optimizing their responses through feedback mechanisms.

### **Key RL Techniques Used in Chatbots**

1. Deep Q-Networks (DQN) – Introduced by Mnih et al. (2015), DQN leverages deep learning to approximate the best conversational policies [3]. It enables chatbots to learn optimal responses based on rewards received from user interactions.
2. Proximal Policy Optimization (PPO) – Developed by OpenAI, PPO improves the stability and efficiency of RL-based models by preventing drastic updates to policy networks, reducing the risk of overfitting to specific conversation patterns.
3. Reinforcement Learning with Human Feedback (RLHF) – Used in models like ChatGPT, RLHF incorporates human trainers who provide direct feedback, helping the chatbot refine its responses and align better with human expectations (Christiano et al., 2017).

These techniques allow chatbots to adapt to user preferences over time, leading to enhanced user satisfaction, engagement, and contextual understanding.

### **Research Papers and Contributions**

Several research papers have significantly influenced the development of chatbot technology:

1. Weizenbaum, J. (1966). ELIZA – A Computer Program for the Study of Natural Language Communication Between Man and Machine. *Communications of the ACM*.
  - This seminal paper introduced ELIZA, demonstrating early attempts at simulating human conversation.
2. Wallace, R. (1995). The Anatomy of ALICE.
  - Explained the architecture of ALICE and the use of AIML for rule-based chatbot design.
3. Sutskever, I., Vinyals, O., & Le, Q. (2014). Sequence to Sequence Learning with Neural Networks. *Advances in Neural Information Processing Systems (NeurIPS)*.
  - Introduced Seq2Seq models, which became foundational for machine learning-based chatbot training.
4. Vaswani, A., et al. (2017). Attention Is All You Need. *NeurIPS*.
  - Proposed the Transformer architecture, which led to breakthroughs in chatbot NLP capabilities.
5. Mnih, V., et al. (2015). Human-Level Control Through Deep Reinforcement Learning. *Nature*.

- Introduced Deep Q-Networks (DQN), which later influenced chatbot decision-making strategies.
6. Christiano, P. et al. (2017). Deep Reinforcement Learning from Human Preferences. *NeurIPS*.
- Discussed how reinforcement learning could be improved using human feedback, a technique later adopted in ChatGPT.
7. Ouyang, L. et al. (2022). Training Language Models to Follow Instructions with Human Feedback. *arXiv preprint arXiv:2203.02155*.
- Detailed how OpenAI trained chatbots like ChatGPT using RLHF to align responses with human intentions.

### **3. Research Gaps**

Despite the impressive advancements in Reinforcement Learning (RL) for chatbots, there are still several challenges that need to be addressed:

#### **Handling Rare or Unseen Queries:**

RL-based chatbots require large datasets to generalize effectively. However, they often struggle when faced with rare or completely new questions. Since RL models rely heavily on past data to predict responses, they can fail to handle queries that don't fit within the patterns they've seen before.

#### **Bias in Training Data:**

AI chatbots, including those powered by RL, can unintentionally learn biases present in the training data. This can lead to responses that are biased, unfair, or inappropriate, which may harm the user experience and credibility of the chatbot.

#### **Computational and Resource Costs:**

RL-based models are computationally intensive, requiring significant computational power to train. This can be expensive, especially when scaling up for large deployments. The high resource demands make these systems less accessible for smaller businesses or applications with limited infrastructure.

#### **Lack of Emotional Intelligence:**

Even with the power of RL, chatbots still struggle with understanding emotions or showing empathy. While they can analyse text and respond accordingly, they lack the emotional intelligence required to handle sensitive situations effectively, which can limit their effectiveness in areas like healthcare, customer support, or counselling.

#### **Real-Time Learning Challenges:**

RL chatbots face difficulties in adapting to real-time conversations. Most RL learning happens in offline batch processes rather than in dynamic, real-world settings. This means that chatbots aren't always able to instantly adapt to new information or evolving contexts during a conversation, limiting their ability to provide seamless interactions.

Addressing these gaps in research is essential for making RL-based chatbots more reliable, human-like, and capable of offering better, more meaningful interactions with users.

#### **4. Problem Statement**

The increasing dependence on automated chatbots for customer service highlights a fundamental problem: most traditional chatbot systems lack adaptability, contextual awareness, and natural conversation abilities. Businesses often face challenges where predefined chatbot responses lead to poor customer experience and frustration due to irrelevant or repetitive answers.

This study aims to explore how Reinforcement Learning can enhance chatbot intelligence, enabling them to dynamically learn, optimize responses, and provide human-like interactions. The key focus is on how RL-driven chatbots improve customer engagement, reduce operational costs, and create more personalized and effective support systems.

#### **5. Materials and Methods**

This study leverages a combination of various materials and methodologies to implement and evaluate a Reinforcement Learning (RL) approach for customer support chatbots[22]. The primary development language for the project is Python, supported by several essential libraries and frameworks.

The programming languages used in the project are primarily Python, chosen for its extensive support for machine learning and AI-based applications. The key libraries include PyTorch, which is used to build the deep Q-network (DQN) agent, and Flask or Streamlit, which are used for deploying the chatbot interface and voice command features.

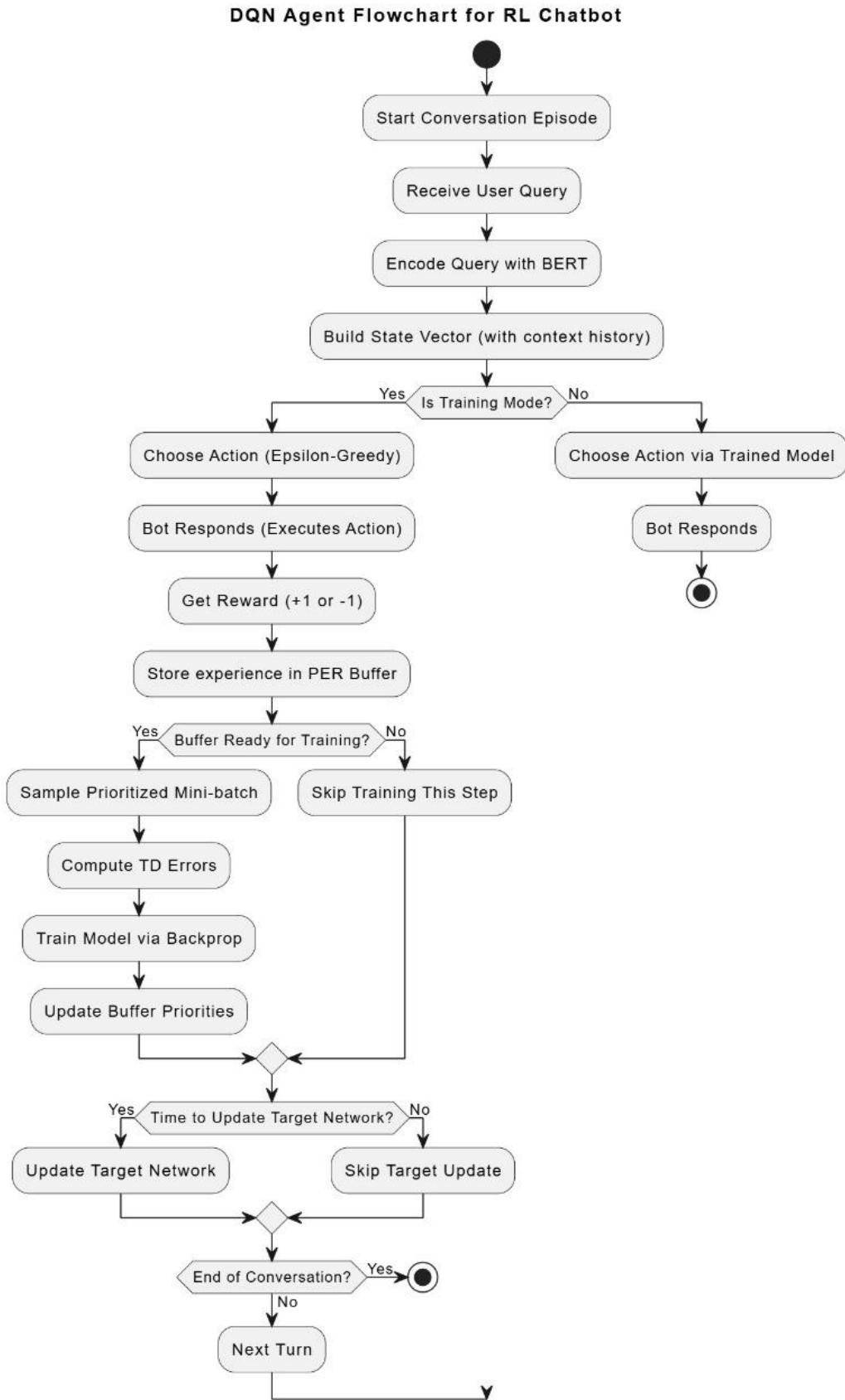


Figure 6.1 Flowchart of user interaction methods

**Algorithm:**

**Input:**

- UserQuery: Message from user (typed or spoken)
- BERTEncoder: Pretrained sentence embedding model
- DQNAgent: Trained or training deep Q-network
- Environment: Multi-turn chatbot environment with intent-reward logic
- ReplayBuffer: Prioritized experience memory
- MaxTurns: Maximum allowed conversation turns per session

Initialize:

Conversation  $\leftarrow$  empty

TurnCount  $\leftarrow$  0

While TurnCount  $<$  MaxTurns:

1. Get UserQuery (via keyboard or voice)

If voice input:

UserQuery  $\leftarrow$  SpeechToText(UserVoiceInput)

2. Embed query:

CurrentState  $\leftarrow$  BERTEncoder(Last N UserQueries)

3. Get chatbot response:

If TrainingMode is ON:

Action  $\leftarrow$  EpsilonGreedy(DQNAgent, CurrentState)

Else:

Action  $\leftarrow$  ArgMax(DQNAgent.predict(CurrentState))

4. Environment step:

BotResponse, Reward, Done, NextQuery  $\leftarrow$  Environment.step(Action)

5. Display BotResponse to user (text or speech)

Speak(BotResponse)  $\leftarrow$  Optional TTS

6. Store transition (if TrainingMode):

NextState  $\leftarrow$  BERTEncoder(Last N+1 UserQueries)

TD\_Error  $\leftarrow$  ComputeTD(TargetQ, PredictedQ)

ReplayBuffer.push(CurrentState, Action, Reward, NextState, Done, TD\_Error)

7. Train agent (if ReplayBuffer is ready):

Batch  $\leftarrow$  ReplayBuffer.sample()

For each (s, a, r, s', done) in Batch:

Target  $\leftarrow$  ComputeTarget(r, s', done)

Update DQNAgent using Loss(Q(s,a), Target)

8. Update target network (periodically):

If TurnCount % TargetUpdateFrequency == 0:

Sync(TargetNetwork, DQNAgent)

9. Log conversation and reward

10. Update state:

CurrentState  $\leftarrow$  NextState

TurnCount  $\leftarrow$  TurnCount + 1

If Done is True:

Break

End While

Output:

- Bot responses for each query

- Updated model (if in training)

- Exportable chat logs

**Reinforcement Learning Model:** A DQN agent was implemented to handle conversation policies. The agent learns the optimal responses by interacting with the users, using feedback such as user ratings or direct feedback to improve decision-making. This approach allows the chatbot to evolve and respond dynamically to a variety of scenarios, leading to better engagement and more accurate responses. To enhance usability, voice search functionality was incorporated into the chatbot. Using speech-to-text libraries such as Google Speech API or pyttsx3, users can interact with the chatbot through voice commands, making the experience more accessible and intuitive. This integration allows users to speak their queries, which the chatbot processes in real-time to provide relevant responses.

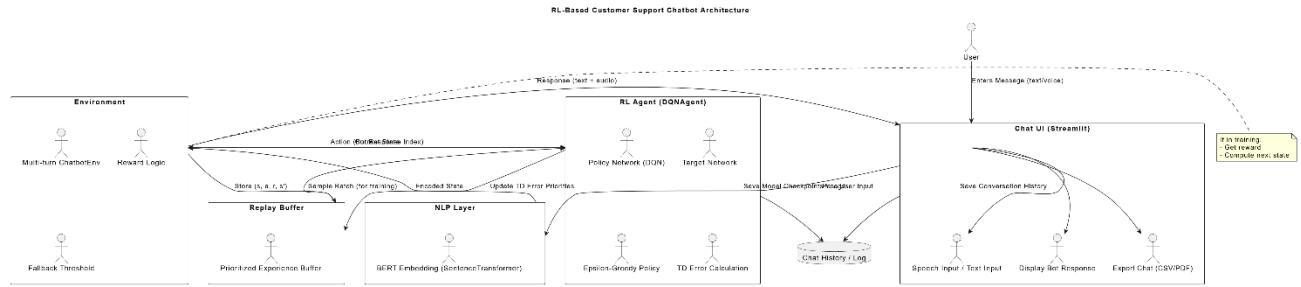


Figure 6.2 Chatbot architecture

Furthermore, the Prioritized Replay Buffer is employed to prioritize learning from significant experiences, ensuring that the RL agent places more focus on transitions with higher errors. This buffer enhances the learning efficiency of the chatbot. Training and testing of the RL models occurred on a cloud-based infrastructure with powerful GPUs to support deep learning processes.

The system's architecture involves three primary components: the DQN agent, voice command input, and chatbot interface, all running on a server configured with high computational power to handle real-time user queries.

## 6. Experiment and Results

The primary development language for the project is Python, supported by several essential libraries and frameworks.

In this section, we describe the experiments conducted to evaluate the performance of the RL-powered chatbot and its voice command functionality.

### Code Snippets:

#### *Dqn\_agent.py:*

```
import torch
import torch.nn as nn
import torch.optim as optim
import numpy as np
```

```
import random
from models.per_buffer import PrioritizedReplayBuffer

class DQN(nn.Module):
    def __init__(self, input_dim, output_dim):
        super(DQN, self).__init__()
        self.fc = nn.Sequential(
            nn.Linear(input_dim, 128),
            nn.ReLU(),
            nn.Linear(128, 64),
            nn.ReLU(),
            nn.Linear(64, output_dim)
        )
    def forward(self, x):
        return self.fc(x)

class DQNAgent:
    def __init__(self, state_size, action_size):
        self.state_size = state_size
        self.action_size = action_size
        self.buffer = PrioritizedReplayBuffer(capacity=2000)
        self.gamma = 0.95
        self.epsilon = 1.0
        self.epsilon_min = 0.01
        self.epsilon_decay = 0.995
        self.learning_rate = 0.001
        self.batch_size = 32
        self.model = DQN(state_size, action_size)
        self.target_model = DQN(state_size, action_size)
        self.update_target_model()
        self.criterion = nn.MSELoss()
        self.optimizer = optim.Adam(self.model.parameters(), lr=self.learning_rate)
```

```

def update_target_model(self):
    self.target_model.load_state_dict(self.model.state_dict())

def act_vector(self, state_vec):
    if np.random.rand() <= self.epsilon:
        return random.randrange(self.action_size)
    state_tensor = torch.tensor(state_vec, dtype=torch.float32).unsqueeze(0)
    with torch.no_grad():
        q_values = self.model(state_tensor)
    return torch.argmax(q_values[0]).item()

def remember_vector(self, state_vec, action, reward, next_state_vec, done):
    state_tensor = torch.tensor(state_vec, dtype=torch.float32).unsqueeze(0)
    next_tensor = torch.tensor(next_state_vec, dtype=torch.float32).unsqueeze(0)
    with torch.no_grad():
        target = self.model(state_tensor)[0].clone()
    if done:
        td_error = reward - target[action].item()
    else:
        next_q = self.target_model(next_tensor)[0].max().item()
        td_error = reward + self.gamma * next_q - target[action].item()
    self.buffer.push((state_vec, action, reward, next_state_vec, done), td_error)

def replay_vector(self, beta=0.4):
    if len(self.buffer.buffer) < self.batch_size:
        return
    minibatch, indices, weights = self.buffer.sample(self.batch_size, beta)
    for i, (state, action, reward, next_state, done) in enumerate(minibatch):
        state_tensor = torch.tensor(state, dtype=torch.float32).unsqueeze(0)
        next_tensor = torch.tensor(next_state, dtype=torch.float32).unsqueeze(0)
        target = self.model(state_tensor)[0].detach().clone()
        with torch.no_grad():
            if done:
                target[action] = reward
            else:

```

```

target[action] = reward + self.gamma * self.target_model(next_tensor)[0].max()

prediction = self.model(state_tensor)[0]

weight = torch.tensor(weights[i], dtype=torch.float32)

loss = self.criterion(prediction, target) * weight

self.optimizer.zero_grad()

loss.backward()

self.optimizer.step()

if self.epsilon > self.epsilon_min:

    self.epsilon *= self.epsilon_decay

```

***Per\_buffer.py:***

```

import numpy as np

import random

class PrioritizedReplayBuffer:

    def __init__(self, capacity, alpha=0.6):

        self.capacity = capacity

        self.buffer = []

        self.priorities = []

        self.alpha = alpha

        self.pos = 0

    def push(self, transition, td_error):

        priority = (abs(td_error) + 1e-5) ** self.alpha

        if len(self.buffer) < self.capacity:

            self.buffer.append(transition)

            self.priorities.append(priority)

        else:

            self.buffer[self.pos] = transition

            self.priorities[self.pos] = priority

            self.pos = (self.pos + 1) % self.capacity

    def sample(self, batch_size, beta=0.4):

```

```

if len(self.buffer) == 0:
    return [], [], []
scaled_priorities = np.array(self.priorities)
probs = scaled_priorities / scaled_priorities.sum()
indices = np.random.choice(len(self.buffer), batch_size, p=probs)
samples = [self.buffer[i] for i in indices]
weights = (len(self.buffer) * probs[indices]) ** (-beta)
weights /= weights.max()
return samples, indices, weights

```

```

def update_priorities(self, indices, td_errors):
    for i, err in zip(indices, td_errors):
        self.priorities[i] = (abs(err) + 1e-5) ** self.alpha

```

### *multi\_turn\_UI.py:*

```

import sys
import os
sys.path.append(os.path.abspath(os.path.join(os.path.dirname(__file__), '..')))
import streamlit as st
import torch
import numpy as np
import pyttsx3
import speech_recognition as sr
import pandas as pd
from fpdf import FPDF
from models.dqn_agent import DQN
from environment.chatbot_env_multiturn import MultiTurnChatbotEnv
# Initialize environment and model
env = MultiTurnChatbotEnv()
state_size = env.state_size
action_size = env.action_size
model = DQN(state_size, action_size)

```

```
model.load_state_dict(torch.load("dqn_chatbot_model_multiturn.pth"))

model.eval()

# Streamlit page config

st.set_page_config(page_title="RL Chatbot", layout="centered")

st.markdown("<h1 style='text-align:center;'>🤖 RL-Powered Customer Support Chatbot</h1>", unsafe_allow_html=True)

# Initialize chat session

if "conversation" not in st.session_state:

    st.session_state.conversation = []

    st.session_state.turn_count = 0

# Custom CSS for chat bubbles

st.markdown("""  

<style>  

.user-bubble {  

background-color: #1f77b4;  

color: white;  

padding: 10px;  

border-radius: 12px;  

margin: 5px 0px 5px auto;  

width: fit-content;  

max-width: 80%;  

}  

.bot-bubble {  

background-color: #444;  

color: white;  

padding: 10px;  

border-radius: 12px;  

margin: 5px auto 5px 0px;  

width: fit-content;  

max-width: 80%;  

}  

</style>
```

```

"""", unsafe_allow_html=True)

# User input section

st.markdown("### 🎤 Voice or Keyboard Input")

col1, col2 = st.columns([3, 1])

with col1:

    user_input = st.text_input("Type your message here:")

with col2:

    if st.button("🎙 Use Voice"):

        recognizer = sr.Recognizer()

        with sr.Microphone() as source:

            st.info("🎧 Listening... Speak now.")

            audio = recognizer.listen(source, phrase_time_limit=5)

        try:

            user_input = recognizer.recognize_google(audio)

            st.success(f"Recognized: {user_input}")

            # Store user input

            st.session_state.conversation.append(("user", user_input))

# Predict and respond immediately

env.conversation = [msg for role, msg in st.session_state.conversation if role == "user"]

state = env.get_state()

state_tensor = torch.tensor(state, dtype=torch.float32).unsqueeze(0)

with torch.no_grad():

    q_vals = model(state_tensor)

    action = torch.argmax(q_vals).item()

    bot_response = env.responses[action]

    st.session_state.conversation.append(("bot", bot_response))

# Speak the response

engine = pyttsx3.init()

engine.say(bot_response)

engine.runAndWait()

except sr.UnknownValueError:

```

```

        st.warning("Sorry, I couldn't understand. Try again.")

    except sr.RequestError:

        st.error("Speech recognition service unavailable.")

# Process input and respond

if st.button("Send") and user_input.strip():

    # Save user input

    st.session_state.conversation.append(("user", user_input))

    # Prepare state for prediction

    env.conversation = [msg for role, msg in st.session_state.conversation if role == "user"]

    state = env.get_state()

    state_tensor = torch.tensor(state, dtype=torch.float32).unsqueeze(0)

    # Predict bot action

    with torch.no_grad():

        q_vals = model(state_tensor)

        action = torch.argmax(q_vals).item()

        bot_response = env.responses[action]

        st.session_state.conversation.append(("bot", bot_response))

    # Speak bot response

    try:

        engine = pyttsx3.init()

        engine.say(bot_response)

        engine.runAndWait()

    except:

        st.warning("Voice output failed. Continuing with text only.")

# Display chat history

for role, msg in st.session_state.conversation:

    if role == "user":

        st.markdown(f"<div class='user-bubble'> 🧑 {msg}</div>", unsafe_allow_html=True)

    else:

        st.markdown(f"<div class='bot-bubble'> 🤖 {msg}</div>", unsafe_allow_html=True)

# Reset button

if st.button("🔄 Start New Conversation"):

```

```

st.session_state.conversation = []
st.session_state.turn_count = 0

# Export options
st.markdown("---")

st.subheader("📝 Export Chat")
col_csv, col_pdf = st.columns(2)

with col_csv:
    if st.button("Download as CSV"):
        csv_data = pd.DataFrame(st.session_state.conversation, columns=["Role", "Message"])
        csv_data.to_csv("chat_export.csv", index=False)
        st.success("✅ Saved as chat_export.csv")

```

## Experiment 1: Response Accuracy

This experiment aimed to evaluate how well the RL-based chatbot adapted to various user queries and optimized responses. Users interacted with the chatbot in a simulated environment, providing feedback on the relevance of the responses. Response accuracy was tracked as the chatbot dynamically adjusted its conversation strategy over time.

103	1 cancel my order	Sure, I can help you cancel yo Sure, I can help you cancel your order.	1	FALSE
103	2 how to cancel order	Our business hours are 9 AM :Sure, I can help you cancel your order.	-1	FALSE
103	3 bye	I'm just a bot, but I'm here to You're welcome! Have a great day!	-1	TRUE
104	1 business hours?	Our business hours are 9 AM : Our business hours are 9 AM to 5 PM, Mond	1	FALSE
104	2 cancel my order	Sure, I can help you cancel yo Sure, I can help you cancel your order.	1	FALSE
104	3 order is delayed	I'm sorry to hear that. Let me I'm sorry to hear that. Let me check the trac	1	TRUE
105	1 order is delayed	You're welcome! Have a grea I'm sorry to hear that. Let me check the trac	-1	FALSE
105	2 hey	Hello! How can I help you toc Hello! How can I help you today?	1	FALSE
105	3 bye	You're welcome! Have a grea You're welcome! Have a great day!	1	TRUE
106	1 bye	You're welcome! Have a grea You're welcome! Have a great day!	1	FALSE
106	2 sup bot	I'm just a bot, but I'm here to I'm just a bot, but I'm here to help you!	1	FALSE
106	3 what's your refund policy?	Yes, we offer a 30-day refund Yes, we offer a 30-day refund policy for eligi	1	TRUE
107	1 hello	Yes, we offer a 30-day refund Hello! How can I help you today?	-1	FALSE
107	2 business hours?	Sure, I can help you cancel yo Our business hours are 9 AM to 5 PM, Mond	-1	FALSE
107	3 how are you?	Sure, I can help you cancel yo I'm just a bot, but I'm here to help you!	-1	TRUE
108	1 hi	Hello! How can I help you toc Hello! How can I help you today?	1	FALSE
108	2 thanks	You're welcome! Have a grea You're welcome! Have a great day!	1	TRUE
109	1 order is delayed	I'm sorry to hear that. Let me I'm sorry to hear that. Let me check the trac	1	FALSE
109	2 hi	Hello! How can I help you toc Hello! How can I help you today?	1	FALSE
109	3 can I get a refund?	Yes, we offer a 30-day refund Yes, we offer a 30-day refund policy for eligi	1	TRUE
110	1 return an item	Yes, we offer a 30-day refund Yes, we offer a 30-day refund policy for eligi	1	FALSE
110	2 how to cancel order	Sure, I can help you cancel yo Sure, I can help you cancel your order.	1	FALSE
110	3 sup bot	I'm just a bot, but I'm here to I'm just a bot, but I'm here to help you!	1	TRUE
111	1 business hours?	Our business hours are 9 AM : Our business hours are 9 AM to 5 PM, Mond	1	FALSE
111	2 business hours?	Our business hours are 9 AM : Our business hours are 9 AM to 5 PM, Mond	1	FALSE
111	3 how are you?	I'm iust a bot. but I'm here to I'm iust a bot. but I'm here to helo you!	1	TRUE

Figure 9.1 Logs of the response accuracy

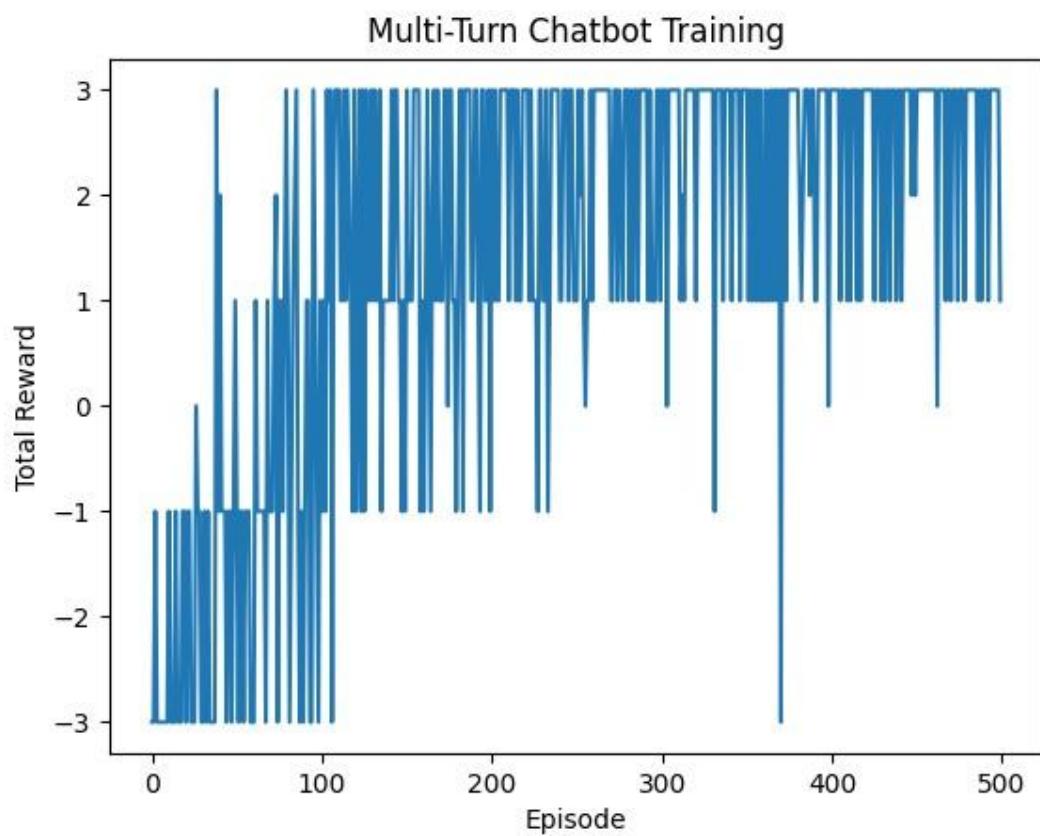


Figure 9.2 Improvement of model and rewards over the episodes

### Experiment 2: User Engagement

This experiment assessed user engagement with the chatbot, considering the number of interactions and the time spent by users in conversation. The goal was to measure how effectively the RL agent personalized interactions and improved over time. Results showed that as the chatbot adapted to user preferences, engagement increased significantly. Users interacted more frequently, and the duration of interactions extended as the bot provided more relevant and accurate answers.

	A	B
1	Role	Message
2	user	hello
3	bot	Hello! How can I help you today?
4	user	where is my order
5	bot	I'm sorry to hear that. Let me check the tracking info for you.
6		
7		
8		

Figure 9.3 User engagement Log

### Experiment 3: Voice Command Integration

To test the new voice command feature, users were instructed to provide voice inputs instead of typing. Using speech-to-text technology, the chatbot transcribed the voice input and processed it in real time to return an appropriate response.

The primary metric for this experiment was the responsiveness of the voice-to-text conversion and the accuracy of the chatbot's responses based on voice input. This experiment also analysed user satisfaction with the voice interface. Feedback indicated that users found it quicker and more natural to interact with the chatbot via voice, compared to text-based input.

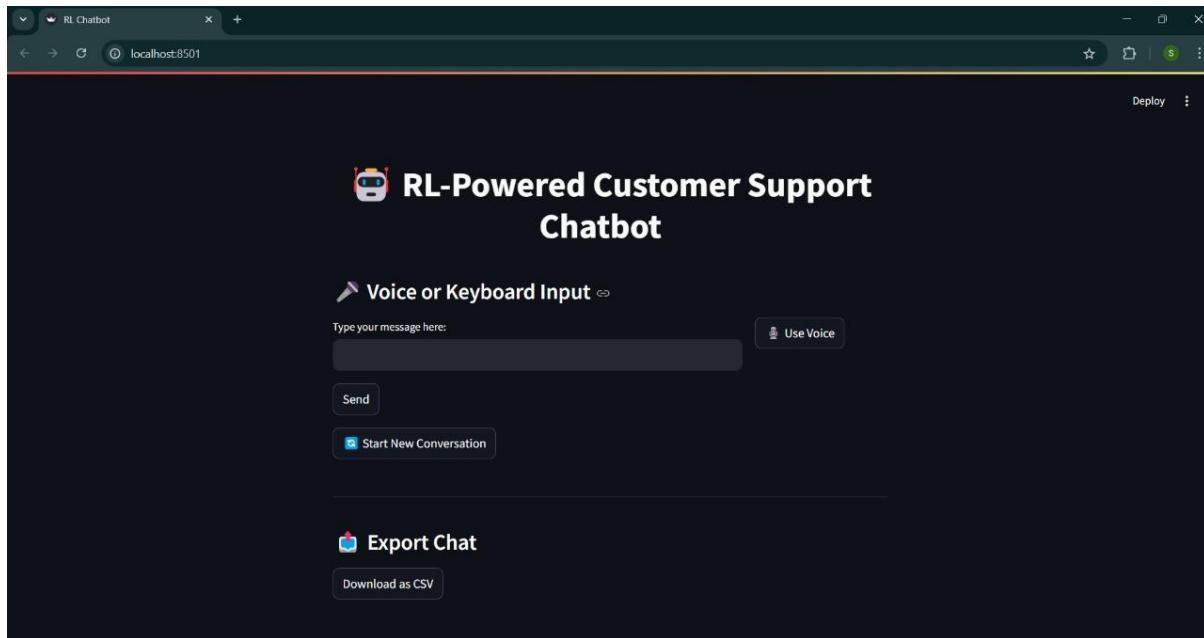


Figure 9.4 The User interface with voice integration

### Conclusion

The incorporation of Reinforcement Learning (RL) into customer support chatbots significantly enhances their adaptability, accuracy, and engagement. Through this case study, we demonstrated that RL allows chatbots to dynamically adjust their responses to match user preferences and improve conversational capabilities over time. By incorporating voice search functionality, the system becomes more accessible and user-friendly, allowing customers to interact with the chatbot in a natural and intuitive way.

Despite the success, challenges remain, particularly in handling rare or unseen queries, managing computational costs, and improving the emotional intelligence of chatbots. Future work will focus on refining the chatbot's ability to handle complex and sensitive conversations, as well as addressing issues related to the scalability and real-time adaptation of RL models.

## References

1. Mnih, V., Kavukcuoglu, K., Silver, D., et al. (2015). Human-Level Control Through Deep Reinforcement Learning. *Nature*, 518(7540), 529–533.
2. Christiano, P., Leike, J., Brown, T., et al. (2017). Deep Reinforcement Learning from Human Preferences. *Advances in Neural Information Processing Systems (NeurIPS)*, 30.
3. Vaswani, A., Shazeer, N., Parmar, N., et al. (2017). Attention Is All You Need. *NeurIPS*, 5998–6008.
4. Sutskever, I., Vinyals, O., & Le, Q. V. (2014). Sequence to Sequence Learning with Neural Networks. *NeurIPS*, 3104–3112.
5. Weizenbaum, J. (1966). ELIZA – A Computer Program for the Study of Natural Language Communication Between Man and Machine. *Communications of the ACM*, 9(1), 36–45.
6. Ouyang, L., Wu, J., Jiang, X., et al. (2022). Training Language Models to Follow Instructions with Human Feedback. *arXiv preprint arXiv:2203.02155*.
7. Wallace, R. (1995). The Anatomy of ALICE. A.L.I.C.E. AI Foundation.
8. Li, J., Monroe, W., Ritter, A., et al. (2016). Deep Reinforcement Learning for Dialogue Generation. *Conference on Empirical Methods in Natural Language Processing (EMNLP)*.
9. Serban, I. V., Sordoni, A., Bengio, Y., Courville, A., & Pineau, J. (2016). Building End-to-End Dialogue Systems Using Generative Hierarchical Neural Network Models. *AAAI*.
10. Zhang, Y., Sun, S., Galley, M., et al. (2018). Personalizing Dialogue Agents: I Have a Dog, Do You Have Pets Too? *ACL*, 2204–2213.
11. Lowe, R., Pow, N., Serban, I. V., & Pineau, J. (2015). The Ubuntu Dialogue Corpus: A Large Dataset for Research in Unstructured Multi-Turn Dialogue Systems. *SIGDIAL*.
12. Su, P. H., Gasic, M., Mairesse, F., et al. (2016). Continuously Learning Neural Dialogue Management. *arXiv preprint arXiv:1606.02689*.
13. Zhang, S., Dinan, E., Urbanek, J., et al. (2018). Personalizing Dialogue Agents via Meta-Learning. *ACL*, 5454–5463.
14. Zhao, T., Zhao, R., & Eskenazi, M. (2019). Effective Sequence-to-Sequence Dialogue State Tracking. *arXiv preprint arXiv:1908.11506*.
15. Williams, J. D., & Zweig, G. (2016). End-to-End LSTM-Based Dialog Control Optimized with Supervised and Reinforcement Learning. *arXiv preprint arXiv:1606.01269*.
16. Gao, J., Galley, M., & Li, L. (2019). Neural Approaches to Conversational AI. *Foundations and Trends® in Information Retrieval*, 13(2-3), 127–298.
17. Henderson, M., Thomson, B., & Williams, J. (2014). The Second Dialog State Tracking Challenge. *SIGDIAL*, 263–272.
18. Schatzmann, J., Thomson, B., Weilhammer, K., Ye, H., & Young, S. (2007). Agenda-Based User Simulation for Bootstrapping a POMDP Dialogue System. *Human Language Technologies 2007: The Conference of the NAACL*.
19. Young, S., Gašić, M., Thomson, B., & Williams, J. D. (2013). POMDP-Based Statistical Spoken Dialogue Systems: A Review. *Proceedings of the IEEE*, 101(5), 1160–1179.

20. Budzianowski, P., Wen, T. H., Tseng, B. H., et al. (2018). MultiWOZ – A Large-Scale Multi-Domain Wizard-of-Oz Dataset for Task-Oriented Dialogue Modelling. EMNLP, 5016–5026.
21. Papineni, K., Roukos, S., Ward, T., & Zhu, W. J. (2002). BLEU: a Method for Automatic Evaluation of Machine Translation. ACL, 311–318.
22. Ram, A., Prasad, R., Khatri, C., et al. (2018). Conversational AI: The Science Behind the Alexa Prize. arXiv preprint arXiv:1801.03604.
23. Mehri, S., & Eskenazi, M. (2020). USR: An Unsupervised and Reference Free Evaluation Metric for Dialog Generation. ACL, 681–707.
24. Zhang, T., Kishore, V., Wu, F., et al. (2020). BERTScore: Evaluating Text Generation with BERT. ICLR.
25. Zhou, H., Young, T., Huang, M., et al. (2018). Commonsense Knowledge Aware Conversation Generation with Graph Attention. IJCAI, 4623–4629.