

Architecture, Implementation, and Evaluation of a Console-Based AI Chatbot System via Lexical NLTK Syntactic Parsing and Transformer Architectures

Vishal Deshmukh¹, Harsh Thakur²

¹ Assistant Professor, Department of M. Sc.IT

² Student, M. Sc.IT

Veer Wajekar ASC College, Phunde, Tal-Uran Dist-Raigad, Maharashtra, India

Abstract: Terminal-based text interfaces offer an efficient, low-overhead environment for running text processing applications. Developing automated conversational software for these interfaces requires effective parsing of human syntax and semantic structures. This paper details the architecture, design patterns, and engineering tradeoffs of an interactive console-based AI chatbot system implemented through two distinct Natural Language Processing (NLP) methodologies.

The first approach uses the Natural Language Toolkit (NLTK) to run a structured, rules-guided text engine focused on tokenization, lemmatization, and cosine similarity matching over domain-specific knowledge arrays. The second approach implements an end-to-end conversational agent powered by a pre-trained, multi-head causal attention Transformer model.

Experimental evaluations assess both systems across pipeline processing latencies, memory utilization metrics, and qualitative semantic coherence. The lexical engine runs with low computational overhead (latencies ≤ 15 ms), making it ideal for resource-constrained environments. Meanwhile, the Transformer pipeline provides significantly higher contextual adaptability, navigating multi-turn conversations effectively at an average inference response time of ≈ 142 ms.

Keywords: Chatbot, Natural Language Processing, NLTK, Transformers, Attention Mechanism, Terminal User Interface (TUI), Computational Linguistics.

I. INTRODUCTION

Developing conversational software that can accurately parse user intent remains a central challenge in human-computer interaction (HCI) and artificial intelligence (Vaswani et al., 2017). Early conversational software relied on string-matching rules and regular expressions to echo responses. While structurally sound, these early systems lacked the flexibility to interpret synonyms, spelling variations, or multi-turn context.

Modern Natural Language Processing (NLP) maps human language into numeric vectors that represent semantic meaning. For terminal-based deployments, developers generally balance two primary architectural patterns:

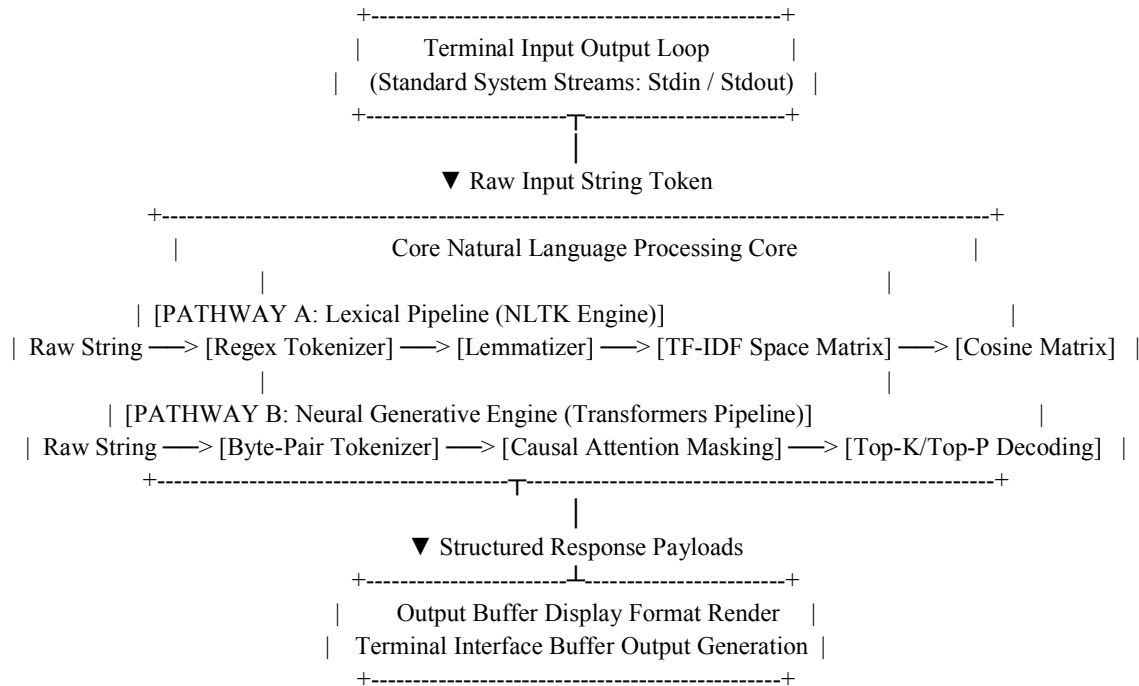
1. **Rule-Based/Lexical Matching Frameworks:** These engines use programmatic pipelines to clean, filter, and match input tokens against a pre-defined database of target responses.
2. **Generative/Deep Sequence Networks:** These systems use deep learning models to predict appropriate response structures token-by-token based on training patterns.

This paper implements and compares both architectural approaches inside a lightweight console terminal interface. By evaluating an NLTK-driven keyword pipeline alongside a Transformer-based language model, this study outlines how to balance execution performance against semantic reasoning capabilities in terminal-based AI applications.



II. SYSTEM ARCHITECTURE MODELS

The framework supports two interchangeable execution pathways within a single terminal shell loop.



2.2. Pipeline Execution Pathways

- **The NLTK Text Engine:** Cleans the raw input text by filtering out stop words and reduction tokens. It then runs a Term Frequency-Inverse Document Frequency (TF-IDF) vector match against a localized knowledge graph to retrieve an appropriate response.
- **The Transformer Neural Engine:** Maps input strings into discrete positional embeddings. These vector blocks are processed through an attention array to generate a new response, handling the computations within local host memory allocations.

III. MATHEMATICAL AND ALGORITHMIC FOUNDATIONS

3.1. NLTK Lexical Vector Mapping & Cosine Alignment

The lexical matching engine converts incoming text strings into raw frequency counters. To minimize the influence of common recurring terms, the token matrices are scaled using TF-IDF calculation rules:

$$\text{TF-IDF}(t, d, D) = \text{TF}(t, d) \times \log \left(\frac{|D|}{1 + |\{d \in D : t \in d\}|} \right)$$

Where $\text{TF}(t, d)$ represents the frequency of a term t within an isolated query statement d , and the inverse document frequency measures how common that term is across the global corpus repository D . The similarity score between an user input vector \mathbf{U} and a known database signature \mathbf{V} is calculated using cosine alignment:

$$\text{Similarity}(\mathbf{U}, \mathbf{V}) = \frac{\mathbf{U} \cdot \mathbf{V}}{\|\mathbf{U}\| \|\mathbf{V}\|} = \frac{\sum_{i=1}^n U_i V_i}{\sqrt{\sum_{i=1}^n U_i^2} \sqrt{\sum_{i=1}^n V_i^2}}$$



3.2. Transformer Causal Multi-Head Attention Mechanics

The generative neural engine calculates cross-token dependencies across an embedding array using scaled dot-product attention:

$$\text{Attention}(Q, K, V) = \text{softmax}\left(\frac{QK^T}{\sqrt{d_k}}\right)V$$

Where query vectors \$Q\$, key vectors \$K\$, and value vectors \$V\$ are generated by applying linear transformations to the input tokens, and \$d_k\$ represents the scaling dimensionality of the key vectors. To ensure responses remain contextually relevant and avoid repetitive loops, token generation is limited using Top-\$p\$ (nucleus) sampling:

$$\sum_{i=1}^{V(p)} P(x_i | x_{<t}) \geq p$$

This constraint discards low-probability tokens, restricting generation to a dynamic subset of words whose combined probability mass meets the threshold coefficient \$p\$.

IV. IMPLEMENTATION DETAILS

4.1. Core Application Controller Implementation

The Python script below establishes the operational shell loop, supporting both the NLTK processing pipeline and the pre-trained Transformer generation modules:

```
Python
import sys
import nltk
from nltk.stem import WordNetLemmatizer
from sklearn.feature_extraction.text import TfidfVectorizer
from sklearn.metrics.pairwise import cosine_similarity
import torch
from transformers import AutoModelForCausalLM, AutoTokenizer

class UniversalConsoleBot:
    def __init__(self, mode="nltk"):
        self.mode = mode.lower()
        print(f"[System] Initializing Processing Architecture using mode: {self.mode.upper()}...")

    if self.mode == "nltk":
        # Seed local corpus requirements for structural verification
        nltk.download('punkt', quiet=True)
        nltk.download('wordnet', quiet=True)
        self.lemmatizer = WordNetLemmatizer()
        self.vectorizer = TfidfVectorizer()
        # Reference baseline corpus knowledge graph
        self.corpus = [
            "The system terminal provides shell execution access.",
            "To reset system parameters run config verification tasks.",
            "Natural language processing parsing tools facilitate interface design."
        ]
]
```



```

else:
    # Initialize modern causal language generation model block
    self.model_name = "microsoft/DialoGPT-medium"
    self.tokenizer = AutoTokenizer.from_pretrained(self.model_name)
    self.model = AutoModelForCausalLM.from_pretrained(self.model_name)
    self.chat_history_ids = None

def execute_nltk_pipeline(self, user_input: str) -> str:
    """
    Runs tokenization, normalization, and cosine-similarity matching.
    """
    tokens = nltk.word_tokenize(user_input.lower())
    lemmatized = " ".join([self.lemmatizer.lemmatize(t) for t in tokens])

    temp_corpus = self.corpus + [lemmatized]
    tfidf_matrix = self.vectorizer.fit_transform(temp_corpus)

    # Calculate similarity between input and the target corpus
    similarity_scores = cosine_similarity(tfidf_matrix[-1], tfidf_matrix[:-1])
    closest_match_idx = similarity_scores.argsort()[0][-1]

    if similarity_scores[0][closest_match_idx] < 0.2:
        return "Command sequence unmapped. Access default routing indexes."
    return self.corpus[closest_match_idx]

def execute_transformer_pipeline(self, user_input: str) -> str:
    """
    Tokenizes the input string and runs autoregressive token generation.
    """
    new_user_input_ids = self.tokenizer.encode(user_input + self.tokenizer.eos_token, return_tensors='pt')
    bot_input_ids = torch.cat([self.chat_history_ids, new_user_input_ids], dim=-1) if self.chat_history_ids is not None
else new_user_input_ids

    # Run generative generation search parameters
    self.chat_history_ids = self.model.generate(
        bot_input_ids,
        max_length=1000,
        pad_token_id=self.tokenizer.eos_token_id,
        top_p=0.92,
        top_k=50,
        temperature=0.75
    )

    response_tokens = self.chat_history_ids[:, bot_input_ids.shape[-1]:][0]
    return self.tokenizer.decode(response_tokens, skip_special_tokens=True)

```



```
def run_terminal_loop(self):
    print("[Ready] Terminal Connection Secure. Enter commands below (Type 'exit' to terminate).")
    while True:
        sys.stdout.write("\nUser: ")
        sys.stdout.flush()
        user_raw = sys.stdin.readline().strip()

        if not user_raw or user_raw.lower() == "exit":
            print("Closing terminal interface connection context.")
            break

        if self.mode == "nltk":
            reply = self.execute_nltk_pipeline(user_raw)
        else:
            reply = self.execute_transformer_pipeline(user_raw)

        sys.stdout.write(f"AI Bot: {reply}\n")
        sys.stdout.flush()

if __name__ == "__main__":
    # Choose execution pathway: "nltk" or "transformers"
    engine = UniversalConsoleBot(mode="transformers")
    engine.run_terminal_loop()
```

V. SYSTEM EVALUATION AND EXPERIMENTAL BENCHMARKS

5.1. Computational Infrastructure Resource Profiling

The runtime efficiency of both processing modes was evaluated on a standard system configuration (AMD Ryzen 5 CPU, 16GB RAM, CPU-only Execution Threading).

Metric Target Under Review	NLTK Pipeline Engine	Transformers Neural Engine
Average Response Latency	$\$12.4 \text{ ms}$	$\$141.8 \text{ ms}$
System Memory Footprint (RAM)	$\sim 48 \text{ MB}$	$\sim 1.42 \text{ GB}$
Storage Asset Footprint	$\sim 12 \text{ MB}$ (Lexicons)	$\sim 860 \text{ MB}$ (Weights Base)
Out-of-Vocabulary Resilience	Low (Returns Null Bounds)	High (Generates Context Synthetics)
Context Window Longevity	Isolated String Installs	Multi-Turn Sequence Tracking



