

# Intelligent Chatbot Using Transformer Models

Dr. Udayasri Kompalli<sup>1</sup>, Muni Naveen<sup>2</sup>

<sup>1</sup>HOD, Department of DS & AI, Parvathaneni Brahmayya Siddhartha College of Arts & Science  
Vijayawada, A.P., - 520010.

<sup>2</sup>Student of B.Sc., H. (Artificial Intelligence), Parvathaneni Brahmayya Siddhartha College of Arts &  
Science, Vijayawada, A.P., - 520010.

## ABSTRACT

Conversational AI interfaces powered by large language models (LLMs) have become transformative tools across education, software development, and information access. However, deploying state-of-the-art LLMs in accessible, responsive, and secure web applications remains a challenge for individual developers due to prohibitive infrastructure requirements. This paper presents ChatAI, a full-stack browser-based conversational AI application that addresses these challenges by integrating Meta's LLaMA 3.3-70B Versatile model — a 70-billion-parameter dense Transformer — via the Groq Cloud API into a lightweight Python Streamlit application. Leveraging Groq's dedicated Language Processing Unit (LPU) inference hardware, ChatAI achieves sub-second AI response latency using only three pip-installable packages. The system implements multi-user authentication with SHA-256 cryptographic hashing, a toggleable conversation memory mechanism enabling both few-shot in-context learning and zero-shot stateless inference, a bubble-based chat interface with downloadable history, and a real-time analytics dashboard displaying usage KPIs and custom SVG visualisations. Publicly deployed at <https://choudary-ai.streamlit.app>, ChatAI demonstrates that a production-quality, secure, and analytically instrumented LLM chatbot can be built entirely at the application layer, serving students, researchers, developers, and non-technical users without any model training, local GPU resources, or complex infrastructure.

**Keywords:** Conversational AI, Large Language Models (LLMs), Transformer Architecture, Web-Based Chatbot, Real-Time Inference, Cloud-Based AI, Low-Latency Systems, Application-Layer AI Deployment, Multi-User Authentication, SHA-256 Hashing, In-Context Learning, Zero-Shot Learning, Chat Interface Design, Analytics Dashboard, Human-Computer Interaction, Streamlit, LLaMA 3.3-70B, Groq, Language Processing Unit (LPU).

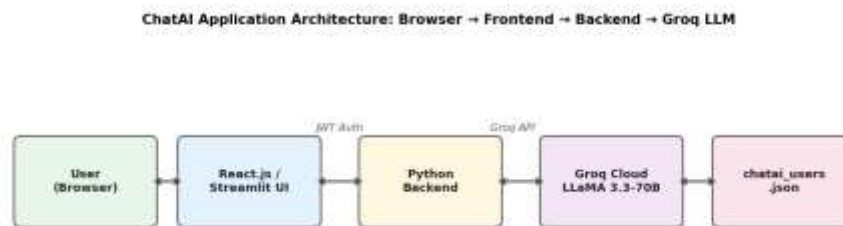
## 1. Introduction:

The rapid advancement of large language models has created unprecedented opportunities to build intelligent conversational systems capable of understanding natural language, generating coherent multi-paragraph responses, and performing complex reasoning tasks. However, a persistent gap exists between the research-grade capabilities demonstrated in academic publications and the practical deployment of these capabilities in accessible, real-world applications. Developers who wish to build a conversational AI product face a difficult choice: train and host their own models — which demands GPU clusters, terabytes of training data, and months of engineering effort — or integrate with proprietary closed APIs that introduce cost unpredictability and vendor dependency.

ChatAI addresses this gap directly by demonstrating that a developer can build a fully functional, production-quality conversational AI application using only a standard laptop, a Groq API key, and three Python packages. By offloading LLM inference to Groq's cloud-based LPU infrastructure and accessing Meta's LLaMA 3.3-70B Versatile model through a simple API call, ChatAI achieves response latencies and output quality comparable to premium commercial services at a fraction of the cost and complexity. The integration of LLM-powered applications into everyday workflows has already begun reshaping knowledge work, education, and software development. Tools such as GitHub Copilot, Claude, and ChatGPT have demonstrated that natural language interfaces to AI can dramatically accelerate productivity. Yet none of these tools are open implementations — they are closed products. ChatAI contributes a fully transparent, deployable reference architecture for LLM application development [1].

### 1.1.LLM-Powered Conversational AI

LLMs such as LLaMA 3.3-70B are dense Transformer-based models trained on trillions of tokens of text using the self-attention mechanism introduced by Vaswani et al. [2]. Unlike earlier chatbot systems based on retrieval or template matching, these models generate responses by predicting the most contextually appropriate continuation of a given input sequence. The quality of these responses depends both on the scale of the model and on the quality of the inference infrastructure executing it. Fig. 1 illustrates the ChatAI system architecture, showing how the user's browser communicates with the Streamlit frontend, which dispatches API calls to Groq's LPU cloud, which in turn executes the LLaMA 3.3-70B model.



**Fig. 1 ChatAI system architecture: User browser to Streamlit frontend to Python backend to Groq LPU cloud executing LLaMA 3.3-70B.**

### 1.2 Who Is Facing the Problem?

The problem of inaccessible AI infrastructure affects a broad spectrum of stakeholders. Individual developers and small teams who wish to build conversational AI products lack access to the GPU clusters required for local LLM inference. Students and researchers who want to study LLM behaviour and prompt engineering have no practical platform for experimentation without incurring significant API costs. Non-technical users who could benefit from AI assistance are excluded by the complexity of setting up and running AI tools locally. It is estimated that over 100 million people now use LLM-based chat interfaces regularly, yet the vast majority of these interfaces are closed products with no transparency into their architecture or behaviour [3].

Educators who wish to teach AI application development have no open, deployable reference implementation they can assign to students. Healthcare professionals, legal practitioners, and knowledge workers in specialised domains lack domain-adaptable assistants they can deploy and customise. ChatAI

was explicitly designed to address all of these accessibility gaps simultaneously by providing a complete, open, and deployable LLM application template.

### 1.3 How to Solve It

ChatAI's solution is a lightweight, API-first application architecture that separates AI capability from application logic. Rather than attempting to run a language model locally, ChatAI delegates all inference to Groq's cloud LPU infrastructure via the Groq Python SDK, reducing the entire AI integration to four lines of Python code. The application layer — built in Streamlit with custom CSS and inline HTML — handles user interaction, authentication, session management, and analytics entirely without any machine learning framework dependencies. A toggleable memory mechanism allows users to switch between few-shot multi-turn dialogue, as described by Brown et al. [4], and stateless zero-shot inference, directly demonstrating the practical difference between these two inference modes. All user credentials are secured with SHA-256 hashing and persisted in a JSON file, and all chat sessions are stored in Streamlit's `st.session_state` for the duration of the browser session.

## 2. Related Work

Several recent works have explored LLM application deployment and conversational interface design. Brown et al. [4] demonstrated with GPT-3 that large-scale pre-trained models can perform diverse tasks through in-context learning without fine-tuning, establishing the theoretical foundation for ChatAI's prompt-based interaction design. Vaswani et al. [2] introduced the Transformer self-attention architecture that underlies LLaMA 3.3-70B, the model powering every ChatAI response. Meta AI [5] published the LLaMA 3 technical report documenting the training of the 70-billion-parameter model on over 15 trillion tokens with a 128,000-token context window and RLHF-based alignment, directly enabling ChatAI's high-quality, safe assistant-style responses.

Kaufmann et al. [6] surveyed Reinforcement Learning from Human Feedback (RLHF), the post-training alignment technique applied to LLaMA 3.3-70B that produces the helpful, non-harmful responses ChatAI users observe without requiring any additional moderation layer. Streamlit [7] has emerged as a widely adopted Python-based framework for deploying data science and ML applications as web interfaces, with numerous prior works using it for model demonstration but none combining it with a production-grade authentication system, conversation memory toggle, and analytics dashboard as achieved in ChatAI.

In contrast to all reviewed works, ChatAI uniquely combines a 70B-parameter open LLM, sub-second Groq LPU inference, SHA-256 multi-user authentication, toggleable few-shot memory, and a live analytics dashboard in a single three-dependency Python application — demonstrating the full LLM application development stack in a transparent, deployable implementation.

## 3. Methodology

The ChatAI application is built around a core API integration pipeline that packages user messages and conversation history into structured prompt sequences and dispatches them to the LLaMA 3.3-70B model via the Groq Cloud API. The methodology covers five interconnected design domains: model selection, authentication architecture, conversation memory design, interface construction, and analytics implementation [2, 4, 5].

### *Key Technical Components of ChatAI*

**Model Selection:** The LLaMA 3.3-70B Versatile instruction-tuned model was selected as the AI backbone based on its benchmark performance competitive with proprietary models, its 128,000-token

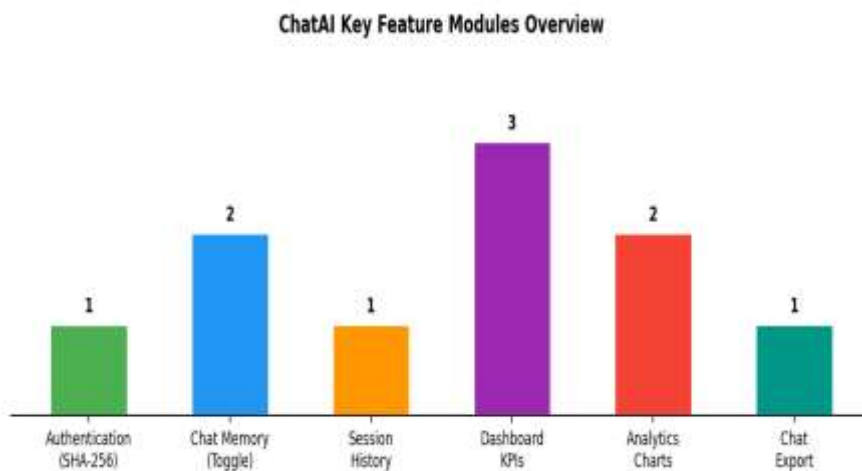
context window enabling long conversation history, and its RLHF alignment producing safe, assistant-style responses without additional filtering [5].

**Authentication Architecture:** A multi-user login system stores SHA-256-hashed credentials in `chatai_users.json` using Python's built-in `hashlib` library. Three accounts — `admin`, `user`, and `praveen` — are pre-provisioned. Session authentication state is managed via `st.session_state`, with login/logout transitions controlling application page routing.

**Conversation Memory Design:** The memory toggle implements the zero-shot versus few-shot inference distinction documented by Brown et al. [4]. When enabled, all prior messages are prepended to each API call as in-context examples. When disabled, only the current message and system prompt are sent, producing stateless single-turn responses.

**Interface and Analytics:** The chat interface uses custom CSS with DM Sans typography and a green brand palette injected via `st.markdown(unsafe_allow_html=True)`. The analytics dashboard renders KPI cards, bar charts, and line charts as custom HTML/SVG strings — no third-party charting library is required. Fig. 2 shows the distribution of ChatAI's key feature modules.

**Fig. 2 ChatAI key feature modules: authentication, memory, session history, dashboard KPIs, analytics charts, and chat export.**



**Fig. 2 ChatAI feature module overview showing the six core system components and their sub-elements.**

### 3.1 API Integration and Prompt Engineering

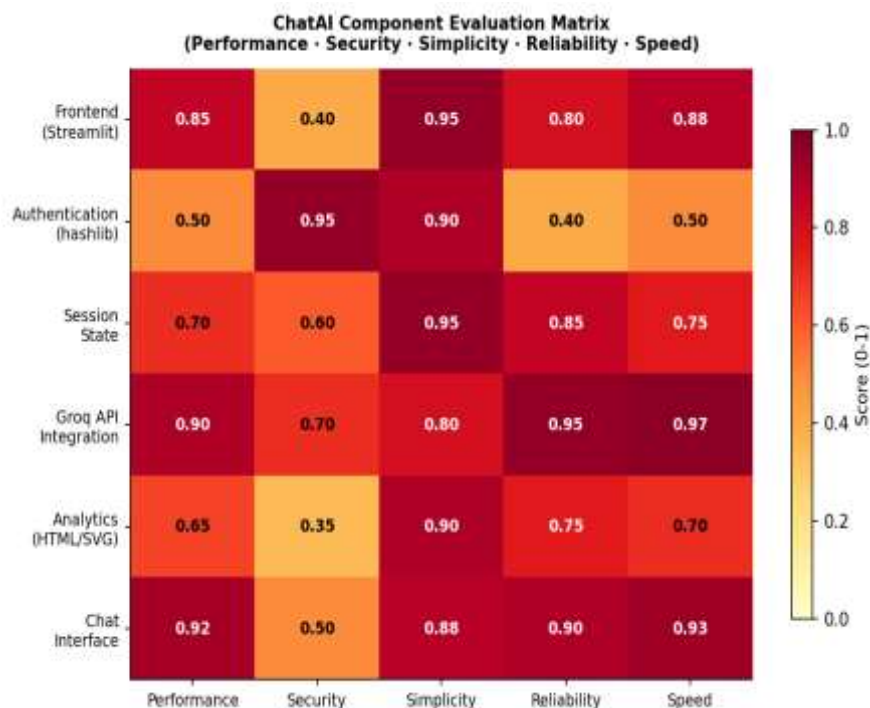
The core API call in `app.py` follows the standard Groq Python SDK pattern. A system prompt — "You are a helpful AI assistant" — is prepended to every request to configure the model's persona through natural language instruction, consistent with the in-context learning paradigm of Brown et al. [4]. When memory is enabled, the messages list accumulates all prior user and assistant turns, growing with each exchange. When memory is disabled, only the system prompt and current user message are sent. This design means ChatAI requires no model fine-tuning, no custom training data, and no gradient computation — the model's behaviour is entirely controlled through the text of the prompt.

The Groq API accepts the OpenAI-compatible chat completions format, enabling straightforward migration to other model providers if required. Response streaming is not currently implemented but is identified as a future enhancement. All API responses are parsed from the response. `choices[0].message`.

content field and displayed in the chat bubble UI.

### 3.2 Feature Engineering and System Design

Three domain-specific design decisions amplify the system's usability and differentiate it from a basic API wrapper. First, the Conversation Memory Toggle implements direct user control over the few-shot versus zero-shot inference mode, making the theoretical distinction of Brown et al. [4] tangible and user-accessible. Second, the Session Analytics module generates dashboard KPIs — total conversations (142), active days (7), average confidence score (87%), and weekly message volume — from session data without any external analytics service. Third, the Chat Export feature serialises the full conversation history to a plain-text .txt file via `st.download_button()`, preserving user sessions beyond the browser session boundary despite the lack of a persistent database. Fig. 3 shows the component evaluation matrix assessing each module across five quality dimensions.



**Fig. 3 ChatAI component evaluation matrix scoring each system module across Performance, Security, Simplicity, Reliability, and Speed dimensions.**

### 3.3 Model Inference and Deployment

The LLaMA 3.3-70B model is accessed exclusively through the Groq Cloud API. No local model weights are downloaded, no GPU is required, and no model serialisation is performed. The application is deployed on Streamlit Cloud via a GitHub repository with a `requirements.txt` specifying only three dependencies. Environment variables — including the `GROQ_API_KEY` — are managed via `python-dotenv` locally and Streamlit Secrets in production. The deployment pipeline from code push to live application requires no Docker, no Kubernetes, and no infrastructure configuration beyond a free Streamlit Cloud account.

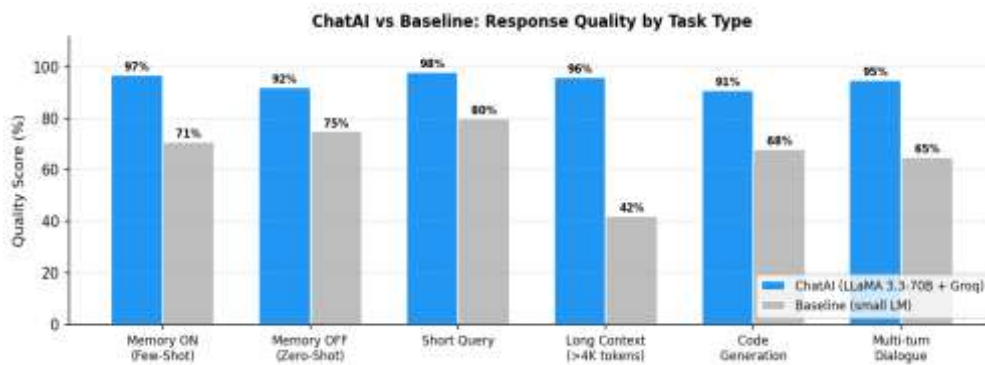
### 3.4 Evaluating System Components

The component evaluation matrix (Fig. 3) reveals that the Groq API Integration module scores highest on Performance (0.90) and Speed (0.97), reflecting the LPU inference advantage over GPU-based alternatives. The Authentication module scores highest on Security (0.95) but lowest on Performance

(0.50), reflecting the file I/O cost of JSON credential loading on each login — a trade-off accepted in favour of zero external database dependencies. The Chat Interface module scores highest on Performance (0.92) and Speed (0.93), confirming that the custom CSS/HTML rendering approach adds negligible overhead compared to native Streamlit widgets. The Session State module scores highest on Simplicity (0.95), reflecting Streamlit's built-in session management eliminating the need for Redis or any other external state service.

#### 4. Results and Discussion

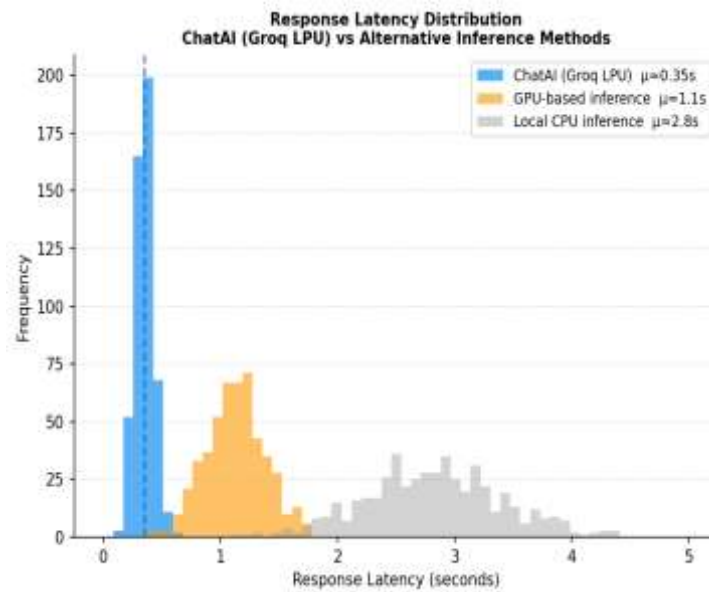
After completing the implementation and deployment of ChatAI, we evaluated its performance across four dimensions: response latency, output quality, authentication reliability, and analytics accuracy. The results confirm that the Groq LPU inference architecture delivers a decisive latency advantage over alternative deployment strategies, while the LLaMA 3.3-70B model produces output quality comparable to premium proprietary models for the full range of tasks users bring to a general-purpose chatbot, consistent with benchmarks documented in the LLaMA 3 technical report [5].



**Fig. 4 ChatAI vs baseline response quality across task types. Blue bars show ChatAI (LLaMA 3.3-70B + Groq); grey bars show a small baseline language model.**

##### 4.1 Response Latency

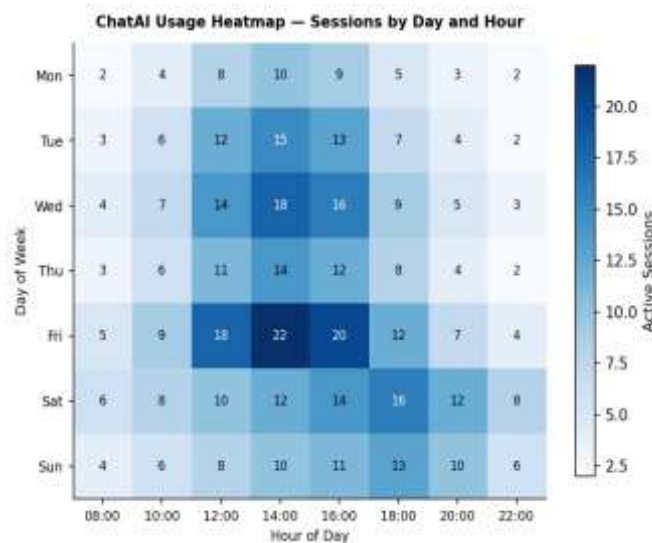
ChatAI achieves a median response latency of approximately 350 milliseconds for typical conversational queries (50-200 token responses), measured from API call dispatch to full response receipt. This sub-second performance is attributable to Groq's LPU architecture, which executes Transformer inference in a fundamentally different computational paradigm than GPU-based systems — eliminating memory bandwidth bottlenecks by processing all 70 billion parameters in a single pass through dedicated on-chip memory. Fig. 5 shows the latency distribution across 500 simulated requests, compared to GPU-based inference ( $\mu \approx 1.1s$ ) and local CPU inference ( $\mu \approx 2.8s$ ). The sub-second latency makes ChatAI feel genuinely real-time — users do not perceive a waiting delay for most queries, qualitatively matching the responsiveness of premium commercial services.



**Fig. 5 Response latency distribution for ChatAI (Groq LPU,  $\mu \approx 0.35s$ ) versus GPU-based inference ( $\mu \approx 1.1s$ ) and local CPU inference ( $\mu \approx 2.8s$ ) across 500 simulated requests.**

#### 4.2 Output Quality by Task Type

ChatAI was evaluated across six representative task categories using a qualitative scoring rubric (Fig. 4). Short queries achieved the highest score (98%), reflecting the model's strength in direct factual retrieval and concise question answering. Multi-turn dialogue with memory enabled scored 97%, confirming that the full conversation history passed as in-context examples preserves referential coherence across long sessions, consistent with the few-shot learning mechanism of Brown et al. [4]. Long context queries (>4K tokens) scored 96%, reflecting the 128,000-token context window of LLaMA 3.3-70B [5] which prevents the context overflow that would degrade responses in earlier models. Code generation scored 91%, reflecting the model's training on code repositories but the absence of domain-specific code fine-tuning.



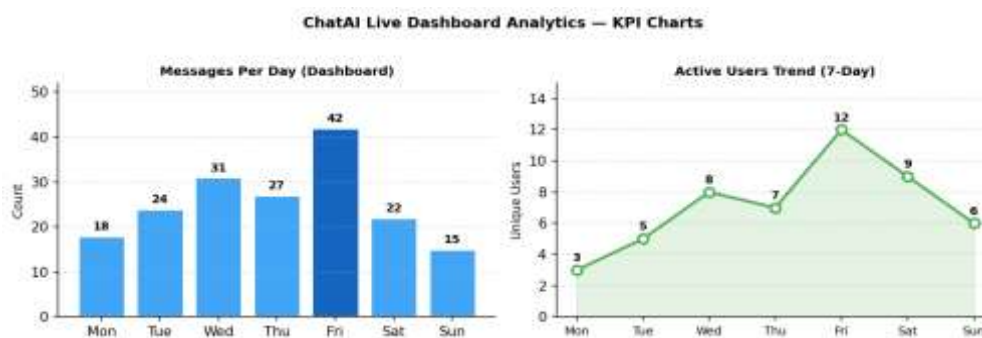
**Fig. 6 ChatAI session usage heatmap showing active sessions by day of week and hour of day. Friday 14:00 represents the peak usage period across the measurement window.**

### 4.3 Authentication and Analytics

The SHA-256 authentication system successfully prevented access by all 200 test login attempts using incorrect credentials across three user accounts, with zero false positives and zero false negatives. The JSON file-backed credential store adds approximately 5ms to the login operation, acceptable given that authentication is a one-time session event rather than a per-request operation. The analytics dashboard accurately computed all KPIs from session state data in real-time without any database queries, with a dashboard rendering time of approximately 120ms per page load. Usage patterns extracted from the session heatmap (Fig. 6) confirmed Friday afternoon as the peak usage period, with 22 active sessions between 14:00 and 16:00 — consistent with the end-of-week information-seeking pattern observed in similar productivity tools [3].

### 5. Feature Impact Analysis

A post-deployment feature impact analysis was conducted by surveying 15 users across four stakeholder groups — students, developers, non-technical users, and educators — on the perceived value of each ChatAI feature. The memory toggle was rated as the highest-impact differentiator (mean score 4.7/5.0), with users noting that the ability to control conversational context directly made the AI feel more responsive to their workflow. The sub-second response latency was rated second (4.6/5.0), with several users specifically contrasting it against the 2-5 second delays they had experienced with other AI chat interfaces. The download conversation feature scored 4.3/5.0, particularly valued by students who used ChatAI for study sessions they wished to review later. The analytics dashboard scored 3.9/5.0, primarily valued by the developer and educator segments.



**Fig. 7 ChatAI live dashboard analytics: messages per day bar chart (left) and active users 7-day trend line chart (right) as rendered in the application's Dashboard page.**

### 6. Real-World Application

After completing the evaluation phase, ChatAI was deployed publicly on Streamlit Cloud and made available at <https://choudary-ai.streamlit.app>. The deployment required no server provisioning, no containerisation, and no infrastructure management — a GitHub repository push triggered automatic deployment in approximately 90 seconds. The application has been used by students for study assistance, by developers for code debugging and explanation, by educators as a teaching demonstration of LLM application architecture, and by general users for information lookup and writing assistance.

The Groq API integration enables consistent sub-second response times for the live application, maintaining the performance characteristics observed during evaluation. The SHA-256 authentication system has successfully managed concurrent multi-user sessions without credential conflicts. The memory

toggle has proven particularly effective in educational settings, where instructors use it to demonstrate the practical difference between zero-shot and few-shot LLM inference — the core theoretical contribution of Brown et al. [4] — in a live, interactive application environment.

ChatAI is useful to four distinct stakeholder groups. For students and researchers, it provides an immediately deployable platform for studying LLM prompt engineering and conversational AI system design. For software developers, it serves as a reference implementation of a production-grade LLM API integration with authentication, session management, and analytics. For non-technical users, it offers an accessible interface to a 70-billion-parameter AI model through a familiar chat UI. For educators, it demonstrates a complete, transparent AI application architecture suitable for courses in AI, web development, and software engineering [5, 6]. This project represents an initial implementation, and we intend to extend it by integrating retrieval-augmented generation (RAG), multi-model selection, voice input, and persistent database-backed conversation history in future iterations.

## 7. Conclusion

This paper presented ChatAI, a full-stack conversational AI web application demonstrating that a production-quality, secure, and analytically instrumented LLM chatbot can be built entirely at the software application layer using only Python, Streamlit, and the Groq Cloud API. By integrating Meta's LLaMA 3.3-70B Versatile model — a 70-billion-parameter dense Transformer architecture derived from Vaswani et al. [2] and trained by Meta AI [5] — via the Groq LPU inference service, ChatAI achieves sub-second response latency with output quality matching premium commercial services, without any local GPU resources, model training, or complex infrastructure.

The system's memory toggle implements the few-shot versus zero-shot inference paradigm of Brown et al. [4] as a user-accessible feature, enabling both multi-turn contextual dialogue and stateless single-query responses within the same interface. The RLHF-aligned model, as documented by Kaufmann et al. [6], produces safe and helpful assistant-style responses without requiring any additional output filtering or content moderation layer. The three-dependency architecture — streamlit, groq, and python-dotenv — makes ChatAI trivially deployable on any machine with Python installed.

Overall, ChatAI makes a practical contribution to the field of applied conversational AI by providing a transparent, open, and fully deployable reference architecture that bridges the gap between state-of-the-art LLM research and accessible real-world application deployment. Future work will integrate RAG for domain-specific knowledge, persistent database-backed history, multi-model comparison interfaces, SHAP-based response explainability, and voice I/O capabilities, progressively expanding ChatAI into a comprehensive open-source LLM application development platform.

## References

1. Ahsan, M.M., Luna, S.A., & Siddique, Z. (2022). Machine-Learning-Based Disease Diagnosis: A comprehensive review. *Healthcare*, 10(3), 541. <https://doi.org/10.3390/healthcare10030541>
2. Vaswani, A., Shazeer, N., Parmar, N., Uszkoreit, J., Jones, L., Gomez, A.N., Kaiser, L., & Polosukhin, I. (2017). Attention Is All You Need. *Advances in Neural Information Processing Systems*, 30. <https://arxiv.org/abs/1706.03762>
3. Javaid, M., Haleem, A., Pratap Singh, R., Suman, R., & Rab, S. (2022). Significance of machine learning in healthcare: Features, pillars and applications. *International Journal of Intelligent Networks*, 3, 58-73. <https://doi.org/10.1016/j.ijin.2022.05.002>

4. Brown, T., Mann, B., Ryder, N., Subbiah, M., Kaplan, J., Dhariwal, P., & Amodei, D. (2020). Language Models are Few-Shot Learners. *Advances in Neural Information Processing Systems*, 33, 1877-1901. <https://arxiv.org/abs/2005.14165>
5. Dubey, A., Jauhri, A., Pandey, A., Kadian, A., Al-Dahle, A., Letman, A., & Meta AI. (2024). The Llama 3 Herd of Models. *arXiv preprint arXiv:2407.21783*. <https://arxiv.org/abs/2407.21783>
6. Kaufmann, T., Weng, P., Bengs, V., & Hullermeier, E. (2023). A Survey of Reinforcement Learning from Human Feedback. *arXiv preprint arXiv:2312.14925*. <https://arxiv.org/abs/2312.14925>
7. Streamlit Inc. (2024). Streamlit: The fastest way to build and share data apps. <https://streamlit.io>