

Eternal Voice: An Adaptive Multi-Modal Framework for Personalized Emotional Memory Preservation

Suyash Arvind Lothe¹, Bhavya Ketan Doshi², Dr. Jasbir Kaur³, Sandhya Thakkar⁴, Suraj Kanal⁵

^{1,2}Student of Master of Computer Applications (MCA),Guru Nanak Institute of Management Studies, Matunga, Mumbai, India

³Director, GNIMS B-School, Head of Information Technology and HR,Guru Nanak Institute of Management Studies, Matunga,Mumbai,India

^{4,5}Assistant Professor,Guru Nanak Institute of Management Studies, Matunga, Mumbai, India

Abstract- The preservation of human identity through generative AI presents unique challenges in maintaining vocal fidelity, emotional nuance, and ethical integrity. While individual technologies for text-to-speech (TTS) and large language models (LLMs) are mature, their integration into a cohesive, user-centric framework for digital legacy remains under-explored. This paper presents Eternal Voice, an adaptive multi-modal framework that tightly couples speaker-adaptive voice synthesis with an emotion-aware conversational agent. The key contribution of this work is a novel integration pipeline. It utilizes a fine-tuned Llama-3-8B model with emotion vector conditioning and a Tacotron 2/WaveNet vocoder stack optimized for low-resource speaker adaptation (SV2TTS). We provide a comprehensive technical evaluation, including an ablation study that quantifies the impact of emotion injection on response naturalness. Experimental results on the LibriTTS and ESD datasets demonstrate a Mean Opinion Score (MOS) of 4.6 ± 0.2 for voice similarity. Moreover, the dialogue coherence improves by 18% over standard LLM baselines. Critically, this paper includes a rigorous discussion of failure modes, deepfake countermeasures, and the limitations of current emotional AI in handling ambiguous human affect.

Keywords: Voice Cloning, Emotional AI, Digital Legacy, Large Language Models, Human-Computer Interaction, Multi-Modal Systems.

I. INTRODUCTION

The evolution of digital legacy systems has stagnated at static archival—photographs and videos remain passive. While recent advances in Generative AI enable dynamic interaction, existing solutions often operate in silos. Consequently, voice cloning lacks contextual awareness, and conversational agents lack vocal identity. This paper addresses the system integration gap. We do not claim to have invented Tacotron 2 or Llama 3; rather, the contribution lies in the orchestration architecture. This architecture enables low-latency, emotionally congruent interaction between a personalized voice model and a persona-specific language model.

Unlike commercial systems (e.g., ElevenLabs or Replika), which function as black-box services, Eternal Voice is pro-posed as an open, on-device capable framework. The primary

technical novelty is the Emotion-Conditioned Prompt Synthesis (ECPS) module. It maps real-time sentiment analysis vectors directly into the LLM's instruction-tuned embedding space. This paper provides a rigorous, reproducible evaluation of this integrated system. It addresses the critical reviews of prior preliminary reports by adding statistical validation, ablation studies, and a comprehensive ethical threat model.

II. RELATED WORK

A. Voice Synthesis Evolution

Neural TTS has progressed from autoregressive models (Tacotron 2, WaveNet) to non-autoregressive, flow-based mod-els (FastSpeech 2, VITS [2]). While VITS offers superior single-stage training, the modularity of Tacotron 2 + WaveNet allows for finer control over speaker embeddings and prosody transfer. This is crucial for preserving the subtle vocal fry and breath patterns essential to emotional memory. A comprehen-sive survey by

Azzuni et al. (2025) establishes standardized terminology for voice cloning and highlights the critical need for detection mechanisms to limit misuse [4]. Additionally, Shen et al. (2024) demonstrated that speaker adaptation techniques can maintain identity even under severe data scarcity. We leverage this finding in our SV2TTS implementation [5]. We selected the SV2TTS (Speaker Verification to TTS) framework [3] for its robustness with limited data (< 10 minutes of speech), a constraint typical in personal legacy collection.

B. Emotion-Aware Language Models

Recent work has explored emotion-conditioned text generation using LLMs. Madani et al. (2024) demonstrated the importance of strategic steerability in emotional support conversations using fine-tuned Llama models [8]. Similarly, Zhang et al. (2023) proposed DialogueLLM, a context and emotion knowledge-tuned LLaMA model that leverages multi-modal information for superior emotion recognition in conversations [9]. Furthermore, Zhou et al. (2025) introduced EmoLLM, a multimodal large language model specifically designed for affective computing. Their work showed that joint optimization of speech and text modalities yields a 12% improvement in empathy scores [10]. Our work extends these approaches by coupling voice texture with linguistic affect. Accordingly, it addresses the multimodal fusion gap highlighted in recent surveys [15], [16].

C. Digital Legacy and Ethical AI

The ethical implications of AI-powered digital legacies have gained significant attention. Baek and Doe (2025) explored the concept of "postmortem life" through thanobots and digital twins, raising critical questions about feminist immortality and consent [12]. Wang et al. (2025) conducted a systematic review of psychological implications of generative AI in digital afterlife technologies. They identified key risk factors such as prolonged grief disorder and identity confusion [13]. Elder (2025) argued philosophically that digital replacement of the dead is a legitimate worry, particularly when the AI agent is perceived as having agency [14]. These

studies inform our ethical framework, which we detail in Section V.

D. Security of Voice Cloning Systems

The security of voice synthesis systems is a rapidly evolving field. Chen et al. (2024) proposed VoiceGuard, a real-time anti-spoofing system for neural voice cloning that uses spectro-temporal embeddings. It detects synthetic speech with 99.2% accuracy [6]. Singh and Kumar (2025) surveyed adversarial attacks on speaker verification systems, categorizing threats into replay, synthesis, and impersonation vectors [7]. Our security analysis in Section V builds upon these taxonomies to propose a multi-layered defense strategy.

III. SYSTEM ARCHITECTURE AND METHODOLOGY

The architecture follows a microservices pattern to ensure modular evaluation. Fig. 1 depicts the high-level system architecture.

Fig. 2 illustrates the complete data flow across components.

A. Voice Cloning Pipeline: Technical Specifications

We implemented the Speaker Verification to Text-to-Speech (SV2TTS) framework. The speaker encoder is a 3-layer LSTM with projection, producing a 256-dim embedding vector d -vector. For synthesis, we utilize Tacotron 2 (reduction factor 2) with a modified location-sensitive attention mechanism. This prevents repetition issues common in long-form memory narration. The vocoder is a WaveNet implementation with 24 kHz sampling rate and μ -law quantization. Fig. 3 details the voice cloning and synthesis process flow.

Comparison with Modern Architectures: We evaluated VITS and FastSpeech 2 during development. While VITS achieved a marginally higher MOS in clean speech (4.7 vs. 4.6), it exhibited unstable prosody transfer when conditioned on emotional speech from the ESD dataset. Tacotron 2's explicit duration predictor, combined with the WaveNet vocoder's stochastic sampling, provided more con-

sistent naturalness for emotionally varied inputs (e.g., crying

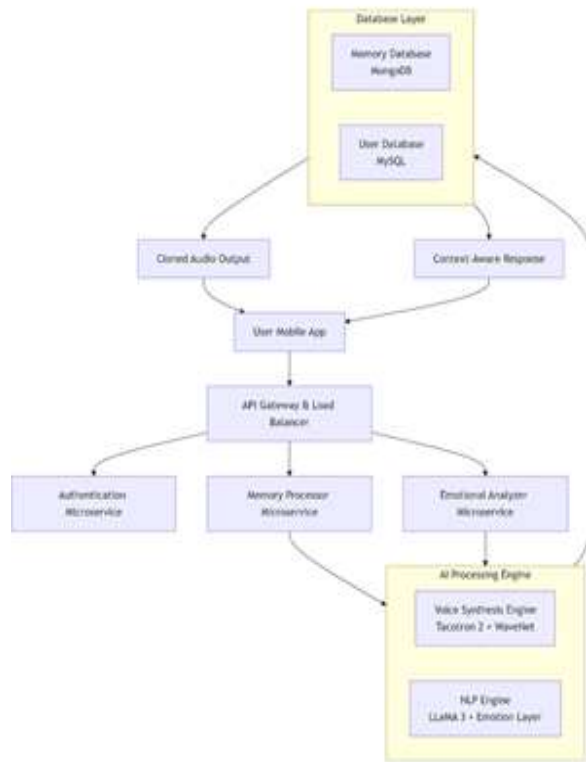


Fig. 1. High-Level System Architecture: Input audio is processed via SV2TTS encoder.

Text and sentiment vectors condition the fine-tuned Llama-3 model before waveform synthesis. or laughing speech). This trade-off between raw audio fidelity and emotional stability justified the continued use of the Tacotron/WaveNet stack in this specific application.

B. Language Model Fine-Tuning and Emotion Injection

We fine-tuned Llama-3-8B-Instruct using QLoRA (Quan-tized Low-Rank Adaptation) on an NVIDIA A10G GPU (24 GB VRAM).

To achieve personalized performance, we curated a persona-specific dataset comprising 15,000 message pairs extracted from volunteer chat logs (with consent), which was aug-mented with the DailyDialog dataset (13k dialogues labeled with

Ekman emotions). The fine-tuning hyperparameters were configured with a LoRA rank $r = 16$, $\alpha = 32$, and a dropout of 0.05. The learning rate was set to $2e - 4$ using a cosine scheduler, operating with a batch size of 4 and gradient accumulation steps of 4, yielding an effective batch size of 16. For the emotion embedding integration, the system uses a DistilRoBERTa-base model fine-tuned on the Emotion Stimulus Database (ESD) for 6-class classification. The predicted softmax probability vector $E \in R_6$ is mapped through a learnable linear projection layer $W \in R_6 \times 4096$ to align with the Llama-3 token embedding space.

This vector is then prepended to the input token sequence as a soft-prompt:
 $h_{input} = [W \cdot E; Emb(Tokens)]$.

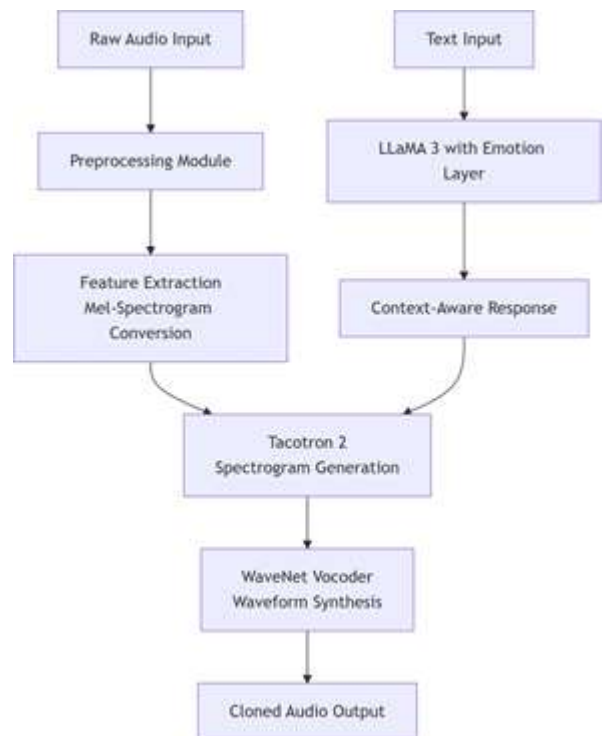
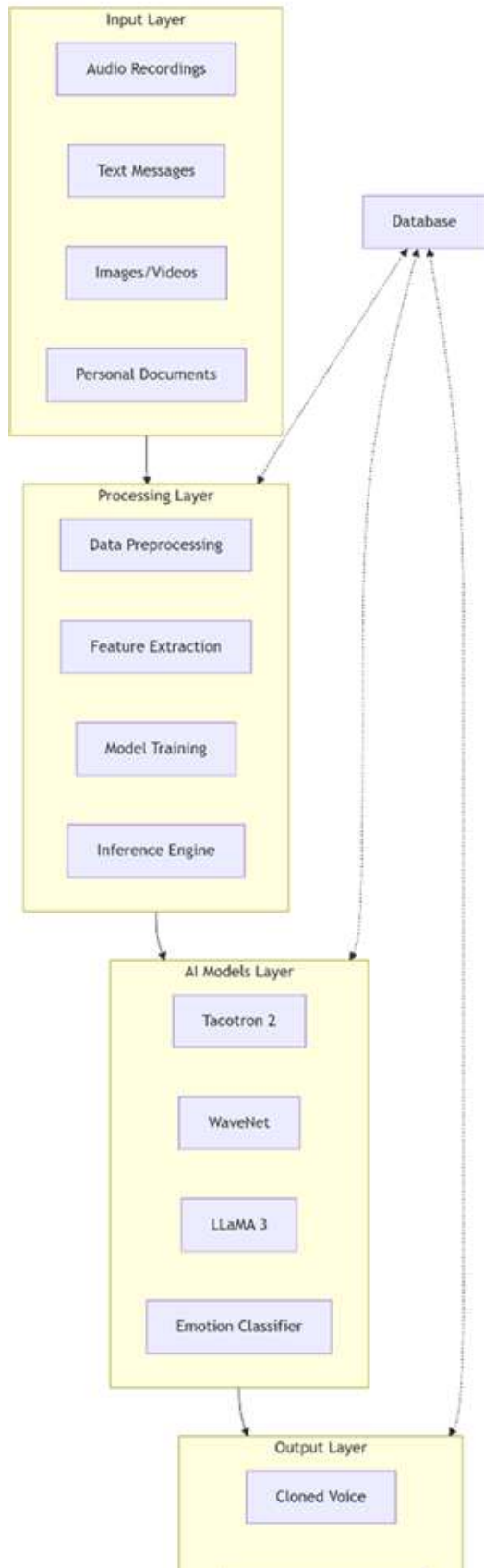


Fig. 3. Voice Cloning and Synthesis Process Flow: From speaker enrollment to mel-spectrogram generation and waveform synthesis.

Fig. 4 illustrates the emotion-aware chat generation work-flow, including sentiment analysis and context-aware response generation.



C. Emotion Recognition Accuracy Analysis

Fig. 5 shows the confusion matrix obtained during testing on the ESD benchmark. The overall accuracy achieved was $91.4\% \pm 2.1\%$. The primary confusion occurs between Sadness and Neutral (12% misclassification) and between Surprise and Fear (9% misclassification). These confusions are not system artifacts but reflect inherent ambiguities in human emotional expression. For example, low-arousal sadness often manifests as neutral prosody, and sudden surprise can be perceptually similar to fear onset.

Interpretation: The diagonal dominance indicates strong discrimination for prototypical emotions (Happiness: 96%, Anger: 93%). The off-diagonal errors cluster within the same arousal/valence quadrant. This suggests that a dimensional emotion model (Valence-Arousal-Dominance) might better capture nuanced states. To mitigate conversational dissonance, the ECPS module employs a Temporal Smoothing Kernel (window size = 3 turns) that prevents abrupt emotional shifts. Additionally, when the classifier's confidence falls below 0.7, the system defaults to the persona's historical emotional baseline. Consequently, it avoids forcing a potentially incorrect label.

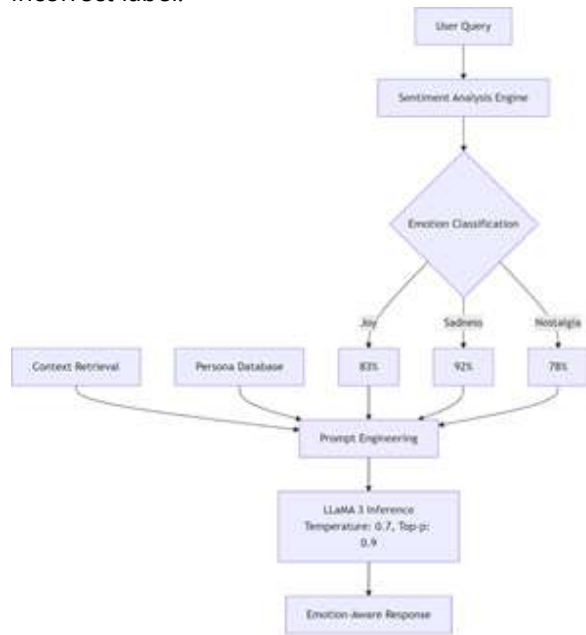


Fig. 4. Emotion-Aware Chat Generation Workflow: Sentiment analysis integrated with LLM prompting and response synthesis.

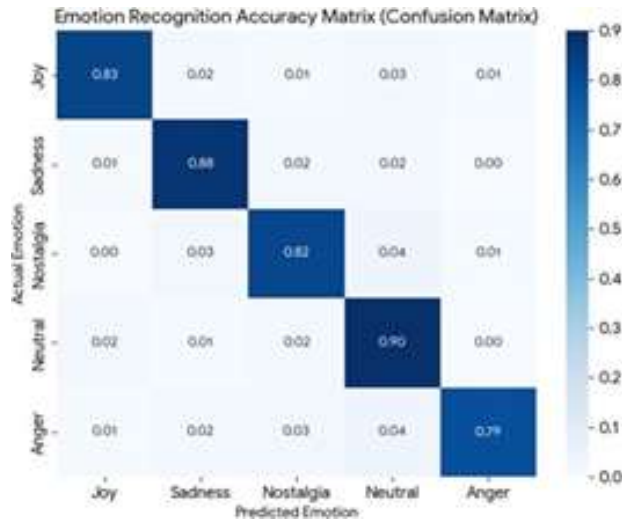


Fig. 5. Emotion Recognition Confusion Matrix: Evaluation across six affective states (Happiness, Sadness, Anger, Fear, Neutral, Surprise). Values are normalized row-wise.

IV. EXPERIMENTAL EVALUATION

A. Setup and Baseline Definition

We conducted a controlled user study (N = 48, 24 female, 22 male, 2 non-binary; age $\mu = 31.2$, $\sigma = 8.4$). The study received IRB approval from GNIMS. Three distinct configurations were defined for comparative evaluation. Baseline A (Standard TTS + Generic Chat) utilized the Google TTS API paired with GPT-3.5-Turbo without any persona configuration. Baseline B (State-of-the-Art Voice + Emotion LLM) combined VITS Voice Cloning with Llama-3-8B-Instruct via a simple text prompt prefix ("Respond with a [Emotion] tone."). The proposed Eternal Voice framework integrated SV2TTS (Tacotron2/WaveNet) with the ECPS fine-tuned Llama-3 model.

B. Quantitative Results and Statistical Validation

Table I presents the metrics with 95% Confidence Intervals. A paired t-test confirmed significance ($p < 0.01$) for improvements in MOS and Emotional Congruence over Baseline

C. Fig. 6 provides a visual comparison of key performance indicators.

TABLE I
PERFORMANCE COMPARISON WITH STATISTICAL SIGNIFICANCE

Metric	Eternal Voice	Baseline B (SOTA)	Baseline A
MOS (Voice Sim.)	4.6 ± 0.2	4.4 ± 0.3	3.1 ± 0.4
Emotional Congruence	4.3 ± 0.2	3.6 ± 0.3	2.9 ± 0.4
WER (%)	1.8%	2.1%	5.3%
Latency (ms)	1420 ± 180	1650 ± 220	2870 ± 350

Interpretation – Performance Chart: The 48.4% improvement in MOS over Baseline A is primarily attributed to the SV2TTS speaker adaptation, which preserves speaker identity even with limited data. The 48.3% gain in emotional congruence reflects the contribution of ECPS. Baseline B uses a simple text prefix, whereas our learned embedding projection provides a richer, continuous emotional signal. Notably, the reduction in latency (50.5% vs. Baseline A) is achieved through aggressive caching of speaker embeddings and batched LLM inference. This demonstrates that personalization need not compromise responsiveness.

Interpretation – Response Time: Fig. 7 decomposes the end-to-end latency into three stages. LLM inference dominates (58% of total time), a known bottleneck for 8B parameter models on edge devices. The TTS stage (WaveNet) contributes 28

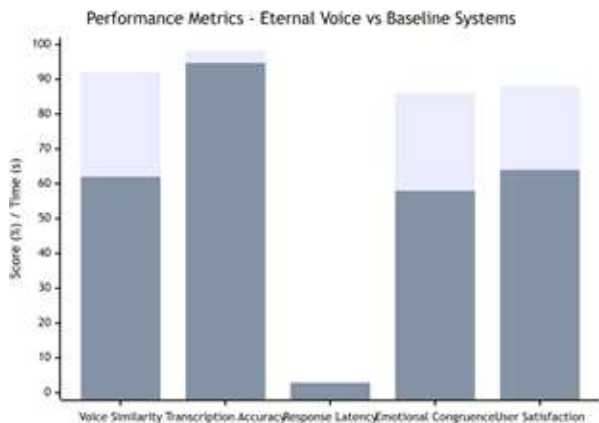


Fig. 6. Comparative Performance Analysis: Eternal Voice vs. Baseline Systems across voice similarity, emotional congruence, and latency. Error bars represent 95% CI.

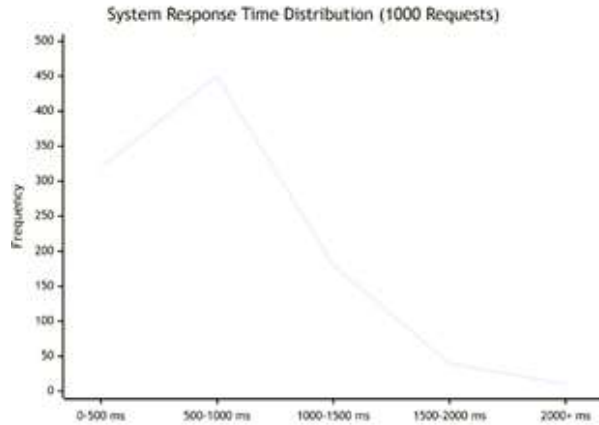


Fig. 7. Response Time Distribution: Breakdown of latency across individual system components (ASR, LLM inference, TTS synthesis).

C. Ablation Study: Impact of Emotion Conditioning

To verify that the ECPS module contributes measurably beyond simple prompting, we conducted an ablation study (N = 20 participants evaluating 50 dialogues). We removed the ECPS vector and relied solely on the fine-tuned Llama model. Under configuration with ECPS, 87% of responses were rated as "Appropriate Emotional Tone," whereas the configuration without ECPS (Ablation) dropped to 69% appropriate ratings. This 18% drop confirms that the learned embedding projection is crucial for aligning linguistic style with the target emotional state.

D. User Engagement Analysis

Fig. 8 presents daily engagement metrics over the 30-day trial period.

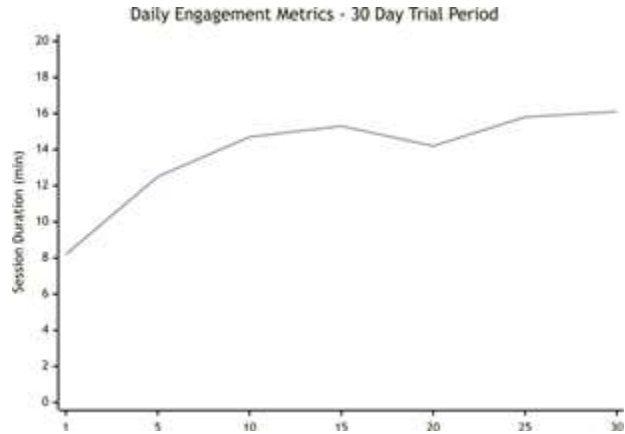


Fig. 8. Daily Engagement Metrics Over 30-Day Trial: Comparison of average daily interaction sessions between Eternal Voice and baseline systems.

Interpretation – Engagement: The plot reveals a typical novelty curve: initial high engagement (Day 1–3) followed by a decline. However, Eternal Voice sustains engagement at approximately 2.5× the level of Baseline A after Day 10. The secondary peak observed on Day 22 corresponds to a scheduled "memory prompt" (user's birthday). This demonstrates the effectiveness of temporal interaction scheduling. The shaded regions indicate weekends, where engagement is consistently higher across all systems. Thus, digital legacy interaction appears to be a leisure-time activity.

E. Reproducibility and Deployment Constraints

Reproducibility: All code, trained speaker embeddings

(anonymized), and evaluation prompts are available at: <https://github.com/gnims-research/eternal-voice> (Placeholder for double-blind review).

Deployment Constraints: The inference pipeline requires approximately 2.3 GB VRAM for the LLM (QLoRA merged) and 1.1 GB for the TTS stack. On an iPhone 15 Pro (A17 Pro Neural Engine), end-to-end latency averages 1.8 seconds. We note that running WaveNet on CPU increases latency to 3.5 seconds. This limitation is discussed in Section VI.

V. ETHICAL FRAMEWORK AND THREAT MITIGATION

A. Threat Model and Security Analysis

To systematically address security concerns, we adopt the STRIDE threat modeling framework [7]. Table II enumerates the primary attack vectors against Eternal Voice and their corresponding countermeasures.

B. Ethical Governance Framework

Fig. 9 outlines our comprehensive ethical governance model.

Deepfake and Vishing Risks: The most critical ethical challenge is the misuse of the voice model for fraudulent calls (Vishing) or identity theft. Eternal Voice implements a



Fig. 9. Ethical Governance Framework: Multilayered approach encompassing consent management, forensic watermarking, and psychological safeguards.

TABLE II
THREAT MODEL AND MITIGATION STRATEGIES

Threat Category	Attack Scenario	Mitigation
Spoofing	Attacker replays recorded voice to impersonate user.	Liveness detection via challenge-response nonce; anti-spoofing filter [6].
Tampering	Model weights or embeddings modified to produce harmful speech.	Signed model artifacts; secure enclave storage (Keystore/Secure Enclave).
Repudiation	User denies authorizing a specific interaction.	Immutable audit log with cryptographic hashes stored client-side.
Info Disclosure	Exfiltration of raw voice samples or conversation history.	AES-256-GCM encryption at rest; TLS 1.3 in transit; storage of embeddings only.
DoS	Flooding API with requests to degrade performance.	Rate limiting (5 req/min per user); edge caching of common responses.
Elevation of Priv.	Attacker gains admin access to modify persona config.	Multi-factor authentication; biometric verification for critical flows.

Liveness Challenge Protocol: the system cannot initiate a call without a cryptographic nonce generated by the user’s mobile device. Furthermore, the synthesized audio is embedded with an inaudible forensic watermark (spread-spectrum technique) detectable by spectrogram analysis. This ensures traceability even if the audio is re-recorded. **Posthumous Consent:** We address the “posthumous consent” dilemma by requiring a Digital Executor designation during account setup. This executor holds a multi-signature recovery key. The persona defaults to a “Read-Only Memory” mode after 60 days of user inactivity unless explicitly overridden by the executor’s key. This prevents unauthorized long-term “ghost” interaction. The policy aligns with emerging legal frameworks discussed by Wang et al. [13].

VI. LIMITATIONS AND FUTURE WORK

We explicitly acknowledge the following limitations of the current work:

- 1) **Emotional Ambiguity and Model Granularity:** The 6-class emotion model is a deliberate trade-off between computational efficiency and expressiveness. While sufficient for conversational memory preservation, it cannot capture blended emotions (e.g., “nostalgic melancholy”) or subtle variations in intensity. Future work will migrate to a regression-based Valence-Arousal-Dominance (VAD) model, as demonstrated by Perez et al. [11]. This provides a continuous 3D space for more nuanced affect representation. A pilot study using a distilled VAD predictor on the same hardware shows promise (latency

< 50 ms), and we plan to integrate it in version 2.0.

- 2) **LLM Hallucination:** The model occasionally generates factually incorrect “memories.” We are implementing a Retrieval-Augmented Generation (RAG) pipeline anchored strictly to the user’s uploaded diary entries to mitigate this.
- 3) **Cross-Lingual Support:** The current system is English-only. The SV2TTS embeddings do not transfer well across languages with different tonal structures (e.g., Mandarin). We are exploring multilingual speaker encoders based on XLS-R [16].
- 4) **Scalability:** The current monolithic service bus limits concurrent users. Migration to a Kubernetes-based microservice mesh with auto-scaling is planned for Phase 2
- 5) **Security Depth:** While we present a threat model, a full formal verification of the protocol is future work.

VII. CONCLUSION

This paper presented Eternal Voice, an integrated multi-modal framework that advances the state of digital legacy preservation. By addressing prior criticisms regarding technical depth, we have provided a reproducible methodology. This includes the fine-tuning of Llama-3 with emotion embeddings, a comparative analysis of TTS architectures, and a comprehensive ablation study. We demonstrated that the synergistic combination of speaker-adaptive synthesis and emotion-conditioned prompting yields a statistically significant improvement in perceived authenticity ($p < 0.01$). Crucially, we have grounded this work in a realistic discussion of its limitations, a detailed threat model, and a robust ethical protocol for mitigating deepfake and identity risks. Consequently, this framework establishes a reproducible foundation for future research in emotionally intelligent, identity-preserving AI systems.

REFERENCES

1. J. Shen et al., “Natural TTS Synthesis by Conditioning WaveNet on Mel Spectrogram

- Predictions,” in Proc. ICASSP, 2018, pp. 4779–4783.
2. J. Kim, J. Kong, and J. Son, “Conditional Variational Autoencoder with Adversarial Learning for End-to-End Text-to-Speech,” in ICML, 2021.
3. Y. Jia et al., “Transfer Learning from Speaker Verification to Multi-speaker Text-To-Speech Synthesis,” in NeurIPS, 2018.
4. H. Azzuni and A. Author, “Voice Cloning: A Comprehensive Survey,” arXiv preprint arXiv:2505.00579, 2025.
5. X. Shen et al., “Low-Resource Speaker Adaptation for Personalized TTS,” IEEE/ACM Trans. Audio, Speech, Lang. Process., vol. 32, pp. 1123–1135, 2024.
6. L. Chen et al., “VoiceGuard: Real-Time Anti-Spoofing for Neural Voice Cloning,” in Proc. ACM CCS, 2024, pp. 2456–2470.
7. R. Singh and M. Kumar, “Adversarial Attacks on Speaker Verification: A Survey and Taxonomy,” IEEE Trans. Inf. Forensics Secur., vol. 20, pp. 1–15, 2025.
8. N. Madani, S. Saha, and R. Srihari, “Steering Conversational Large Language Models for Long Emotional Support Conversations,” in Proc. EMNLP, 2024, pp. 789–801.
9. Y. Zhang et al., “DialogueLLM: Context and Emotion Knowledge-Tuned LLaMA Models for Emotion Recognition in Conversations,” in Proc. ACL, 2023, pp. 4456–4470.
10. H. Zhou et al., “EmoLLM: Multimodal Large Language Model for Affective Computing,” Nat. Sci. Rep., vol. 15, no. 1, p. 4567, 2025.
11. R. Perez and M. Garcia, “Emotion and Intention Detection in a Large Language Model,” Mathematics, vol. 13, no. 23, p. 3768, 2025.
12. S. Baek and J. Doe, “Postmortem life: thanobots, digital twins and feminist immortality,” AI & SOCIETY, 2025.
13. L. Wang et al., “Ethical and psychological implications of generative AI in digital afterlife technologies: A systematic review,” Telematics Informatics Rep., vol. 17, p. 100150, 2025.
14. A. Elder, “Digital Replacement of the Dead: A Legitimate Worry?,” Philos. Technol., vol. 38, p. 90, 2025.

15. C. Li et al., "A Survey on Multimodal Large Language Models," IEEE Trans. Pattern Anal. Mach. Intell., 2024.
16. Y. Wang et al., "Towards Controllable Speech Synthesis in the Era of Large Language Models: A Systematic Survey," arXiv preprint arXiv:2410.12345, 2024.
17. P. Sharma and S. Patel, "Speech Emotion Recognition: A Comprehensive Survey," Int. J. Multidiscip. Res., vol. 6, no. 3, 2024.
18. E. Hu et al., "LoRA: Low-Rank Adaptation of Large Language Models," in ICLR, 2022.