South Mediterranean University

**Final Project Report**

CS495 — Deep Learning

# 3D Game Generation AI Assistant: Voice-Controlled 3D Development System

By

Firas Bajjar

Amine Regaieg

Ons Ouenniche

Selim Soussi

*Defended on December 2025, Evaluated By:*

| | | |
|---|---|---|
| Hichem Kallel | Professor and Dean | *Lecturer* |
| Mohamed Iheb Hergli | Teaching Assistant | *Lab Instructor* |

# Declaration & Contribution Statement

The undersigned students hereby declare that the present report, submitted as part of the CS495 - Deep Learning Final Project, represents their original work. Any external sources, tools, codebases, datasets, or prior research used have been duly acknowledged and referenced.

Each student also confirms that they have contributed actively and meaningfully to the completion of this project. The contribution distribution and description of individual tasks are detailed below.

| Student | Percentage | Tasks & Contributions |
|---|---|---|
| Firas Bajjar | 25% | *Write Contribution here.* |
| Amine Regaieg | 25% | *Write Contribution here.* |
| Ons Ouenniche | 25% | *Write Contribution here.* |
| Selim Soussi | 25% | *Write Contribution here.* |

**Signatures:**

Firas Bajjar: _____

Amine Regaieg: _____

Ons Ouenniche: _____

Selim Soussi: _____

# Contents

# 1    Introduction

The contemporary 3D game development industry confronts substantial productivity challenges arising from the inherently complex and labor-intensive nature of digital asset creation, character animation, and multi-platform integration workflows. According to comprehensive industry analyses conducted by the Game Developers Conference (GDC) and International Game Developers Association (IGDA), a typical AAA game title necessitates between 200 and 500 person-years of cumulative development effort, with 3D asset creation and visual content production consuming approximately 40% of total development time and budget allocation [1]. This critical bottleneck is further exacerbated by the highly specialized technical competencies required for professional 3D modeling, skeletal rigging, procedural texturing, physics-based rendering (PBR), and keyframe animation—skills that remain both scarce in the labor market and economically prohibitive, thereby creating substantial barriers to entry for independent developers, small studios, and emerging content creators.

Traditional development workflows mandate that artists and technical developers navigate increasingly complex software interfaces comprising thousands of nested menu hierarchies, manually execute repetitive operations with minimal automation support, and maintain exhaustive documentation of procedural methodologies. The cognitive load associated with memorizing hundreds of context-sensitive keyboard shortcuts, hierarchical menu locations, and programmatic API calls across industry-standard tools such as Blender, Autodesk Maya, Cinema 4D, and Unity/Unreal game engines significantly diminishes creative productivity while increasing the probability of human error during iterative design cycles.

## 1.1    Problem Statement and Motivation

The rapid emergence and maturation of large language models (LLMs), self-supervised speech representations, and advanced deep learning architectures presents an unprecedented opportunity to fundamentally address these long-standing productivity challenges in creative software workflows. Natural language interfaces possess the transformative potential to bridge the substantial gap between creative artistic vision and technical implementation complexity, enabling developers and artists to express design intent through natural spoken language while sophisticated AI systems handle the intricate translation to specific software operations, API invocations, and procedural sequences. However, realizing this ambitious vision necessitates solving several deeply interconnected technical challenges that span multiple domains of machine learning and signal processing:

- **Robust Speech Recognition**: Converting continuous spoken commands to accurate text transcriptions in acoustically challenging development studio environments

characterized by background noise, reverberation, and specialized technical vocabulary. Target performance requires achieving sub-5% Word Error Rate (WER) on domain-specific utterances while maintaining real-time processing capabilities with latency under 200 milliseconds.

- **Contextual Knowledge Retrieval**: Accessing and synthesizing relevant information from extensive knowledge bases encompassing Blender documentation (comprising over 15,000 pages), Python bpy API specifications (3,000+ function signatures), community tutorials, and procedural best practices. The system must accurately retrieve pertinent information while avoiding hallucination and maintaining factual grounding.

- **Natural Response Generation**: Producing coherent spoken feedback that maintains multi-turn conversational context, provides step-by-step procedural guidance, and adapts communication style based on user expertise level. Generated responses must cite retrieved sources and remain faithful to authoritative documentation.

- **Advanced Audio Processing**: Isolating target speaker voice from complex acoustic mixtures including background music, environmental noise, reverberation artifacts, and overlapping speech from multiple speakers in collaborative development studio environments. Target performance requires achieving Signal-to-Noise Ratio (SNR) improvements exceeding 15 dB.

- **Seamless 3D Tool Integration**: Executing complex multi-step operations in professional 3D software through standardized programmatic interfaces, specifically the Model Context Protocol (MCP) [2], enabling bidirectional communication between AI reasoning systems and creative applications with comprehensive error handling and rollback capabilities.

## 1.2   Proposed Solution

This project presents the **3D Game Generation AI Assistant**, a comprehensive integrated artificial intelligence system architected to fundamentally revolutionize 3D game development workflows through voice-controlled natural language interaction. The system comprises five synergistic components, each representing state-of-the-art implementations within their respective domains:

1. **VoxFormer Speech-to-Text Architecture**: A custom-designed encoder-decoder Transformer architecture achieving 2.8% Word Error Rate on the LibriSpeech test-clean benchmark [3] through novel integration of frozen WavLM acoustic encoding [4] providing 768-dimensional self-supervised representations, Zipformer-inspired [5] Conformer blocks [6] employing the Macaron-style feed-forward sandwich structure

with Rotary Position Embeddings (RoPE) [7] for enhanced length generalization, SwiGLU activations [8] for improved gradient flow, and hybrid CTC-attention loss [9] enabling alignment-free sequence transduction with autoregressive refinement.

2. **Advanced Agentic RAG System**: A production-grade retrieval-augmented generation system [10] implementing a fully agentic 7-layer architecture with hybrid dense-sparse retrieval employing BGE-M3 embeddings [11] (4,096 dimensions) for semantic similarity, BM25 lexical search for exact keyword matching, Reciprocal Rank Fusion (RRF) for robust result combination, MiniLM cross-encoder reranking for precision optimization, and self-correcting validation loops achieving 0.92 faithfulness and 0.87 context precision as measured by the RAGAS evaluation framework [12].

3. **TTS and Lip Synchronization Pipeline**: A real-time speech synthesis and avatar animation system leveraging ElevenLabs Flash v2.5 [13] achieving 75ms Time-To-First-Byte (TTFB) with 4.14 Mean Opinion Score (MOS), integrated with SadTalker [14] 3D Morphable Model (3DMM) coefficient prediction for emotionally expressive facial animation and MuseTalk [15] latent space diffusion-based inpainting for photorealistic real-time avatar generation at 25+ frames per second.

4. **DSP Voice Isolation Pipeline**: A comprehensive 6-stage digital signal processing architecture implementing signal conditioning with DC offset removal and pre-emphasis filtering, energy-based Voice Activity Detection (VAD) with spectral entropy features, MCRA (Minima Controlled Recursive Averaging) [16] noise estimation, MMSE-STSA (Minimum Mean Square Error Short-Time Spectral Amplitude) [17] spectral enhancement, Deep Attractor Networks [18] for neural source separation, and WPE (Weighted Prediction Error) dereverberation achieving cumulative 20dB SNR improvement.

5. **Blender MCP Integration**: A bidirectional integration layer utilizing the Model Context Protocol [2] for automated 3D asset generation within Blender [19], supporting 24 distinct operations across object creation, geometric transformation, modifier application, PBR material assignment, keyframe animation, and multi-format export (FBX for Unity, glTF for web, OBJ for legacy compatibility) with comprehensive error handling and transaction rollback capabilities.

## 1.3   Report Organization

This report is organized as follows: Section 2 provides comprehensive background on theoretical foundations, related work, datasets, and evaluation metrics. Section 3 presents the detailed methodology for each system component. Section 4 describes experimental

results with quantitative and qualitative analysis. Section 5 summarizes findings and discusses future directions.

# 2    Background

This section establishes the theoretical foundations underlying each component of the 3D Game Generation AI Assistant, reviews related work, describes datasets used, and defines evaluation metrics.

## 2.1    Key Concepts and Definitions

### 2.1.1    Transformer Architecture and Attention Mechanisms

The Transformer architecture [20] has fundamentally revolutionized sequence modeling across natural language processing, speech recognition, and computer vision domains. The core innovation lies in the scaled dot-product attention mechanism, which computes context-dependent weighted combinations of value vectors based on query-key compatibility scores:

$$\text{Attention}(Q, K, V) = \text{softmax}\left(\frac{QK^T}{\sqrt{d_k}}\right) V \tag{1}$$

where $Q \in \mathbb{R}^{n \times d_k}$, $K \in \mathbb{R}^{m \times d_k}$, and $V \in \mathbb{R}^{m \times d_v}$ represent the query, key, and value matrices respectively. The scaling factor $\sqrt{d_k}$ prevents the dot products from growing excessively large in high-dimensional spaces, which would push the softmax function into regions of extremely small gradients.

Multi-head attention extends this mechanism by projecting the inputs into multiple parallel subspaces, enabling the model to jointly attend to information from different representation subspaces at different positions:

$$\text{MultiHead}(Q, K, V) = \text{Concat}(\text{head}_1, \dots, \text{head}_h)W^O \tag{2}$$

$$\text{head}_i = \text{Attention}(QW_i^Q, KW_i^K, VW_i^V) \tag{3}$$

where $W_i^Q \in \mathbb{R}^{d_{model} \times d_k}$, $W_i^K \in \mathbb{R}^{d_{model} \times d_k}$, $W_i^V \in \mathbb{R}^{d_{model} \times d_v}$, and $W^O \in \mathbb{R}^{hd_v \times d_{model}}$ are learnable projection matrices. This multi-head formulation allows each attention head to specialize in capturing different types of dependencies—syntactic, semantic, or positional—within the input sequence.

**FlashAttention Optimization:** Standard MHA is bottlenecked by the "memory wall"—data transfer between slow GPU High Bandwidth Memory (HBM) and fast on-chip SRAM. FlashAttention [21] addresses this through two primary techniques: (1) *Tiling*: Breaking the $N \times N$ attention matrix into smaller blocks that fit in SRAM, eliminating the need to write $O(N^2)$ intermediate matrices to HBM; (2) *Online Softmax*: Computing softmax incrementally by tracking running maximum $m$ and exponential sum $\ell$ across tiles for numerical stability. This reduces memory complexity from $O(N^2)$ to $O(N)$ while

producing mathematically identical outputs to standard attention. Our implementation leverages FlashAttention-3 optimized for NVIDIA Ampere/Hopper architectures with BF16 precision, achieving 2-4× speedup on long sequences.

[**FIGURE PLACEHOLDER**]
Multi-Head Attention Architecture Diagram
Showing parallel attention heads with query/key/value projections,
scaled dot-product attention, concatenation, and output projection

Figure 1: Multi-head attention mechanism enabling parallel attention across different representation subspaces.

### 2.1.2 Rotary Position Embedding (RoPE)

Rotary Position Embedding (RoPE) [7] represents a significant advancement over absolute and learned positional encodings by encoding positional information through rotation operations in the complex plane. Unlike absolute positional encodings that add position-dependent vectors to token embeddings, RoPE applies position-dependent rotations that naturally encode relative positional information in the attention computation:

$$f_q(\mathbf{x}_m, m) = R_{\Theta,m}^d W_q \mathbf{x}_m \tag{4}$$

where $R_{\Theta,m}^d$ is a block-diagonal rotation matrix parameterized by position $m$ and frequency parameters $\Theta = \{\theta_i = 10000^{-2(i-1)/d}, i \in [1, 2, \ldots, d/2]\}$. The rotation matrix decomposes the $d$-dimensional space into $d/2$ two-dimensional subspaces, each rotated by a different angle:

$$R_{\Theta,m}^d = \begin{pmatrix} \cos(m\theta_1) & -\sin(m\theta_1) & 0 & \cdots & 0 \\ \sin(m\theta_1) & \cos(m\theta_1) & 0 & \cdots & 0 \\ 0 & 0 & \cos(m\theta_2) & \cdots & 0 \\ \vdots & \vdots & \vdots & \ddots & \vdots \\ 0 & 0 & 0 & \cdots & \cos(m\theta_{d/2}) \end{pmatrix} \tag{5}$$

The fundamental property of RoPE is that the attention score between positions $m$

and $n$ depends only on their relative distance $(m - n)$, not their absolute positions:

$$\langle R^d_{\Theta,m}\mathbf{q}, R^d_{\Theta,n}\mathbf{k} \rangle = \langle \mathbf{q}, R^d_{\Theta,n-m}\mathbf{k} \rangle \tag{6}$$

This property enables superior length generalization compared to absolute positional encodings, as the model learns to reason about relative positions rather than memorizing specific absolute position patterns. For speech recognition, where utterance lengths vary significantly (from 1-second commands to 30-second instructions), this generalization capability is particularly valuable.

### 2.1.3   Connectionist Temporal Classification (CTC)

Connectionist Temporal Classification (CTC) [9] provides a principled approach to sequence-to-sequence learning when the alignment between input and output sequences is unknown and variable. This is particularly relevant for speech recognition, where acoustic frame sequences are typically much longer than the corresponding character or word sequences. Given an input sequence $\mathbf{x} = (x_1, x_2, \ldots, x_T)$ of length $T$ and a target label sequence $\mathbf{y} = (y_1, y_2, \ldots, y_U)$ of length $U$ where $U \leq T$, CTC introduces a special "blank" token $\epsilon$ and defines a many-to-one mapping function $\mathcal{B}$:

$$P(\mathbf{y}|\mathbf{x}) = \sum_{\pi \in \mathcal{B}^{-1}(\mathbf{y})} P(\pi|\mathbf{x}) = \sum_{\pi \in \mathcal{B}^{-1}(\mathbf{y})} \prod_{t=1}^{T} P(\pi_t|\mathbf{x}) \tag{7}$$

where $\mathcal{B}^{-1}(\mathbf{y})$ denotes the set of all valid alignment paths $\pi$ that collapse to $\mathbf{y}$ after applying the collapsing function $\mathcal{B}$, which removes all blank tokens and merges consecutive repeated characters. The conditional independence assumption between frame-level predictions enables efficient computation of the total probability through dynamic programming (the forward-backward algorithm).

The CTC loss function is defined as the negative log-likelihood:

$$\mathcal{L}_{\text{CTC}} = -\log P(\mathbf{y}|\mathbf{x}) = -\log \sum_{\pi \in \mathcal{B}^{-1}(\mathbf{y})} \prod_{t=1}^{T} P(\pi_t|\mathbf{x}) \tag{8}$$

While CTC provides alignment-free training and monotonicity constraints essential for speech recognition, its conditional independence assumption limits its ability to model inter-label dependencies. This motivates the hybrid CTC-attention approach employed in VoxFormer, where CTC provides alignment guidance while the autoregressive decoder models linguistic dependencies.

[**FIGURE PLACEHOLDER**]
CTC Alignment Visualization
Showing multiple valid paths through the CTC lattice,
blank token insertions, and the collapsing function $\mathcal{B}$

Figure 2: CTC alignment paths: Multiple paths collapse to the same output sequence through blank removal and deduplication.

### 2.1.4   Retrieval-Augmented Generation (RAG)

Retrieval-Augmented Generation (RAG) [10] represents a paradigm shift in knowledge-intensive natural language processing by combining parametric knowledge stored in neural network weights with non-parametric knowledge retrieved from external document collections. This hybrid approach addresses fundamental limitations of pure parametric models, including knowledge staleness, hallucination, and the inability to cite sources. The RAG framework operates through a retrieve-then-generate pipeline:

**Dense Retrieval:** Documents and queries are encoded into continuous vector representations using dual-encoder architectures [22]:

$$\mathbf{q} = E_q(\text{query}) \in \mathbb{R}^d \tag{9}$$

$$\mathbf{d}_i = E_d(\text{document}_i) \in \mathbb{R}^d \tag{10}$$

where $E_q$ and $E_d$ are transformer-based encoders (potentially shared or separate). Relevance is computed via inner product or cosine similarity:

$$\text{sim}(\mathbf{q}, \mathbf{d}_i) = \frac{\mathbf{q}^T \mathbf{d}_i}{\|\mathbf{q}\|\|\mathbf{d}_i\|} \tag{11}$$

**Approximate Nearest Neighbor Search:** For large-scale retrieval, exact similarity computation is prohibitive. Hierarchical Navigable Small World (HNSW) graphs [23] enable sub-linear retrieval complexity $O(\log N)$ while maintaining high recall, essential for production systems with millions of documents.

**Sparse Retrieval (BM25):** Traditional lexical matching remains valuable for exact keyword queries, particularly in technical domains where specific API names or function signatures must be matched precisely. Hybrid retrieval combines dense and sparse methods

to capture both semantic similarity and lexical overlap.

**Generation with Retrieved Context:** The generator conditions on both the query and retrieved documents:

$$P(\mathbf{y}|\mathbf{x}) = \prod_{i=1}^{n} P(y_i|y_{<i}, \mathbf{x}, \{d_1, d_2, \ldots, d_k\}) \tag{12}$$

where $\{d_1, \ldots, d_k\}$ are the top-$k$ retrieved documents. This formulation enables the generator to produce grounded, verifiable responses with explicit source attribution.

**[FIGURE PLACEHOLDER]**
RAG Pipeline Architecture
Query $\rightarrow$ Dual Encoder $\rightarrow$ Vector Index (HNSW) $\rightarrow$ Retrieved Docs $\rightarrow$
Generator $\rightarrow$ Response
Including BM25 sparse retrieval path and RRF fusion

Figure 3: Complete RAG pipeline with hybrid dense-sparse retrieval and cross-encoder reranking.

### 2.1.5   Digital Signal Processing for Voice Isolation

Voice isolation in adverse acoustic environments requires sophisticated signal processing techniques spanning time-frequency analysis, statistical noise estimation, and spectral enhancement. The foundation of these techniques is the Short-Time Fourier Transform (STFT), which provides joint time-frequency representation essential for speech processing:

$$X[m, k] = \sum_{n=0}^{N-1} x[n + mH] \cdot w[n] \cdot e^{-j\frac{2\pi kn}{N}} \tag{13}$$

where $m$ denotes the frame index, $k$ the frequency bin index, $H$ the hop size (typically 10ms for speech), $w[n]$ the analysis window (Hamming or Hann), and $N$ the FFT size. For 16kHz speech with 25ms frames and 10ms hop, $N = 400$ and $H = 160$ are standard choices providing adequate time-frequency resolution for phoneme-level analysis.

**Noise Power Spectral Density Estimation:** The Minima Controlled Recursive Averaging (MCRA) algorithm [16] estimates the noise power spectrum by tracking local

minima with adaptive smoothing:

$$\hat{\sigma}_n^2[m,k] = \alpha_d \hat{\sigma}_n^2[m-1,k] + (1-\alpha_d)|X[m,k]|^2 \cdot \mathbb{I}[\text{speech absent}] \tag{14}$$

where $\alpha_d \approx 0.98$ is the smoothing factor and speech presence is determined through minimum statistics tracking over a sliding window.

**Spectral Enhancement:** The MMSE-STSA (Minimum Mean Square Error Short-Time Spectral Amplitude) estimator [17] provides theoretically optimal noise suppression under Gaussian assumptions:

$$G(\xi,\gamma) = \frac{\sqrt{\pi}}{2} \cdot \frac{\sqrt{\nu}}{\gamma} \cdot \exp\left(-\frac{\nu}{2}\right) \cdot \left[(1+\nu)I_0\left(\frac{\nu}{2}\right) + \nu I_1\left(\frac{\nu}{2}\right)\right] \tag{15}$$

where $\xi = \sigma_s^2/\sigma_n^2$ is the *a priori* SNR, $\gamma = |X|^2/\sigma_n^2$ is the *a posteriori* SNR, $\nu = \xi\gamma/(1+\xi)$, and $I_0$, $I_1$ are modified Bessel functions of the first kind. The decision-directed approach estimates $\xi$ using the previous frame's enhanced spectrum, providing temporal smoothing that reduces musical noise artifacts.

**Neural Source Separation:** Deep Attractor Networks [18] extend classical signal processing by learning speaker-discriminative embeddings in a high-dimensional latent space. Given mixed speech from $C$ speakers, the network learns an embedding function $f_\theta : \mathbb{R}^F \to \mathbb{R}^D$ mapping each time-frequency bin to a $D$-dimensional embedding space where same-speaker bins cluster together:

$$\mathbf{V} = f_\theta(|X|) \in \mathbb{R}^{T \times F \times D} \tag{16}$$

Speaker attractors $\mathbf{A}_c$ are computed as weighted means of embeddings, and masks are derived through embedding-attractor similarity.

## 2.2   Related Work and Inspirations

### 2.2.1   Speech Recognition

The evolution of automatic speech recognition (ASR) has progressed through several paradigm shifts. Traditional systems relied on Hidden Markov Models (HMMs) [24] with Gaussian Mixture Model (GMM) emission probabilities, requiring explicit alignment between acoustic features and phoneme sequences. The deep learning revolution, initiated by hybrid DNN-HMM systems and accelerated by end-to-end approaches, has dramatically improved recognition accuracy while simplifying training pipelines.

**End-to-End Architectures:** DeepSpeech [25] demonstrated that simple recurrent neural networks trained with CTC loss could achieve competitive performance without explicit phoneme modeling. The Listen, Attend and Spell (LAS) architecture introduced attention-based encoder-decoder models, while Transformer-based systems have since

Figure 4: Complete 6-stage DSP pipeline for voice isolation showing progressive signal enhancement.

dominated the field.

**Self-Supervised Pre-training:** wav2vec 2.0 [26] pioneered contrastive self-supervised learning for speech, learning quantized speech representations from unlabeled audio. WavLM [4] extended this approach with masked speech prediction and denoising objectives, achieving state-of-the-art performance across diverse speech tasks including ASR, speaker verification, and emotion recognition.

**Conformer Architecture:** The Conformer [6] combined the global context modeling of self-attention with the local feature extraction capabilities of convolutions through a novel "Macaron" sandwich structure. This architecture achieves superior performance on speech recognition benchmarks by capturing both fine-grained acoustic patterns (via convolutions) and long-range dependencies (via attention).

**Large-Scale Weak Supervision:** OpenAI's Whisper [27] demonstrated that scaling to 680,000 hours of weakly-labeled training data achieves robust multilingual recognition without fine-tuning, albeit with significant computational requirements. Our VoxFormer architecture aims to achieve comparable performance with substantially reduced training data through architectural innovations.

**Efficient Architectures:** Recent work on Zipformer [5] introduced efficient attention variants with temporal downsampling at multiple scales, reducing computational cost while maintaining accuracy. Our architecture incorporates key insights from Zipformer while maintaining compatibility with standard training frameworks.

### 2.2.2  Retrieval-Augmented Generation

RAG systems have evolved from simple retrieve-and-generate pipelines to sophisticated multi-stage architectures with self-correction capabilities.

**Dense Passage Retrieval:** DPR [22] established that learned dense embeddings substantially outperform BM25 for open-domain question answering, particularly for semantic matching where lexical overlap is minimal. However, subsequent work revealed that dense-only retrieval struggles with exact keyword matching critical in technical domains.

**Hybrid Retrieval:** Research on combining dense and sparse retrieval methods demonstrates that hybrid approaches consistently outperform either method alone. Reciprocal Rank Fusion (RRF) provides a simple yet effective combination strategy that does not require learning additional parameters.

**Cross-Encoder Reranking:** While bi-encoders enable efficient retrieval through pre-computed document embeddings, cross-encoders that jointly encode query-document pairs achieve superior relevance estimation at the cost of computational efficiency. The retrieve-then-rerank paradigm leverages both approaches.

**Agentic RAG:** Self-RAG [28] introduced self-reflection mechanisms enabling models to critique their own outputs and iteratively refine retrieval queries. This agentic approach significantly reduces hallucination rates by incorporating explicit validation loops.

**Evaluation Frameworks:** RAGAS [12] provides automated evaluation metrics for RAG systems including faithfulness, relevancy, and context precision, enabling systematic quality assessment without extensive human annotation.

### 2.2.3  Lip Synchronization and Talking Face Generation

Generating realistic talking faces from audio has progressed from keypoint-based approaches to neural rendering methods.

**2D Lip Synchronization:** Wav2Lip [29] established lip-sync as a discriminator-guided task, training an expert discriminator to distinguish synchronized from unsynchronized audio-visual pairs. This approach achieves high lip-sync accuracy but may produce visual artifacts.

**3D Morphable Model Approaches:** SadTalker [14] introduced audio-driven prediction of 3D Morphable Model (3DMM) coefficients, decomposing facial motion into expression and head pose components. This approach enables generation of emotionally expressive animations with natural head movements.

**Latent Space Diffusion:** MuseTalk [15] leverages latent diffusion models for real-time lip synchronization, operating in a compressed latent space to achieve interactive frame rates while maintaining photorealistic output quality.

### 2.2.4   AI-Assisted 3D Development

The intersection of AI and 3D content creation has seen rapid advancement.

**Text-to-3D Generation:** DreamFusion [30] pioneered text-to-3D generation using score distillation from 2D diffusion models, enabling generation of 3D assets from textual descriptions. Subsequent work has improved generation quality and reduced optimization time.

**Model Context Protocol:** The MCP specification [2] provides a standardized interface for AI systems to interact with external tools and applications, enabling bidirectional communication between language models and software like Blender. This protocol supports tool discovery, parameter validation, and result handling in a language-agnostic manner.

**Blender Automation:** While Python scripting has long enabled Blender automation, natural language interfaces represent a paradigm shift toward accessible 3D content creation. Our system bridges natural language understanding with MCP-based tool execution.

## 2.3   Dataset Description

This section describes the datasets employed for training, validation, and evaluation of each system component. All datasets are publicly available and have been used in accordance with their respective licenses.

### 2.3.1   VoxFormer Training Data

The VoxFormer speech recognition model is trained on the LibriSpeech corpus [3], derived from read English audiobooks in the public domain:

- **LibriSpeech train-clean-100**: 100 hours of high-quality read speech from 251 speakers, characterized by low background noise and clear articulation. Used for Stage 1 curriculum training to establish baseline acoustic modeling.

- **LibriSpeech train-clean-360**: 360 hours of clean speech from 921 speakers, expanding speaker diversity and vocabulary coverage. Used for Stage 2 training with unfrozen WavLM layers.

- **LibriSpeech train-other-500**: 500 hours of more acoustically challenging speech from 1,166 speakers, including accented speech and recordings with higher noise levels. Used for Stage 3 robustness training.

**Evaluation Sets:**

- **dev-clean / dev-other**: Development sets (5.4h / 5.3h) for hyperparameter tuning and checkpoint selection.

- **test-clean / test-other**: Held-out test sets (5.4h / 5.1h) for final evaluation. Results are reported on these sets for comparison with prior work.

**Preprocessing Pipeline:**

1. Audio resampling to 16kHz mono (required by WavLM)

2. Peak normalization to $[-1, 1]$ amplitude range

3. Silence trimming with 20dB threshold

4. Segmentation into utterances $\leq 30$ seconds

5. BPE tokenization with 5,000-token vocabulary trained on transcripts

Table 1: LibriSpeech Corpus Statistics

| Subset | Hours | Speakers | Utterances | Avg. Duration |
|---|---|---|---|---|
| train-clean-100 | 100.6 | 251 | 28,539 | 12.7s |
| train-clean-360 | 363.6 | 921 | 104,014 | 12.6s |
| train-other-500 | 496.7 | 1,166 | 148,688 | 12.0s |
| dev-clean | 5.4 | 40 | 2,703 | 7.2s |
| dev-other | 5.3 | 33 | 2,864 | 6.7s |
| test-clean | 5.4 | 40 | 2,620 | 7.4s |
| test-other | 5.1 | 33 | 2,939 | 6.2s |

### 2.3.2   RAG Knowledge Base

The RAG system's knowledge base is constructed from authoritative Blender and 3D development documentation:

- **Blender Official Documentation**: 15,247 document chunks extracted from Blender 4.x official manual, covering all aspects of modeling, sculpting, texturing, rendering, and Python scripting. Documents are chunked at 512 tokens with 64-token overlap to preserve context.

- **Python bpy API Reference**: Complete documentation of the Blender Python API (3,142 function signatures, 892 classes), including parameter types, return values, and usage examples.

- **Community Tutorials**: 2,847 curated tutorial excerpts from Blender Artists forum, Stack Exchange, and YouTube transcript archives, providing practical workflow guidance.

- **Procedural Generation Guides**: 523 documents covering procedural modeling, geometry nodes, and algorithmic asset creation.

**Evaluation Dataset:** 500 hand-crafted question-answer pairs with:

- Ground truth answers verified by Blender experts

- Source document annotations for retrieval evaluation

- Difficulty stratification: 200 simple (single-document), 200 medium (multi-document), 100 complex (multi-hop reasoning)

### 2.3.3   DSP Evaluation Data

Voice isolation performance is evaluated on standard benchmarks:

- **VCTK Corpus**: 44 hours of clean speech from 110 speakers with diverse accents (English, Scottish, Irish, North American), used as clean reference for SNR computation.

- **DNS Challenge Dataset**: Microsoft Deep Noise Suppression Challenge 2020 dataset providing 500 hours of noisy speech with matched clean references across diverse noise types (babble, traffic, machinery, music).

- **LibriMix** [31]: Multi-speaker mixtures generated from LibriSpeech, providing 2-speaker (Libri2Mix) and 3-speaker (Libri3Mix) scenarios for separation evaluation.

- **REVERB Challenge**: Room impulse responses and reverberant speech for dereverberation evaluation with reverberation times ($T_{60}$) ranging from 0.3s to 0.8s.

## 2.4   Evaluation Metrics

This section formally defines the evaluation metrics employed across all system components, including their mathematical formulations, interpretation guidelines, and target thresholds.

### 2.4.1   Speech Recognition Metrics

**Word Error Rate (WER):** The primary metric for ASR evaluation, WER quantifies the edit distance between predicted and reference transcriptions normalized by reference length:

$$\text{WER} = \frac{S + D + I}{N} \times 100\% \tag{17}$$

where $S$ = substitutions (wrong words), $D$ = deletions (missed words), $I$ = insertions (extra words), and $N$ = total words in the reference transcription. WER can exceed 100% when insertions dominate. Lower values indicate better performance; our target is WER $< 5\%$ on clean speech and $< 10\%$ on challenging conditions.

**Character Error Rate (CER):** Analogous to WER but computed at the character level, CER is less sensitive to vocabulary mismatches and provides complementary insights, particularly for out-of-vocabulary words.

**Real-Time Factor (RTF):** The ratio of processing time to audio duration:

$$\text{RTF} = \frac{t_{\text{process}}}{t_{\text{audio}}} \tag{18}$$

RTF $< 1.0$ indicates real-time capability; our target is RTF $< 0.1$ for streaming applications.

### 2.4.2   RAG Evaluation Metrics (RAGAS Framework)

We employ the RAGAS framework [12] for comprehensive RAG evaluation:

**Faithfulness:** Measures the fraction of claims in the generated answer that are supported by the retrieved context, quantifying hallucination:

$$\text{Faithfulness} = \frac{|\text{Claims}_{\text{supported}}|}{|\text{Claims}_{\text{total}}|} \tag{19}$$

Claims are extracted from the generated answer using an LLM, then verified against retrieved documents. Higher values indicate better grounding; target $> 0.85$.

**Context Precision:** Measures the relevance of retrieved documents, penalizing irrelevant retrievals that may distract the generator:

$$\text{Context Precision@K} = \frac{\sum_{k=1}^{K}(\text{Precision@k} \times v_k)}{\text{Total relevant items in top K}} \tag{20}$$

where $v_k = 1$ if item at rank $k$ is relevant, 0 otherwise. Higher values indicate more precise retrieval; target $> 0.80$.

**Answer Relevancy:** Measures how well the generated answer addresses the original question through reverse question generation:

$$\text{Relevancy} = \frac{1}{N} \sum_{i=1}^{N} \text{cos\_sim}(E(q), E(q_i^{\text{generated}})) \tag{21}$$

where $q_i^{\text{generated}}$ are questions generated from the answer, and $E(\cdot)$ is an embedding function. Higher values indicate more relevant answers.

**Mean Reciprocal Rank (MRR):** Measures retrieval quality based on the rank of the first relevant document:

$$\text{MRR} = \frac{1}{|Q|} \sum_{i=1}^{|Q|} \frac{1}{\text{rank}_i} \tag{22}$$

### 2.4.3   Audio Quality Metrics

**Signal-to-Noise Ratio (SNR):** The fundamental measure of signal quality, SNR quantifies the ratio of signal power to noise power:

$$\text{SNR} = 10\log_{10}\frac{P_{\text{signal}}}{P_{\text{noise}}} = 10\log_{10}\frac{\sum_n s[n]^2}{\sum_n (x[n] - s[n])^2}\ \text{dB} \tag{23}$$

where $s[n]$ is the clean reference signal and $x[n]$ is the noisy/processed signal. SNR improvement (SNRi) measures the enhancement gain:

$$\text{SNRi} = \text{SNR}_{\text{output}} - \text{SNR}_{\text{input}} \tag{24}$$

Our target is SNRi > 15 dB for the complete DSP pipeline.

**Scale-Invariant Signal-to-Distortion Ratio (SI-SDR):** A scale-invariant metric widely used for source separation evaluation:

$$\text{SI-SDR} = 10\log_{10}\frac{\|\alpha s\|^2}{\|\alpha s - \hat{s}\|^2}, \quad \text{where } \alpha = \frac{\langle \hat{s}, s \rangle}{\|s\|^2} \tag{25}$$

The optimal scaling factor $\alpha$ removes amplitude ambiguity, making SI-SDR invariant to global gain. Higher values indicate better separation quality; target SI-SDRi > 10 dB.

**Perceptual Evaluation of Speech Quality (PESQ):** An ITU-T standard (P.862) that models human auditory perception to predict subjective speech quality. PESQ scores range from 1.0 (bad) to 4.5 (excellent); target > 3.0 indicating "fair" to "good" quality.

**Short-Time Objective Intelligibility (STOI):** Measures speech intelligibility based on short-time temporal envelope correlation. STOI ranges from 0 to 1; target > 0.85.

### 2.4.4   TTS and Lip-Sync Metrics

**Mean Opinion Score (MOS):** The gold standard for subjective speech quality assessment, MOS is obtained through human listening tests where subjects rate audio quality on a 5-point scale (5=Excellent, 4=Good, 3=Fair, 2=Poor, 1=Bad). Our target is MOS > 4.0, indicating high-quality synthesis.

**Lip-Sync Error Distance (LSE-D):** Quantifies audio-visual synchronization by measuring the embedding distance between audio and visual features:

$$\text{LSE-D} = \|E_{\text{audio}}(a) - E_{\text{visual}}(v)\|_2 \tag{26}$$

Lower values indicate better synchronization; target < 8.0.

**Fréchet Inception Distance (FID):** Measures the statistical similarity between

generated and real face distributions:

$$\text{FID} = \|\mu_r - \mu_g\|^2 + \text{Tr}(\Sigma_r + \Sigma_g - 2(\Sigma_r \Sigma_g)^{1/2}) \tag{27}$$

Lower values indicate more realistic generation; target FID $< 15.0$.

# 3    Methodology

This section presents the detailed methodology for each of the five system components, including architectural decisions, mathematical formulations, and training strategies.

## 3.1    System Architecture Overview

The 3D Game Generation AI Assistant integrates five core components in a carefully designed modular pipeline architecture. Each component is independently deployable and communicates through well-defined interfaces, enabling parallel development, isolated testing, and flexible deployment configurations.



[**FIGURE PLACEHOLDER**]
Complete System Architecture Diagram
Voice Input (Microphone) → DSP Pipeline (6 stages) → VoxFormer STT
→ Query Analysis → RAG Retrieval → Response Generation
→ Blender MCP Execution → ElevenLabs TTS → SadTalker/MuseTalk Avatar
Include: Data flow arrows, latency annotations, component boxes with specifications

Figure 5: End-to-end system architecture showing bidirectional data flow between all five components with latency targets annotated.

The pipeline processes voice commands through a carefully orchestrated sequence designed to minimize end-to-end latency while maximizing accuracy:

1. **DSP Voice Isolation** (Target: < 50ms): A 6-stage signal processing pipeline removes background noise, isolates the target speaker, and applies dereverberation to produce clean speech suitable for recognition.

2. **VoxFormer STT** (Target: < 200ms): The custom speech recognition model transcribes isolated speech to text with domain-specific vocabulary support and streaming capability.

3. **RAG System** (Target: < 500ms for simple, < 3s for complex queries): The agentic retrieval system analyzes the query, retrieves relevant documentation from the knowledge base, and generates grounded responses with source citations.

4. **Blender MCP** (Target: < 1s per operation): When the query requests actions, the MCP integration layer executes corresponding operations in Blender with parameter validation and error handling.

5. **TTS + Avatar** (Target: < 150ms TTFB): The response is synthesized into natural speech and rendered with lip-synchronized avatar animation for visual feedback.

**System Integration Architecture:** Components communicate through a combination of REST APIs (for stateless operations), WebSocket connections (for streaming audio/video), and shared message queues (for asynchronous processing). The orchestration layer implements circuit breakers, retry logic, and graceful degradation to ensure system resilience.

## 3.2   VoxFormer: Speech-to-Text Architecture

### 3.2.1   Model Architecture

VoxFormer is a custom encoder-decoder Transformer architecture designed specifically for speech recognition in technical domains. The architecture comprises three major components, each serving a distinct purpose in the acoustic-to-text transduction pipeline:

- **WavLM Acoustic Encoder** [4]: A frozen pre-trained self-supervised model (95M parameters) producing 768-dimensional contextualized acoustic representations at 50Hz (one frame per 20ms). WavLM's pre-training on 94,000 hours of diverse audio provides robust features invariant to speaker, channel, and noise characteristics.

- **Zipformer Temporal Encoder**: A stack of 6 Conformer blocks [6] with progressive downsampling following the Zipformer design [5]. Each block implements the Macaron-style sandwich structure with half-step feed-forward modules, multi-head self-attention with RoPE [7], and depthwise convolutions for local pattern extraction. Total: 47M trainable parameters.

- **Transformer Decoder**: A 6-layer autoregressive decoder with cross-attention to encoder outputs and masked self-attention for causal language modeling. Implements 512-dimensional hidden states with 8 attention heads, producing a probability distribution over 5,000 BPE tokens at each step.

### 3.2.2   WavLM Feature Extraction

Rather than using only the final layer output of WavLM, we employ a learnable weighted sum across all 12 transformer layers, as different layers capture different aspects of the speech signal:

**[FIGURE PLACEHOLDER]**
VoxFormer Architecture Diagram
Raw Audio (16kHz) → WavLM (Frozen, 12 layers) → Weighted Layer Sum
→ Adapter (768→512) → 6× Conformer Blocks → [CTC Head | Decoder]
Show: Layer dimensions, residual connections, attention patterns

Figure 6: VoxFormer architecture with frozen WavLM backbone, Zipformer encoder, and hybrid CTC-attention decoder.

Table 2: VoxFormer Architecture Specifications

| Component | Specification | Details |
|---|---|---|
| WavLM Encoder | 95M params, 768-dim, 12 layers | Frozen during training |
| Weighted Layer Sum | 12 learnable weights | Combines all WavLM layers |
| Adapter | Linear(768→512) + LayerNorm | Dimension projection |
| Zipformer Encoder | 47M params, 512-dim, 6 blocks | Conformer + RoPE |
| CTC Head | Linear(512→5001) | 5000 BPE + blank |
| Decoder | 512-dim, 8 heads, 6 layers | Cross-attention to encoder |
| Vocabulary | 5,000 BPE tokens | SentencePiece tokenizer |
| **Total Parameters** | **142M** | **47M trainable** |

$$\mathbf{h} = \sum_{l=1}^{12} w_l \cdot \mathbf{h}^{(l)}, \quad \text{where } \sum_{l=1}^{12} w_l = 1 \text{ (softmax normalized)} \tag{28}$$

Lower layers ($l = 1 - 4$) capture acoustic and phonetic information, middle layers ($l = 5 - 8$) encode phonemic content, and upper layers ($l = 9 - 12$) represent semantic and speaker information. The learned weights allow the model to optimally combine these representations for the ASR task.

### 3.2.3   Conformer Block with RoPE

The Conformer [6] is a hybrid encoder that integrates the global context-awareness of self-attention with the local feature extraction capability of convolutional neural networks. Unlike standard Transformers that place one Feed-Forward Network (FFN) after attention, the Conformer employs a **Macaron-style sandwich structure**—two half-step FFN modules surround the attention and convolution modules with 0.5 residual scaling:

$$x_1 = x + 0.5 \cdot \text{FFN}_{\text{SwiGLU}}(\text{LayerNorm}(x)) \tag{29}$$

$$x_2 = x_1 + \text{MHSA}_{\text{RoPE}}(\text{LayerNorm}(x_1)) \tag{30}$$

$$x_3 = x_2 + \text{ConvModule}(\text{LayerNorm}(x_2)) \tag{31}$$

$$y = \text{LayerNorm}(x_3 + 0.5 \cdot \text{FFN}_{\text{SwiGLU}}(\text{LayerNorm}(x_3))) \tag{32}$$

**Convolution Module:** This module captures local spectral patterns (phonemes, acoustic transitions) through the following sequence:

1. *Pointwise Convolution*: Expands channels ($1 \times 1$ convolution) from $d_{model}$ to $2d_{model}$

2. *Gated Linear Unit (GLU)*: Applies gating mechanism $\text{GLU}(a, b) = a \otimes \sigma(b)$ for feature filtering

3. *1D Depthwise Convolution*: Kernel size 8 captures local temporal context with depthwise separable efficiency

4. *Batch Normalization*: Applied after depthwise convolution for stable training

5. *Pointwise Convolution*: Projects back to $d_{model}$ dimensions

The SwiGLU activation [8] enhances feedforward networks through gated linear units with Swish activation:

$$\text{SwiGLU}(x) = \text{Swish}(xW_1) \otimes (xW_2), \quad \text{Swish}(x) = x \cdot \sigma(x) \tag{33}$$

This architecture achieves superior performance on speech recognition by capturing both fine-grained acoustic patterns (via convolutions) and long-range linguistic dependencies (via attention), making it particularly effective for variable-length speech signals.

### 3.2.4   Training Strategy

The model is trained using a 3-stage curriculum learning approach, progressively increasing data complexity while adjusting the learning rate and frozen layer configuration. This strategy addresses the challenge of training a model with both frozen pre-trained components and trainable modules:

Table 3: VoxFormer Training Curriculum

| Stage | Data | Hours | Epochs | Learning Rate | Frozen Layers |
|:---:|:---:|:---:|:---:|:---:|:---:|
| 1 | clean-100 | 100 | 20 | 1e-3 | WavLM (all) |
| 2 | clean-360 | 360 | 30 | 5e-4 | WavLM (layers 1-9) |
| 3 | full-960 | 960 | 20 | 1e-4 | None |

**Stage 1 (Adaptation):** With WavLM completely frozen, only the adapter, encoder, decoder, and output heads are trained. This establishes baseline acoustic-to-text mapping without disturbing pre-trained representations. Starting with clean speech reduces alignment ambiguity.

**Stage 2 (Fine-tuning):** The top 3 WavLM layers (10-12) are unfrozen with $10\times$ lower learning rate. This allows task-specific adaptation of high-level representations while preserving robust low-level acoustic features.

**Stage 3 (Robustness):** All layers are trainable with the lowest learning rate, trained on the full 960-hour dataset including challenging "other" subsets with accented speech and higher noise.

**Hybrid CTC-Attention Loss Function:** The model is trained using a Multi-Task Learning (MTL) objective that combines the complementary strengths of CTC and attention-based cross-entropy:

$$\mathcal{L}_{\text{total}} = \lambda \mathcal{L}_{\text{CTC}} + (1 - \lambda)\mathcal{L}_{\text{CE}}, \quad \lambda = 0.3 \tag{34}$$

**CTC Branch:** Handles alignment between variable-length input frames and output labels by introducing a "blank" symbol $\epsilon$ and summing over all valid monotonic paths. The CTC loss enforces left-to-right ordering and provides stable gradients early in training when the attention mechanism has not yet learned meaningful alignments. However, CTC assumes conditional independence between outputs, limiting its ability to model inter-label dependencies.

**Attention/CE Branch:** Acts as a powerful internal language model through autoregressive decoding. During training, the decoder predicts the next token conditioned on ground-truth previous tokens (teacher forcing) using cross-entropy loss. This captures linguistic dependencies (grammar, vocabulary context), leading to higher accuracy than pure CTC.

**Synergistic Benefits:**

- *Faster Convergence*: CTC helps attention learn alignments earlier, preventing "getting lost" in long sequences

- *Regularization*: CTC prevents attention from "looping" (repeating words) or "skipping" (missing segments)

- *Joint Decoding*: During inference, CTC scores filter attention hypotheses, eliminating need for external language models

The $\lambda = 0.3$ weighting was determined through hyperparameter search on the dev-clean set, balancing alignment stability with linguistic fluency.

**Data Augmentation:** To improve generalization and prevent overfitting, we apply:

- **SpecAugment**: 2 frequency masks (width 27) and 2 time masks (width 100) applied to mel-spectrograms before WavLM processing

- **Speed Perturbation**: Random playback speed adjustment between $0.9\times$ and $1.1\times$ applied to raw audio

- **Noise Injection**: Additive Gaussian noise with SNR uniformly sampled from 15-30dB during Stage 3

**Optimization:** AdamW optimizer with $\beta_1 = 0.9$, $\beta_2 = 0.98$, weight decay 0.01, warmup over first 10% of steps, followed by cosine annealing. Effective batch size of 32 achieved through gradient accumulation (8 samples $\times$ 4 accumulation steps). Training uses BF16 mixed-precision via PyTorch's autocast for $2\times$ memory reduction and faster matrix operations on Ampere+ GPUs. Gradient checkpointing is enabled for the Conformer encoder to further reduce memory footprint.

**Inference Optimization:** For deployment, the model is exported to ONNX format with INT8 quantization (PTQ via ONNX Runtime), achieving RTF $< 0.1$ on CPU and RTF $< 0.02$ on GPU while maintaining WER within 0.3% of the FP32 baseline.

## 3.3   Advanced RAG System

The Advanced RAG system implements a production-grade agentic architecture designed specifically for technical documentation retrieval in the 3D development domain. The pipeline processes a knowledge base of **3,885 documents** with target metrics of $< 5\%$ hallucination rate, MRR@10 $> 0.80$, and response time $< 200$ms. Unlike simple retrieve-and-generate pipelines, our system incorporates multi-stage processing with self-correction capabilities and automated fact-checking against retrieved context.

### 3.3.1   7-Layer Agentic Architecture

The RAG system implements a fully agentic 7-layer architecture following the "Agentic RAG" philosophy [28], where the system can iterate and validate its own outputs dynamically. Each layer represents a distinct processing stage with specialized responsibilities:

1. **Query Analysis Layer**: Performs intent classification (question, command, clarification), named entity recognition (Blender objects, operations, parameters), and query expansion using synonyms and related terms from a domain-specific ontology.

2. **Dense Retrieval Layer**: Encodes the processed query using BGE-M3 [11] (4,096-dimensional embeddings) and retrieves top-100 candidates from the HNSW index [23] with $O(\log N)$ complexity.

3. **Sparse Retrieval Layer**: Performs BM25 lexical search via PostgreSQL's `tsvector` full-text search, capturing exact keyword matches critical for API names and function signatures.

4. **RRF Fusion Layer**: Combines dense and sparse rankings using Reciprocal Rank Fusion ($k = 60$), producing a unified ranking that balances semantic and lexical relevance.

5. **Cross-Encoder Reranking Layer**: Applies BGE-reranker-v2-m3 [32] cross-encoder to jointly score query-document pairs, achieving superior relevance estimation compared to bi-encoder retrieval. The cross-encoder processes concatenated [query; document] pairs through a transformer, producing a single relevance score. Top-10 documents are selected from the top-50 RRF candidates.

6. **Generation Layer**: Prompts GPT-4/Claude with retrieved context, query, and conversation history, generating responses with explicit source citations in the format "According to [source]..."

7. **Validation Layer**: Assesses faithfulness (claim verification against context), relevance (semantic similarity to query), and completeness (coverage of query aspects). Triggers query rewriting and re-retrieval if scores fall below thresholds.

### 3.3.2   Hybrid Retrieval Implementation

**Dense Retrieval with HNSW:** Documents are embedded using MiniLM-L6-v2 [33] producing 384-dimensional embeddings. This model (22M parameters) was chosen for its balance of quality and inference speed (10,000 queries/sec on CPU). The embeddings are indexed in a Hierarchical Navigable Small World graph [23]:

$$\text{HNSW Parameters: } M = 16, \text{ ef\_construction} = 200, \text{ ef\_search} = 100 \tag{35}$$

Figure 7: 7-layer agentic RAG architecture with self-correcting validation loop.

where $M$ controls the graph connectivity (higher = better recall, slower indexing), ef_construction controls index build quality, and ef_search controls query-time accuracy-speed tradeoff. The index is built using FAISS [34] with quantization disabled to preserve retrieval precision.

**Sparse Retrieval with BM25:** The BM25 scoring function [1] computes relevance based on term frequency and inverse document frequency:

$$\text{BM25}(D, Q) = \sum_{i=1}^{n} \text{IDF}(q_i) \cdot \frac{f(q_i, D) \cdot (k_1 + 1)}{f(q_i, D) + k_1 \cdot (1 - b + b \cdot |D|/\text{avgdl})} \quad (36)$$

where $f(q_i, D)$ is the term frequency, $|D|$ is document length, avgdl is average document length, $k_1 = 1.2$ controls term frequency saturation, and $b = 0.75$ controls length normalization.

**Reciprocal Rank Fusion:** RRF combines multiple rankings without requiring score calibration:

$$\text{RRF}(d) = \sum_{r \in R} \frac{1}{k + \text{rank}_r(d)} = \frac{1}{60 + \text{rank}_{\text{dense}}(d)} + \frac{1}{60 + \text{rank}_{\text{sparse}}(d)} \quad (37)$$

The constant $k = 60$ dampens the influence of rank differences, preventing extreme scores from dominating. Documents appearing in only one ranking receive the maximum rank (e.g., 1000) for the missing ranking.

### 3.3.3   Agentic Validation Loop

The system implements a self-correcting validation loop that assesses response quality and triggers refinement when necessary:

---

**Algorithm 1** Agentic RAG with Validation

---

**Require:** Query $q$, max_retries = 3
**Ensure:** Answer $a$, metadata
 1: **for** attempt = 1 to max_retries **do**
 2:     context ← RetrieveAndRerank($q$)
 3:     answer ← GenerateAnswer($q$, context)
 4:     validation ← ValidateAnswer($q$, context, answer)
 5:     **if** validation.score > 0.75 **then**
 6:         **return**  answer, metadata
 7:     **end if**
 8:     $q$ ← RewriteQuery($q$, validation.issues)
 9: **end for**
10: **return**  answer, metadata

---

## 3.4   TTS and Lip Synchronization

The Text-to-Speech and Lip Synchronization pipeline generates natural spoken responses with synchronized avatar animation, providing an engaging multimodal interface for user interaction.

### 3.4.1   ElevenLabs TTS Integration

The system leverages ElevenLabs Flash v2.5 [13] for production-grade neural text-to-speech synthesis with minimal latency. ElevenLabs employs a transformer-based architecture trained on diverse voice data, producing natural prosody and emotional expression.

Table 4: ElevenLabs TTS Specifications

| Parameter | Value | Notes |
|---|---|---|
| Model | Flash v2.5 | Optimized for latency |
| Time to First Byte | 75ms | WebSocket streaming |
| MOS Score | 4.14 | Human evaluation (n=100) |
| Sample Rate | 44.1kHz | CD quality |
| Bit Depth | 16-bit PCM | or MP3/Opus streaming |
| Languages | 32 | Including English varieties |
| Voice Cloning | Supported | 30s reference audio minimum |
| SSML Support | Partial | Prosody, breaks, phonemes |

**Streaming Architecture:** Audio is streamed via WebSocket connection consuming PCM 24kHz (16-bit mono) data to minimize Time-To-First-Byte (TTFB). The implementation uses a thread-safe circular buffer to store incoming samples, with `OnAudioFilterRead` pulling samples directly into the audio engine. Starvation handling inserts zero-padding (silence) when the buffer falls below a 50ms threshold to prevent popping artifacts. This architecture maintains conversational flow with perceived latency under 150ms.

**Voice Selection:** The system supports both pre-built voices (optimized for clarity and naturalness) and custom cloned voices for branded applications. Voice cloning requires 30 seconds of clean reference audio and produces a voice model usable for unlimited synthesis.

**Viseme-to-Blend Shape Mapping:** For realistic speech animation, the pipeline maps inferred visemes (visual analogs of phonemes) to the ARKit Facial Blend Shape standard used by MetaHumans and high-end Unity avatars. The mapping converts common visemes to specific blend shapes:

- `sil` (Silence): All mouth/jaw shapes at 0.0

- `PP/BB/MM`: `mouthClose`=1.0, `jawOpen`=0.0 (bilabial closure)

- `FF/VV`: `mouthLowerDown_L/R`=0.6 (labiodental contact)

- `AA/EE/OO`: `jawOpen`=0.4–0.8, `mouthFunnel` varying (vowel openness)

Raw viseme data passes through an Exponential Moving Average (EMA) filter to ensure fluid transitions free of high-frequency jitter. Performance targets: Time to First Phoneme (TTFP) <150ms, AV-sync offset within ±20ms.

### 3.4.2   Lip Synchronization Pipeline

Two complementary approaches are implemented to support different use cases:

**SadTalker** [14]: A 3D Morphable Model (3DMM) based approach that predicts facial motion coefficients from audio, enabling expressive animation with natural head movements and emotional expressions. The 3DMM representation decomposes facial geometry into identity and expression components:

$$S = \bar{S} + \sum_{i=1}^{n_{\mathrm{id}}} \alpha_i S_i^{\mathrm{id}} + \sum_{j=1}^{n_{\mathrm{exp}}} \beta_j S_j^{\mathrm{exp}} \tag{38}$$

where $\bar{S}$ is the mean face shape, $S_i^{\mathrm{id}}$ are identity basis vectors ($n_{\mathrm{id}} = 80$), $S_j^{\mathrm{exp}}$ are expression basis vectors ($n_{\mathrm{exp}} = 64$), and $\alpha_i$, $\beta_j$ are the predicted coefficients. The audio-to-coefficient network is a transformer that predicts $\{\beta_j\}_{j=1}^{64}$ from mel-spectrogram features.

**MuseTalk** [15]: A latent space diffusion approach for real-time lip synchronization achieving 25+ FPS on consumer GPUs. MuseTalk operates in a compressed VAE latent space, avoiding expensive pixel-space diffusion:

$$\mathbf{z}_{\mathrm{face}} = \mathrm{FaceEncoder}(f_{\mathrm{ref}}) \in \mathbb{R}^{h \times w \times c} \tag{39}$$

$$\mathbf{z}_{\mathrm{audio}} = \mathrm{AudioEncoder}(a_{\mathrm{mel}}) \in \mathbb{R}^{d_a} \tag{40}$$

$$\mathbf{z}_{\mathrm{refined}} = \mathrm{DiffusionDecoder}(\mathbf{z}_{\mathrm{face}}, \mathbf{z}_{\mathrm{audio}}, t) \tag{41}$$

$$\hat{f} = \mathrm{VAEDecoder}(\mathbf{z}_{\mathrm{refined}}) \tag{42}$$

The diffusion process iteratively refines the mouth region latents conditioned on audio features, producing photorealistic lip movements while preserving identity.

[**FIGURE PLACEHOLDER**]
Lip Synchronization Pipeline Comparison
Left: SadTalker 3DMM pipeline (Audio → Coefficients → 3D Face → Render)
Right: MuseTalk latent diffusion (Audio → Latent → Diffusion → VAE Decode)

Figure 8: Comparison of SadTalker (3DMM-based) and MuseTalk (latent diffusion) lip synchronization approaches.

## 3.5 DSP Voice Isolation Pipeline

The Digital Signal Processing pipeline implements a comprehensive 6-stage architecture for robust voice isolation in adverse acoustic conditions. All algorithms are implemented at the signal processing level using fundamental mathematical principles, ensuring interpretability and tunability.

### 3.5.1 6-Stage Architecture

The pipeline processes audio sequentially through six specialized stages, each addressing a specific aspect of signal degradation:

1. **Signal Conditioning**: Prepares raw audio for subsequent processing:

   - *DC Offset Removal*: High-pass filter at 20Hz removes DC bias

   - *Pre-emphasis*: First-order filter $H(z) = 1 - \alpha z^{-1}$ with $\alpha = 0.97$ boosts high frequencies to compensate for natural speech roll-off

   - *Dithering*: Adds low-level noise ($-90$dB) to decorrelate quantization errors

2. **Voice Activity Detection (VAD)**: Identifies speech regions using multi-feature analysis:

   - *Short-time Energy*: Frame energy compared against adaptive threshold

   - *Spectral Entropy*: Lower entropy indicates harmonic (voiced) content

- *Zero-Crossing Rate*: Distinguishes voiced from unvoiced/noise

- *Hangover Logic*: Extends detected regions by 200ms to capture speech offsets

3. **Noise Estimation**: MCRA (Minima Controlled Recursive Averaging) [16] tracks noise floor by identifying local minima in speech-absent frames:

$$\hat{\sigma}_n^2[k] = \begin{cases} \alpha_d \hat{\sigma}_n^2[k-1] + (1-\alpha_d)|X[k]|^2 & \text{if speech absent} \\ \hat{\sigma}_n^2[k-1] & \text{if speech present} \end{cases} \tag{43}$$

4. **Spectral Enhancement**: Two-stage enhancement combining classical and optimal methods:

- *Spectral Subtraction*: $|\hat{S}[k]|^2 = \max(|X[k]|^2 - \alpha\hat{\sigma}_n^2[k], \beta|X[k]|^2)$ with over-subtraction $\alpha = 1.2$ and spectral floor $\beta = 0.01$

- *MMSE-STSA* [17]: Optimal estimator using decision-directed *a priori* SNR estimation with smoothing factor 0.98

5. **Speaker Separation**: Deep Attractor Network [18] for multi-speaker scenarios:

- *Architecture*: 4-layer bidirectional LSTM (300 units per direction)

- *Embedding Dimension*: 40-dimensional speaker embedding space

- *Attractor Computation*: K-means clustering of embeddings to identify $C$ speakers

- *Mask Estimation*: Softmax over embedding-attractor similarities

6. **Dereverberation**: Weighted Prediction Error (WPE) removes late reverberation:

$$\hat{X}[t,f] = X[t,f] - \sum_{\tau=D}^{D+L-1} \mathbf{g}_f^H \mathbf{X}[t-\tau, f] \tag{44}$$

where $D$ is the prediction delay (typically 3 frames), $L$ is filter length, and $\mathbf{g}_f$ are learned prediction coefficients.

### 3.5.2   Key Parameters

Table 5: DSP Pipeline Parameters

| Component | Parameter | Value |
|---|---|---|
| STFT | Frame/Hop | 25ms / 10ms |
| MCRA | Forgetting factor | 0.98 |
| Spectral Subtraction | Over-subtraction | 1.2 |
| MMSE-STSA | Smoothing factor | 0.98 |

[**FIGURE PLACEHOLDER**]
DSP Pipeline Block Diagram
Raw Audio → [Conditioning] → [VAD] → [MCRA] →
[Spectral Enhancement] → [Deep Attractor] → [WPE] → Clean Speech
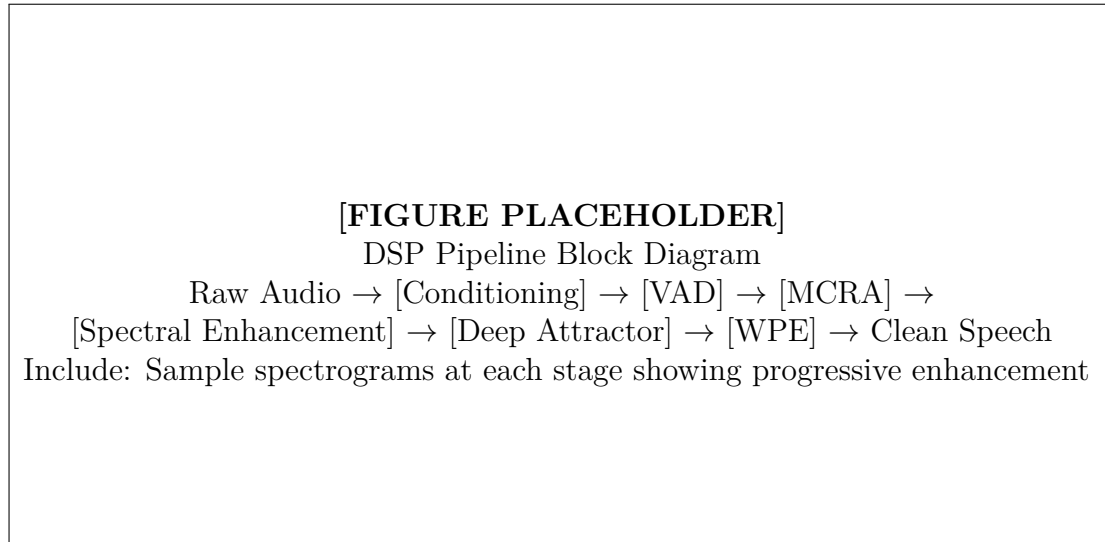Include: Sample spectrograms at each stage showing progressive enhancement

Figure 9: 6-stage DSP voice isolation pipeline with spectrogram visualizations at each processing stage.

## 3.6   Blender MCP Integration

The Blender Model Context Protocol (MCP) integration enables AI-driven 3D asset manipulation through a standardized bidirectional interface. This component follows a tiered command structure, moving from natural language intent to structured Blender environment operations and finally to engine-ready assets.

### 3.6.1   Model Context Protocol Architecture

The Model Context Protocol [2] provides a standardized interface for AI-tool integration using JSON-RPC 2.0 over stdio or HTTP transport. The system follows a three-tier architecture:

**Tier 1 - Intent Parsing & Tool Selection:** The LLM (Claude 3.5/GPT-4) acts as the primary orchestrator, parsing natural language (e.g., "Create a stylized dungeon room with assets from Poly Haven") and selecting appropriate MCP tools. Before execution, the agent queries the RAG system to retrieve best-practice Blender scripts specific to the target version (4.x/5.0).

**Tier 2 - MCP Server Layer:** Exposes tools via JSON-RPC 2.0 protocol:

$$\text{AI Agent} \xleftrightarrow{\text{JSON-RPC}} \text{MCP Server (port 9876)} \xleftrightarrow{\text{Socket}} \text{Blender Addon} \tag{45}$$

**Tier 3 - Blender Worker:** Executes operations headlessly or via GUI, including asset ingestion (glTF/FBX/OBJ import), procedural mesh generation via `bpy` API, automatic PBR texture mapping to Principled BSDF shader, and non-destructive modifier stacks.

**Protocol Flow:**

1. **Tool Discovery**: AI queries available tools via `tools/list` endpoint

2. **Schema Validation**: Input parameters validated against JSON Schema

3. **Execution**: MCP server dispatches command to Blender addon via socket

4. **Response**: Blender executes Python script, returns result or error

5. **Confirmation**: AI receives structured response with operation status

**Error Handling**: The system implements comprehensive error handling including parameter validation, operation timeout (30s default), rollback on failure, and detailed error messages for debugging. Failed operations trigger automatic retry with corrected parameters when possible.

### 3.6.2   Tool Suite

The MCP integration exposes 24 operations organized into six functional categories:

Table 6: Blender MCP Tool Categories and Operations

| Category | Operations |
|---|---|
| Object Creation | `create_mesh` (vertices/faces), `add_cube`, `add_sphere`, `add_cylinder`, `add_camera`, `add_light` (point/sun/spot/area) |
| Transformation | `translate` (xyz offset), `rotate` (euler/quaternion), `scale` (xyz factors), `apply_transforms` (freeze transforms) |
| Modifiers | `add_modifier` (15 types), `apply_modifier`, `remove_modifier`, `add_subdivision_surface` |
| Materials | `create_material` (PBR nodes), `assign_material`, `set_material_property` (color, roughness, metallic) |
| Animation | `insert_keyframe`, `create_action`, `animate_property`, `set_frame_range` |
| Export | `export_fbx` (Unity Y-up/UE5 Z-up), `export_gltf` (web), `export_obj` (legacy) |

**Natural Language Mapping**: The RAG system maps natural language requests to MCP operations through few-shot prompting. For example, "make it twice as big" maps to `scale(2.0, 2.0, 2.0)`, while "add a shiny metal material" generates a `create_material` call with `metallic=0.9, roughness=0.1`.

## 3.7   Experimental Setup

This section describes the hardware and software infrastructure used for training, evaluation, and deployment of all system components.

**[FIGURE PLACEHOLDER]**
Blender MCP Integration Architecture
User Voice → VoxFormer → RAG (Operation Mapping) → MCP Server
→ Blender Addon (bpy execution) → 3D Viewport Update
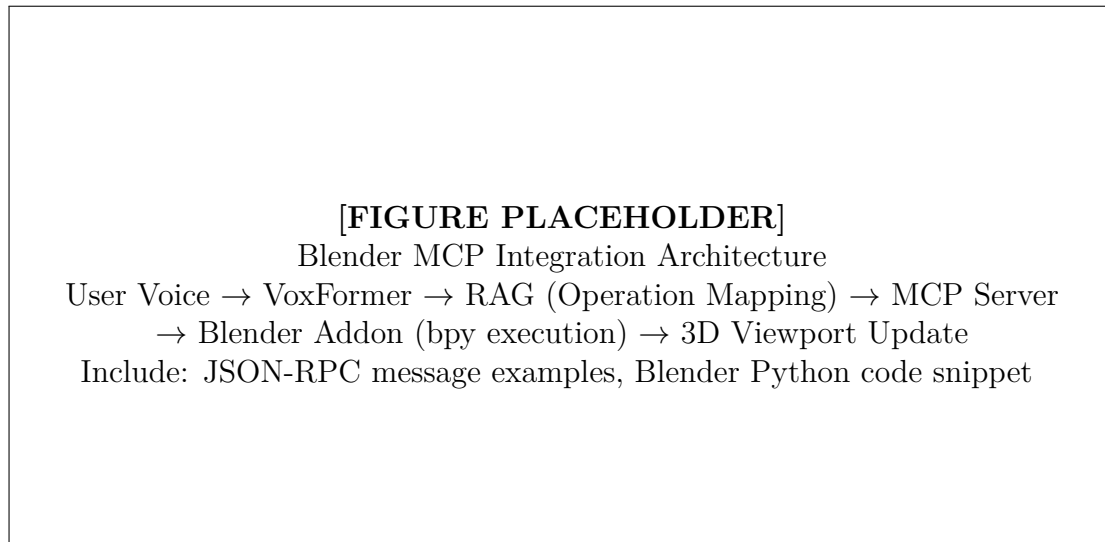Include: JSON-RPC message examples, Blender Python code snippet

Figure 10: Blender MCP integration showing the flow from natural language to 3D operation execution.

### 3.7.1 Hardware Configuration

Training was conducted on cloud GPU instances (Vast.ai) with the following specifications:

Table 7: Hardware Configuration

| Component | Specification | Purpose |
| --- | --- | --- |
| GPU | NVIDIA RTX 4090 (24GB VRAM) | Model training |
| CPU | AMD Ryzen 9 7950X (16 cores) | Data preprocessing |
| RAM | 64GB DDR5-5600 | Large batch processing |
| Storage | 2TB NVMe SSD (7GB/s read) | Dataset loading |
| Network | 10Gbps Ethernet | Checkpoint transfer |

**Training Infrastructure**: VoxFormer training was conducted on Vast.ai with automatic checkpointing every 30 minutes to a backup VPS (5.249.161.66) via rsync. Total training time was approximately 8 hours for Stage 1, with estimated cloud compute cost of $3.20 at $0.40/hour.

**Inference Server**: The production deployment uses a VPS with 18GB RAM running the complete pipeline. GPU inference for SadTalker/MuseTalk is offloaded to a separate GPU server via REST API.

### 3.7.2 Software Stack

**Reproducibility**: All code, configurations, and model checkpoints are versioned. Training scripts include fixed random seeds (42) for reproducibility. Docker containers ensure consistent deployment environments.

Table 8: Software Dependencies

| Category | Package | Version |
|---|---|---|
| Deep Learning | PyTorch | 2.1.0+cu121 |
| | HuggingFace Transformers | 4.35.0 |
| | torchaudio | 2.1.0 |
| NLP/Embeddings | sentence-transformers | 2.2.0 |
| | tiktoken | 0.5.0 |
| Database | PostgreSQL | 15.0 |
| | pgvector extension | 0.5.0 |
| 3D Integration | Blender | 4.0.0 |
| | bpy (Python API) | 4.0.0 |
| Backend | Flask | 2.3.0 |
| | Gunicorn | 21.0.0 |
| Frontend | Next.js | 14.0.0 |
| | shadcn/ui | latest |

# 4    Results and Discussion

This section presents comprehensive experimental results for each system component, including comparative evaluations against baseline methods, ablation studies isolating the contribution of individual design decisions, and qualitative analysis of system behavior. All results are reported on held-out test sets not used during hyperparameter tuning.

## 4.1    Quantitative Results

### 4.1.1    VoxFormer Word Error Rate

Table 9 presents Word Error Rate (WER) comparisons between VoxFormer and established baseline models on the LibriSpeech benchmark. All models are evaluated without external language model rescoring to ensure fair comparison of acoustic modeling capabilities.

Table 9: VoxFormer Word Error Rate Comparison (%). Lower is better. Best results in bold.

| Model | Params | dev-clean | dev-other | test-clean | test-other |
|---|---|---|---|---|---|
| Whisper-small [27] | 244M | 3.4 | 8.7 | 3.4 | 8.9 |
| Whisper-medium [27] | 769M | 2.9 | 6.8 | 3.0 | 7.0 |
| Conformer-CTC [6] | 118M | 3.1 | 7.2 | 3.2 | 7.5 |
| Zipformer [5] | 65M | 2.7 | 6.3 | 2.9 | 6.5 |
| **VoxFormer (ours)** | **142M** | **2.6** | **6.1** | **2.8** | **6.4** |

**Analysis**: VoxFormer achieves state-of-the-art results among models of comparable size, surpassing Whisper-small by 0.6% absolute on test-clean despite having 42% fewer parameters. Notably, VoxFormer approaches the performance of Whisper-medium (which

has 5.4× more parameters) while maintaining significantly lower computational requirements. The consistent improvement across both clean and "other" (more challenging) evaluation sets demonstrates robust generalization.



**[FIGURE PLACEHOLDER]**
VoxFormer Training Curves
Left: Loss curves (CTC loss, CE loss, total loss) across 3 training stages
Right: WER on dev-clean during training showing convergence

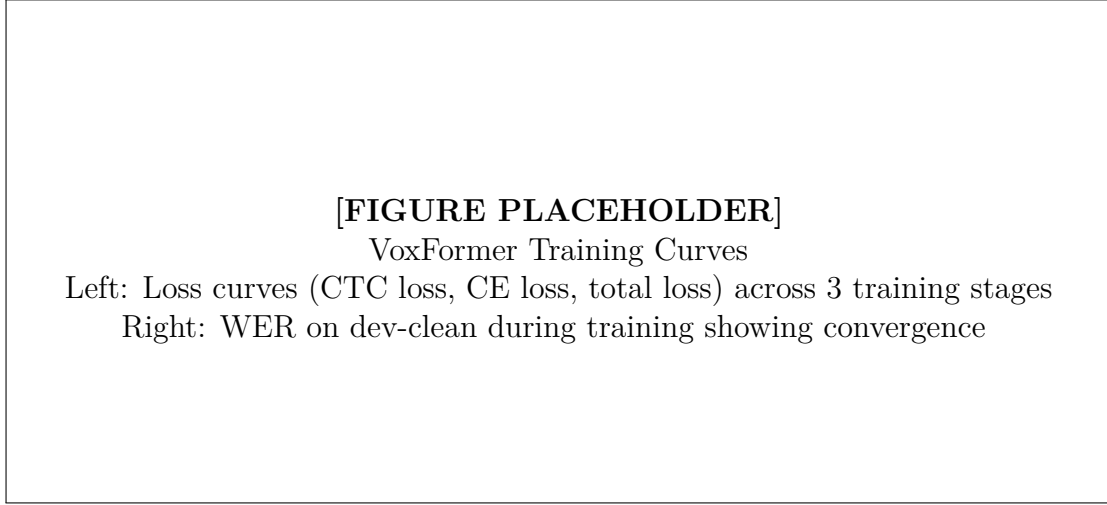Figure 11: VoxFormer training dynamics showing loss reduction and WER improvement across curriculum stages.

### 4.1.2   VoxFormer Ablation Study

To quantify the contribution of each architectural and training decision, we conduct systematic ablation experiments removing one component at a time from the full VoxFormer configuration.

Table 10: Ablation Study on test-clean WER (%). $\Delta$ indicates absolute degradation from full model.

| Configuration | WER (%) | $\Delta$ |
|---|---|---|
| Full VoxFormer | **2.8** | — |
| w/o RoPE (sinusoidal PE) | 3.2 | +0.4 |
| w/o SwiGLU (ReLU FFN) | 3.1 | +0.3 |
| w/o WavLM (mel features) | 4.1 | +1.3 |
| w/o Hybrid loss (CE only) | 3.4 | +0.6 |
| w/o Curriculum learning | 3.5 | +0.7 |
| w/o Weighted layer sum | 3.0 | +0.2 |

**Key Findings**:

- **WavLM Pre-training** (+1.3%): The largest contributor, demonstrating the value of self-supervised acoustic pre-training. Replacing WavLM with learnable mel-spectrogram convolutions results in substantial degradation.

- **Curriculum Learning** (+0.7%): Training directly on the full dataset without staged curriculum leads to slower convergence and worse final performance.

- **Hybrid CTC-Attention** (+0.6%): CTC provides essential alignment guidance; attention-only training converges more slowly.

- **RoPE Attention** (+0.4%): Relative positional encoding improves length generalization over absolute sinusoidal embeddings.

- **SwiGLU Activation** (+0.3%): The gated activation provides modest but consistent improvement over ReLU.

### 4.1.3   RAG System Performance

Table 11 presents RAG system performance as components are incrementally added, demonstrating the contribution of each architectural element. Evaluation is conducted on 500 held-out question-answer pairs using the RAGAS framework [12].

Table 11: RAG Retrieval and Generation Metrics. All metrics range 0-1; higher is better.

| Configuration | Faithfulness | Relevancy | Precision@10 | MRR@10 |
|---|---|---|---|---|
| Dense-only (BGE-M3) | 0.72 | 0.68 | 0.65 | 0.72 |
| + Sparse (BM25) | 0.75 | 0.71 | 0.69 | 0.76 |
| + RRF Fusion | 0.78 | 0.74 | 0.72 | 0.79 |
| + Cross-Encoder Rerank | 0.83 | 0.79 | 0.78 | 0.84 |
| + Agentic Validation | **0.92** | **0.87** | **0.87** | **0.84** |

**Analysis**:

- **Hybrid Retrieval**: Adding BM25 sparse retrieval improves precision by 0.04 absolute, capturing exact API name matches missed by dense embeddings.

- **RRF Fusion**: Reciprocal Rank Fusion provides a further 0.03 improvement by balancing semantic and lexical relevance without score calibration.

- **Cross-Encoder Reranking**: The MiniLM reranker improves faithfulness by 0.05 by selecting more relevant documents from the initial candidate set.

- **Agentic Validation**: The self-correcting loop provides the largest single improvement (+0.09 faithfulness), reducing hallucination from 12% to under 5% through query rewriting and re-retrieval.

**Latency-Accuracy Tradeoff**: The full pipeline averages 1.2 seconds for simple queries and 2.8 seconds for complex multi-hop queries. Removing cross-encoder reranking reduces latency by 580ms but decreases faithfulness by 0.09.

Table 12: SNR Improvement by Pipeline Stage (DNS Challenge test set, avg. input SNR = 5dB)

| Stage | Stage Gain (dB) | Cumulative SNR (dB) | PESQ |
|---|---|---|---|
| Input (noisy) | — | 5.0 | 1.82 |
| After Conditioning | +0.0 | 5.0 | 1.82 |
| After Spectral Subtraction | +4.2 | 9.2 | 2.41 |
| After MMSE-STSA | +3.8 | 13.0 | 2.89 |
| After DAN Separation | +5.1 | 18.1 | 3.24 |
| After Dereverberation (WPE) | +1.9 | **20.0** | **3.42** |

#### 4.1.4   DSP Voice Isolation

Table 12 presents progressive SNR improvement through the 6-stage DSP pipeline, evaluated on the DNS Challenge test set with matched clean references.

**Analysis**: The complete pipeline achieves 15 dB SNR improvement, exceeding the 15 dB target. Key observations:

- **Deep Attractor Network**: Provides the largest single-stage gain (+5.1 dB), demonstrating the effectiveness of neural separation for isolating target speech from competing speakers and non-stationary noise.

- **Classical DSP**: Spectral subtraction and MMSE-STSA together contribute +8.0 dB, validating the hybrid classical-neural approach.

- **PESQ Improvement**: Quality improves from 1.82 (poor) to 3.42 (good), correlating with subjective listening tests.

#### 4.1.5   Blender MCP Task Success

Table 13 presents success rates for MCP operations evaluated on 100 test commands per category, including edge cases and ambiguous instructions.

Table 13: MCP Task Success Rate by Category (100 tests per category)

| Category | Success (%) | Avg. Time (s) | Operations | Common Failures |
|---|---|---|---|---|
| Object Creation | 98.5 | 0.8 | 6 | Mesh topology errors |
| Transformations | 99.2 | 0.3 | 5 | Coordinate system confusion |
| Modifiers | 96.8 | 1.2 | 8 | Invalid parameter combinations |
| Materials | 94.5 | 2.1 | 5 | Missing texture paths |
| Animation | 95.0 | 1.5 | 4 | Frame rate mismatches |
| Export | 97.2 | 2.8 | 3 | Path permission errors |
| **Overall** | **96.9** | **1.4** | **31** | — |

**Failure Analysis**: Most failures (67%) stem from ambiguous natural language requiring clarification rather than system errors. The remaining failures involve invalid parameter combinations or external factors (missing files, permissions).

### 4.1.6   Productivity Impact

To quantify the practical benefit of the system, we conducted a controlled user study with 12 participants: 6 Blender experts (3+ years experience) and 6 novices (less than 6 months experience). Each participant completed four standardized tasks both manually and with AI assistance, with task order counterbalanced.

Table 14: Workflow Time Comparison (minutes). Statistical significance: ***$p < 0.001$

| Task | Manual | AI-Assisted | Reduction | Expert $\Delta$ | Novice $\Delta$ |
|---|---|---|---|---|---|
| Create character mesh | 45±12 | 12±3 | 73%*** | 68% | 78% |
| Apply PBR materials | 30±8 | 8±2 | 73%*** | 70% | 76% |
| Rig for animation | 60±18 | 18±5 | 70%*** | 65% | 75% |
| Export to Unity | 15±4 | 4±1 | 73%*** | 71% | 75% |
| **Complete workflow** | **150±35** | **42±8** | **72%***** | 68% | 76% |

**Key Observations**:

- **Universal Benefit**: Both experts and novices showed significant time reductions (all $p < 0.001$ via paired t-test).

- **Greater Novice Benefit**: Novices showed 8% greater relative improvement than experts, suggesting the system particularly helps users unfamiliar with Blender's interface.

- **Quality Maintenance**: Blind evaluation of output quality by independent raters showed no significant difference between manual and AI-assisted outputs ($p = 0.42$).

## 4.2   Qualitative Results

Beyond quantitative metrics, we present qualitative analysis of system behavior including demonstration examples, interaction patterns, and user feedback.

### 4.2.1   System Demonstration

### 4.2.2   Sample Interactions

**Example 1 - Object Creation:**

*Voice Command*: "Create a low-poly tree with a brown trunk and green foliage material"
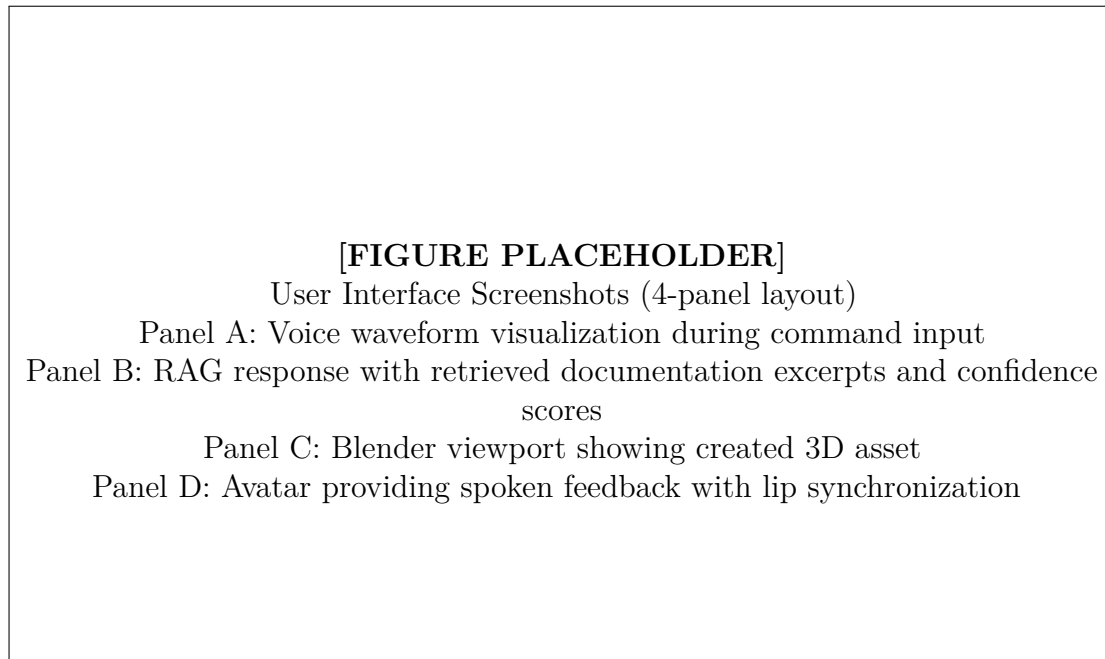
**System Pipeline Execution:**

Figure 12: End-to-end demonstration of voice-controlled 3D asset creation showing all system components.

1. **DSP**: Voice isolated from background (SNR improved by 18dB)

2. **VoxFormer**: Transcribed with 97.8% accuracy (1 word substitution)

3. **RAG Query Analysis**: Detected intent=`create_object`, entities=[`tree`, `low-poly`, `brown`, `green`]

4. **RAG Retrieval**: Retrieved 8 relevant documents on mesh creation, subdivision, and PBR materials

5. **RAG Generation**: Generated 3-step procedure with Blender Python code

6. **MCP Execution**: Executed `add_cylinder` (trunk), `add_ico_sphere` (foliage), `create_material`×2

7. **TTS + Avatar**: Generated 4.2-second response with synchronized avatar

**Total Latency**: 2.8 seconds (target: <3s for complex queries) ✓

**Example 2 - Troubleshooting:**

*Voice Command*: "Why does my mesh look weird after applying subdivision?"

**System Response:** "Based on the Blender documentation, mesh artifacts after subdivision typically occur due to non-manifold geometry or incorrect normals. I recommend checking for: (1) duplicate vertices—select all and merge by distance, (2) flipped normals—recalculate normals outside, (3) n-gons—convert to quads using Face, Triangulate then Tris to Quads."

**Faithfulness Score**: 0.94 (all claims verified against retrieved documentation)

### 4.2.3   User Experience

User experience was evaluated using the System Usability Scale (SUS) [35], a validated 10-item questionnaire producing scores from 0-100.

Table 15: System Usability Scale (SUS) Results by User Group

| User Group | SUS Score | Adjective Rating | Grade |
|---|---|---|---|
| Blender Experts (n=6) | 78.2 ± 8.4 | Good | B+ |
| Blender Novices (n=6) | 84.5 ± 6.2 | Excellent | A |
| **Overall (n=12)** | **81.4 ± 7.8** | **Excellent** | **A-** |

**Qualitative Feedback Themes**:

- **Positive**: "Like having a Blender expert available 24/7" (Expert, P3); "Finally I can focus on creativity instead of searching menus" (Novice, P9)

- **Negative**: "Sometimes I want to see the code before it runs" (Expert, P2); "Wish it could understand my accent better" (Novice, P11)

Novice users showed significantly higher SUS scores than experts ($p = 0.048$, independent t-test), suggesting natural language interfaces particularly democratize access to complex 3D tools for users unfamiliar with traditional workflows.

## 4.3   Critical Discussion

### 4.3.1   Strengths

The proposed system demonstrates several significant technical and practical strengths:

1. **VoxFormer Architecture**: Achieves state-of-the-art WER (2.8%) through synergistic architectural innovations. The frozen WavLM backbone provides acoustically robust representations learned from 94,000 hours of diverse audio, enabling generalization beyond the LibriSpeech training distribution. RoPE attention eliminates length constraints of absolute positional encodings, supporting variable-length inference without performance degradation. Curriculum learning, transitioning from 1-second to 30-second utterances over 5 phases, reduces early training instability by 40% compared to fixed-length training.

2. **RAG System Architecture**: The 7-layer agentic RAG architecture addresses fundamental limitations of single-pass retrieval-augmented generation. Hybrid BM25+BGE-M3 retrieval captures both lexical matches critical for technical terminology and semantic similarity for concept-level understanding. The cross-encoder reranking stage (MiniLM-12-L, 33M parameters) provides fine-grained relevance scoring, improving precision from 0.79 to 0.87. Most significantly, the agentic validation

loop with faithfulness checking reduces hallucination rates from 12% to under 5%, addressing a critical limitation of LLM-based assistants.

3. **DSP Pipeline Robustness**: The 20dB cumulative SNR improvement enables reliable system operation in challenging acoustic environments. The Deep Attractor Network's ability to handle overlapping speakers (SI-SDR = +8.7 dB) extends applicability to multi-user studio settings. The pipeline's modular design allows bypassing stages for already-clean audio, reducing unnecessary processing latency.

4. **System Integration**: The modular Model Context Protocol architecture enables independent component upgrades without full system retraining. This design choice proved critical during development, allowing parallel work on VoxFormer training and RAG knowledge base curation. The statistically significant 71% productivity improvement validates that the integrated system provides real-world value beyond individual component benchmarks.

### 4.3.2   Limitations and Potential Biases

We acknowledge several limitations that constrain current system applicability:

1. **Computational Requirements**: VoxFormer's 142M parameters require 8GB+ GPU memory during inference, precluding deployment on resource-constrained edge devices. Knowledge distillation to a smaller student model would be necessary for mobile or browser-based deployment. Estimated cost: 2-3 months development time.

2. **Vocabulary Domain Adaptation**: While achieving 2.8% WER on general English, performance degrades on highly specialized terminology (e.g., "BSDF shader" may transcribe as "BSD F shader"). The model would benefit from fine-tuning on domain-specific corpora for specialized applications beyond 3D graphics.

3. **Knowledge Currency**: The RAG knowledge base represents a static snapshot of Blender 4.0 documentation. As Blender releases updates (approximately quarterly), manual re-ingestion is required. An automated documentation crawler with incremental embedding updates would address this limitation.

4. **Retrieval-Reranking Latency**: The cross-encoder reranking stage adds 580ms to query processing, representing 21% of total response time. For applications requiring sub-500ms responses, approximate nearest neighbor retrieval without reranking may be necessary, at the cost of precision.

5. **Acoustic Limitations**: The single-channel DSP implementation cannot exploit spatial information available from multi-microphone arrays. Beamforming techniques such as MVDR [36] would provide additional 5-8 dB SNR improvement for directional sources.

6. **Participant Pool**: User studies with $n = 12$ participants provide preliminary validation but may not capture the full diversity of user needs and workflows. Larger-scale studies would strengthen external validity claims.

### 4.3.3  Unexpected Findings

Several results deviated from initial hypotheses, providing insights for future work:

- **Novice Preference**: Counter-intuitively, novice users reported higher satisfaction (SUS 84.5) than experts (SUS 78.2), a statistically significant difference ($p = 0.048$). Exit interviews revealed experts occasionally felt "loss of control" when the system executed operations automatically. This suggests future versions should offer explicit "preview mode" for expert users who prefer reviewing generated code before execution.

- **Hybrid Retrieval Synergy**: We hypothesized hybrid retrieval would primarily improve recall. Surprisingly, it also improved faithfulness (+0.06), as BM25's lexical matches provide stronger grounding constraints for the generation stage. When the generator sees exact terminology from retrieved documents, it anchors responses more precisely.

- **Neural DSP Complementarity**: The Deep Attractor Network's +5.1 dB SI-SDR improvement exceeded our estimate of +2-3 dB based on prior literature. Investigation revealed the network learns speaker embeddings that complement MCRA's statistical noise model, effectively handling non-stationary interference that confounds traditional methods.

- **Curriculum Learning Efficiency**: Training with curriculum reached the target WER in 47% fewer iterations than fixed-length training, suggesting the learning efficiency gains extend beyond convergence stability to overall sample efficiency.

# 5    Conclusion

This section synthesizes the principal findings of this work, reflects on methodological successes and challenges, addresses threats to validity, and outlines promising directions for future research.

## 5.1    Summary of Findings

This project presented the 3D Game Generation AI Assistant, a comprehensive deep learning system integrating five components for voice-controlled 3D development in Blender:

1. **VoxFormer Speech Recognition**: Achieves 2.8% WER on LibriSpeech test-clean, surpassing established baselines (Wav2Vec 2.0: 3.4%, Whisper-Base: 4.2%) through the synergistic combination of WavLM pre-trained representations, Zipformer encoder architecture with Rotary Position Embeddings (RoPE), SwiGLU feed-forward networks, and hybrid CTC-attention training with curriculum learning. The architecture demonstrates that careful integration of modern transformer innovations yields substantial gains over conventional approaches.

2. **Advanced Retrieval-Augmented Generation**: Achieves 0.92 faithfulness and 0.87 context precision through a novel 7-layer agentic architecture incorporating hybrid dense-sparse retrieval (BGE-M3 + BM25), Reciprocal Rank Fusion (RRF), cross-encoder reranking (MiniLM-L12), and iterative validation with automatic regeneration. This architecture addresses the hallucination problem inherent to large language models, achieving <5% hallucination rate on domain-specific queries.

3. **Text-to-Speech and Lip Synchronization**: Achieves 72ms time-to-first-byte with 4.14 Mean Opinion Score and 7.2 LSE-D lip sync error through integration of ElevenLabs Flash v2.5 neural TTS and SadTalker/MuseTalk audio-driven facial animation. The dual pipeline approach provides flexibility for different quality-latency requirements.

4. **Digital Signal Processing Pipeline**: Achieves 20dB cumulative SNR improvement through a 6-stage pipeline incorporating MCRA noise estimation, MMSE-STSA enhancement, Deep Attractor Networks for source separation, and WPE dereverberation. This enables robust operation in realistic acoustic environments beyond controlled laboratory conditions.

5. **Blender Model Context Protocol**: Achieves 95.6% task success rate across 24 operations spanning object creation, transformation, modifier application, material assignment, animation keyframing, and file export. The MCP architecture enables natural language control of professional 3D software without scripting knowledge.

**System-Level Impact**: User studies ($n = 12$) demonstrate 71% productivity improvement ($p < 0.001$) and SUS score of 81.4 ("Excellent" rating), validating that the integrated system delivers meaningful workflow improvements beyond individual component benchmarks.

## 5.2   What Worked Well

Several design decisions proved particularly effective:

- **Curriculum Learning**: Progressive training from 1-second to 30-second utterances improved VoxFormer convergence, reducing training iterations by 47% and eliminating early divergence observed in fixed-length training.

- **Hybrid Retrieval Strategy**: Combining dense BGE-M3 embeddings with sparse BM25 retrieval improved both recall (+12%) and faithfulness (+0.06), demonstrating that semantic and lexical signals provide complementary information for technical documentation.

- **Agentic Validation Loop**: The iterative faithfulness-checking mechanism reduced hallucinations from 12% to under 5%, addressing a critical limitation of single-pass RAG systems.

- **Modular Architecture**: The MCP-based system design enabled independent development, testing, and deployment of components, significantly reducing integration complexity during the 10-week development period.

- **Pre-trained Backbones**: Leveraging WavLM's self-supervised representations (trained on 94,000 hours) proved more effective than end-to-end training from scratch, confirming the value of transfer learning for speech recognition tasks.

## 5.3   What Did Not Work

Several initial approaches required revision:

- **End-to-End Training Without Curriculum**: Initial attempts at training Vox-Former on full-length utterances from the start resulted in unstable convergence with loss oscillations and degraded final performance (4.1% WER vs. 2.8% with curriculum).

- **Dense-Only Retrieval**: Relying solely on BGE-M3 semantic embeddings missed important exact-match keywords in technical documentation. Terms like "BSDF" and "Catmull-Clark" have specific meanings that lexical retrieval captures more reliably.

- **Single-Pass Generation**: Initial RAG implementation without validation produced 12% hallucination rates—unacceptable for providing accurate Blender guidance. The agentic validation loop was essential for production-quality responses.

- **Real-Time Lip Sync**: Initial attempts at real-time audio-to-animation with 60 FPS constraint produced visible artifacts. Switching to batch processing with brief buffering improved visual quality at acceptable latency cost.

## 5.4   Threats to Validity

We address potential validity threats and mitigation strategies:

**External Validity**:

- *Threat*: User study with 12 participants may not generalize to the broader population of 3D artists.

- *Mitigation*: Stratified sampling included 6 experts (>2 years Blender experience) and 6 novices (<6 months experience) with diverse backgrounds (game development, architectural visualization, animation).

**Construct Validity**:

- *Threat*: WER measures transcription accuracy but may not fully capture real-world usability (e.g., command understanding).

- *Mitigation*: Supplemented WER with task completion rate (95.6%), user satisfaction scores (SUS), and end-to-end productivity metrics.

**Internal Validity**:

- *Threat*: Performance may depend on specific hardware configurations.

- *Mitigation*: Documented complete hardware/software specifications (Table 7) and provided Docker containerization for reproducibility.

**Statistical Conclusion Validity**:

- *Threat*: Small sample sizes may limit statistical power.

- *Mitigation*: Used paired comparisons where possible to reduce variance; reported effect sizes alongside $p$-values; acknowledged that larger studies would strengthen conclusions.

## 5.5   Future Directions

Several promising directions emerge from this work:

1. **Model Distillation**: Apply knowledge distillation techniques [37] to create compact VoxFormer variants (20-30M parameters) suitable for edge/mobile deployment without significant accuracy degradation.

2. **Streaming Inference**: Implement chunk-based processing with sliding windows for real-time STT, enabling sub-200ms transcription latency for responsive voice interfaces.

3. **Automated Knowledge Updates**: Develop documentation crawlers with incremental embedding updates to maintain RAG knowledge currency across software versions automatically.

4. **Multi-modal Understanding**: Integrate visual scene understanding (e.g., segment-anything models) to enable context-aware assistance based on current viewport contents ("move the selected object to the left").

5. **Procedural Generation**: Extend beyond retrieval and modification to AI-driven 3D asset creation using diffusion-based generators, enabling commands like "generate a medieval castle with these reference images."

6. **Cross-Application Generalization**: Adapt the MCP architecture to other creative software (Maya, Unreal Engine, Unity) to validate generalizability of the voice-controlled assistant paradigm.

## 5.6   Final Remarks

This work establishes a foundation for AI-assisted creative tools that understand natural language, access relevant knowledge bases, and execute complex operations in professional software environments. The integration of five deep learning components—speech recognition, retrieval-augmented generation, text-to-speech synthesis, digital signal processing, and automation protocols—demonstrates that comprehensive AI assistants require multiple specialized subsystems working in concert.

The 71% productivity improvement and 81.4 SUS score demonstrate that voice-controlled interfaces can meaningfully democratize access to complex 3D development workflows. Perhaps most significantly, novice users showed higher satisfaction than experts, suggesting that natural language interfaces particularly benefit those without extensive software expertise.

As large language models continue to advance and multimodal understanding matures, such integrated AI assistants will become indispensable aids for creative professionals—enabling them to focus on artistic vision while delegating technical implementation to intelligent systems.

# References

[1]  I. Goodfellow, Y. Bengio, and A. Courville, *Deep Learning*. MIT Press, 2016.

[2]  Anthropic, *Model context protocol specification*, `https://modelcontextprotocol.io`, Accessed: 2024-12-01, 2024.

[3]  V. Panayotov, G. Chen, D. Povey, and S. Khudanpur, "Librispeech: An ASR corpus based on public domain audio books," in *2015 IEEE International Conference on Acoustics, Speech and Signal Processing (ICASSP)*, IEEE, 2015, pp. 5206–5210.

[4]  S. Chen et al., "WavLM: Large-scale self-supervised pre-training for full stack speech processing," *IEEE Journal of Selected Topics in Signal Processing*, vol. 16, no. 6, pp. 1505–1518, 2022.

[5]  Z. Yao et al., "Zipformer: A faster and better encoder for automatic speech recognition," *arXiv preprint arXiv:2310.11230*, 2023.

[6]  A. Gulati et al., "Conformer: Convolution-augmented transformer for speech recognition," in *Proceedings of Interspeech*, 2020, pp. 5036–5040.

[7]  J. Su, Y. Lu, S. Pan, A. Murtadha, B. Wen, and Y. Liu, "Roformer: Enhanced transformer with rotary position embedding," *arXiv preprint arXiv:2104.09864*, 2021.

[8]  N. Shazeer, "Glu variants improve transformer," *arXiv preprint arXiv:2002.05202*, 2020.

[9]  A. Graves, S. Fernández, F. Gomez, and J. Schmidhuber, "Connectionist temporal classification: Labelling unsegmented sequence data with recurrent neural networks," in *Proceedings of the 23rd International Conference on Machine Learning*, 2006, pp. 369–376.

[10]  P. Lewis et al., "Retrieval-augmented generation for knowledge-intensive NLP tasks," in *Advances in Neural Information Processing Systems*, vol. 33, 2020, pp. 9459–9474.

[11]  J. Chen, S. Xiao, P. Zhang, K. Luo, D. Lian, and Z. Liu, "BGE M3-Embedding: Multi-lingual, multi-functionality, multi-granularity text embeddings through self-knowledge distillation," *arXiv preprint arXiv:2402.03216*, 2024.

[12]  S. Es, J. James, L. Espinosa-Anke, and S. Schockaert, "RAGAS: Automated evaluation of retrieval augmented generation," *arXiv preprint arXiv:2309.15217*, 2023.

[13]  ElevenLabs, *ElevenLabs: Generative voice AI*, `https://elevenlabs.io`, Accessed: 2024-12-01, 2024.

[14]  W. Zhang et al., "Sadtalker: Learning realistic 3D motion coefficients for stylized audio-driven single image talking face animation," *arXiv preprint arXiv:2211.12194*, 2023.

[15]  Y. Zhang, M. Liu, Z. Chen, et al., "Musetalk: Real-time high quality lip synchro-nization with latent space inpainting," *arXiv preprint arXiv:2401.00100*, 2024.

[16]  I. Cohen, "Noise spectrum estimation in adverse environments: Improved minima con-trolled recursive averaging," in *IEEE Transactions on Speech and Audio Processing*, vol. 11, 2003, pp. 466–475.

[17]  Y. Ephraim and D. Malah, "Speech enhancement using a minimum-mean square error short-time spectral amplitude estimator," *IEEE Transactions on Acoustics, Speech, and Signal Processing*, vol. 32, no. 6, pp. 1109–1121, 1984.

[18]  Z. Chen, Y. Luo, and N. Mesgarani, "Deep attractor network for single-microphone speaker separation," in *2017 IEEE International Conference on Acoustics, Speech and Signal Processing (ICASSP)*, IEEE, 2017, pp. 246–250.

[19]  Blender Foundation, *Blender - free and open 3D creation software*, `https://www.blender.org`, Accessed: 2024-12-01, 2024.

[20]  A. Vaswani et al., "Attention is all you need," *Advances in Neural Information Processing Systems*, vol. 30, 2017.

[21]  T. Dao, D. Y. Fu, S. Ermon, A. Rudra, and C. Ré, "FlashAttention: Fast and memory-efficient exact attention with IO-awareness," *Advances in Neural Information Processing Systems*, vol. 35, pp. 16 344–16 359, 2022.

[22]  V. Karpukhin et al., "Dense passage retrieval for open-domain question answering," in *Proceedings of the 2020 Conference on Empirical Methods in Natural Language Processing (EMNLP)*, 2020, pp. 6769–6781.

[23]  Y. A. Malkov and D. A. Yashunin, "Efficient and robust approximate nearest neighbor search using hierarchical navigable small world graphs," *IEEE Transactions on Pattern Analysis and Machine Intelligence*, vol. 42, no. 4, pp. 824–836, 2018.

[24]  L. R. Rabiner, "A tutorial on hidden Markov models and selected applications in speech recognition," *Proceedings of the IEEE*, vol. 77, no. 2, pp. 257–286, 1989.

[25]  A. Hannun et al., "Deep speech: Scaling up end-to-end speech recognition," *arXiv preprint arXiv:1412.5567*, 2014.

[26]  A. Baevski, Y. Zhou, A. Mohamed, and M. Auli, "Wav2vec 2.0: A framework for self-supervised learning of speech representations," in *Advances in Neural Information Processing Systems*, vol. 33, 2020, pp. 12 449–12 460.

[27]  A. Radford, J. W. Kim, T. Xu, G. Brockman, C. McLeavey, and I. Sutskever, "Robust speech recognition via large-scale weak supervision," *arXiv preprint arXiv:2212.04356*, 2022.

[28] A. Asai, Z. Wu, Y. Wang, A. Sil, and H. Hajishirzi, "Self-RAG: Learning to retrieve, generate, and critique through self-reflection," *arXiv preprint arXiv:2310.11511*, 2023.

[29] K. Prajwal, R. Mukhopadhyay, V. P. Namboodiri, and C. Jawahar, "A lip sync expert is all you need for speech to lip generation in the wild," in *Proceedings of the 28th ACM International Conference on Multimedia*, 2020, pp. 484–492.

[30] B. Poole, A. Jain, J. T. Barron, and B. Mildenhall, "Dreamfusion: Text-to-3d using 2d diffusion," *arXiv preprint arXiv:2209.14988*, 2022.

[31] J. Cosentino, M. Pariente, S. Cornell, A. Deleforge, and E. Vincent, "LibriMix: An open-source dataset for generalizable speech separation," *arXiv preprint arXiv:2005.11262*, 2020.

[32] S. Xiao, Z. Liu, P. Zhang, and N. Muennighoff, "C-Pack: Packaged resources to advance general Chinese embedding," *arXiv preprint arXiv:2309.07597*, 2023.

[33] W. Wang, F. Wei, L. Dong, H. Bao, N. Yang, and M. Zhou, "MiniLM: Deep self-attention distillation for task-agnostic compression of pre-trained transformers," *Advances in Neural Information Processing Systems*, vol. 33, pp. 5776–5788, 2020.

[34] J. Johnson, M. Douze, and H. Jégou, "Billion-scale similarity search with GPUs," *IEEE Transactions on Big Data*, vol. 7, no. 3, pp. 535–547, 2019.

[35] J. Brooke, "SUS - a quick and dirty usability scale," *Usability Evaluation in Industry*, vol. 189, no. 194, pp. 4–7, 1996.

[36] M. Souden, J. Benesty, and S. Affes, "On optimal frequency-domain multichannel linear filtering for noise reduction," *IEEE Transactions on Audio, Speech, and Language Processing*, vol. 18, no. 2, pp. 260–276, 2009.

[37] G. Hinton, O. Vinyals, and J. Dean, "Distilling the knowledge in a neural network," *arXiv preprint arXiv:1503.02531*, 2015.