# Lightweight Time-Domain Audio-Visual Speech Enhancement Model

Fazal E Wahab, Nasir Saleem, Kia Dashtipour, Arif Reza Anwary
Adeel Hussain, Mandar Gogate, Amir Hussain

July 6, 2025

## Introduction

This document presents a lightweight architecture for audio-visual speech enhancement that combines visual cues from lip movements with corrupted audio signals to produce clean speech. The model is designed for real-time applications with limited computational resources.

## 1 Model Architecture

The proposed architecture consists of four main modules: Visual Feature Extrac-tor, Audio Encoder, Cross-Modal Fusion, and Audio Decoder.

### 1.1 Visual Feature Extractor

The visual module processes lip movement frames to extract temporal-spatial features:

- **Input**: Sequence of 112×112 RGB lip region crops (T frames)

- **3D Convolution Block**: Single 3D convolution (kernel 3×5×5) with BatchNorm and ReLU

- **ResNet-18 Lite**: Modified ResNet-18 with:

  - Reduced channel counts (32, 64, 128, 256)
  - 2D temporal average pooling after conv5
  - Output shape: T×256

- **Temporal ConvNet**: Two 1D convolutions (kernel 3) with dilation factors 1 and 2 to capture temporal dynamics

- **Output**: Visual features $V \in R^{T \times D_v}$ where $D_v = 128$

## 1.2 Audio Encoder

Processes noisy speech to extract spectral features:

- **Input**: Noisy speech STFT $X \in R^{F \times T}$ (F=257, T=100 for 1s at 16kHz)

- **Conv Blocks**: Two 2D convolution layers (kernel 3×3) with BatchNorm and PReLU

- **GRU Layer**: Bidirectional GRU with 64 hidden units

- **Attention**: Temporal attention layer to weight important frames

- **Output**: Audio features $A \in R^{T \times D_a}$ where $D_a = 128$

## 1.3 Cross-Modal Fusion

Combines visual and audio features effectively:

- **Cross-Attention**: Multi-head attention (4 heads) between visual and audio features

$$F_{fusion} = \text{LayerNorm}(A + \text{MultiHead}(A, V, V)) \tag{1}$$

- **Gated Fusion**: Learnable weights combine modalities:

$$F_{final} = \alpha \cdot F_{fusion} + (1 - \alpha) \cdot A \tag{2}$$

  where $\alpha$ is a learned parameter (sigmoid-activated)

- **Temporal Conv**: Two 1D convolutions to smooth fused features

## 1.4 Audio Decoder

Reconstructs clean speech from fused features:

- **GRU Layer**: Uni-directional GRU with 128 hidden units

- **Conv Blocks**: Two transposed convolutions to upsample features

- **Mask Prediction**: 1×1 convolution to estimate complex ideal ratio mask (cIRM)

- **Output**: Enhanced STFT $\hat{Y} = X \odot M$ where $M$ is the predicted mask

Table 1: Model parameter counts

| Module | Parameters |
|---|---|
| Visual Feature Extractor | 1.2M |
| Audio Encoder | 0.8M |
| Cross-Modal Fusion | 0.3M |
| Audio Decoder | 0.9M |
| Total | 3.2M |

# 2  Implementation Details

## 2.1  Model Parameters

## 2.2  Training Strategy

- **Loss Function**: Combination of spectral convergence and magnitude loss:

$$\mathcal{L} = \||Y| - |\hat{Y}|\|_1 + \lambda\|\frac{|Y| - |\hat{Y}|}{|Y|}\|_2 \tag{3}$$

- **Optimizer**: AdamW with learning rate 3e-4

- **Regularization**: Dropout (0.2) and weight decay (1e-4)

# 3  Advantages

- **Lightweight**: Only 3.2M parameters (12MB)

- **Efficient**: Processes 1s audio in 15ms on mobile CPU

- **Robust**: Works well with various noise types (SNR 0-20dB)

- **Adaptive**: Gated fusion automatically adjusts to input quality