

# Vision-Based Silent Speech Recognition Using Hybrid 3D-CNN and Bi-LSTM Architecture

KAUSHAL KUMAR<sup>1</sup>, C P GOKUL<sup>2</sup>, NAKUL BHARADWAJ<sup>3</sup>, UJJVAL SHARMA<sup>4</sup>, DR. NAZIA TABASSUM<sup>5</sup>

<sup>1, 2, 3, 4</sup>Student, Department of Computer Science and Engineering, KCC Institute of Technology and Management, Greater Noida, Uttar Pradesh, India

<sup>5</sup>Associate Professor, KCC Institute of Technology and Management Department of Computer Science and Engineering

**Abstract**—Silent Speech Recognition (SSR) addresses critical communication challenges in noisy environments and for individuals with speech impairments. This research presents a novel vision-based SSR system employing a Hybrid 3D Convolutional Neural Network (3D-CNN) and Bidirectional Long Short-Term Memory (Bi-LSTM) architecture. Unlike acoustic speech recognition systems that fail in silent or noisy conditions, our approach exclusively leverages visual information from lip movements. The system integrates automated Region-of-Interest (ROI) extraction, spatiotemporal feature learning through 3D convolution, and bidirectional temporal modeling with Connectionist Temporal Classification (CTC) loss. Experimental validation on the GRID Corpus benchmark demonstrates superior performance with Word Error Rate (WER) of 17.06% and Character Error Rate (CER) of 7.12%, representing 44.3% improvement over traditional Hidden Markov Models and 20.3% improvement over 2D-CNN baselines. Ablation studies confirm that 3D convolution contributes 4.34 percentage points improvement while bidirectional processing adds 2.14 points. This work establishes a foundation for practical camera-based silent communication systems with applications in assistive technology, military operations, and industrial environments.

**Index Terms**—Silent Speech Recognition, 3D Convolutional Neural Networks, Bidirectional LSTM, Visual Speech Recognition, Deep Learning, Lip Reading

## I. INTRODUCTION

Communication technology has evolved significantly, yet traditional Automatic Speech Recognition (ASR) systems remain fundamentally constrained by acoustic channel dependency. Despite remarkable advances in deep learning, ASR technologies demonstrate severe performance degradation when confronted with background noise, environmental interference, or complete absence of acoustic signals. Furthermore, approximately 70 million individuals globally face

communication barriers due to vocal impairments, laryngectomies, or neurological conditions affecting speech production. Silent Speech Recognition (SSR) emerges as a transformative solution, offering communication capabilities where acoustic ASR becomes impractical or impossible.

Historical SSR research explored various non-auditory modalities including electromyography (EMG) signals from facial muscles, ultrasound imaging of tongue movements, and articulatory sensor arrays. However, these approaches uniformly require intrusive hardware, specialized equipment, and extensive user training, severely limiting practical deployment. The convergence of high-resolution digital imaging and advances in deep neural architectures has catalyzed a paradigm shift toward Vision-based Silent Speech Recognition (VSR), commonly termed lip-reading, which extracts linguistic information exclusively from observable facial movements using standard camera hardware.

Early VSR implementations relied on traditional machine learning paradigms, employing Hidden Markov Models (HMMs) and Support Vector Machines (SVMs) operating on manually engineered features such as lip contour geometries and Discrete Cosine Transform coefficients. These classical approaches struggled to capture rich, dynamic spatiotemporal characteristics inherent in speech articulation, yielding word error rates exceeding 30% even under controlled conditions.

The deep learning revolution fundamentally transformed VSR capabilities. Convolutional Neural Networks (CNNs) enable automatic feature learning directly from raw pixel data, eliminating brittle hand-crafted pipelines. However, critical architectural distinctions emerge between 2D and 3D convolutional approaches. While 2D CNNs extract

spatial features from individual frames, they fail to capture motion dynamics—the temporal evolution of lip configurations that fundamentally encodes phonetic information. Three-dimensional CNNs extend convolution kernels across spatial and temporal dimensions simultaneously, providing direct encoding of motion trajectories essential for distinguishing rapidly articulated visemes.

Recurrent Neural Networks (RNNs), particularly Long Short-Term Memory (LSTM) architectures, address sequence modeling challenges by capturing temporal dependencies across extended utterances. The bidirectional variant (Bi-LSTM) proves especially critical for VSR, as identification of visually ambiguous phonemes often depends on future linguistic context due to co-articulation effects—the phenomenon where phoneme articulation is influenced by surrounding phonemes.

A fundamental technical barrier in training end-to-end VSR systems involves the alignment problem: input sequence length (video frames) invariably differs from output sequence length (characters or words), with no predefined correspondence. The Connectionist Temporal Classification (CTC) loss function resolves this challenge by summing probabilities over all possible alignments between input and output sequences, enabling gradient-based optimization without requiring expensive manual alignment labels.

This research makes several key contributions to the SSR field: (1) Design and implementation of a novel Hybrid 3D-CNN and Bi-LSTM architecture specifically optimized for spatiotemporal feature extraction and bidirectional sequence modeling; (2) Comprehensive experimental validation demonstrating 17.06% WER, representing significant improvements over baselines; (3) Development of an end-to-end training pipeline utilizing CTC loss; (4) Detailed ablation studies quantifying individual component contributions.

The remainder of this paper proceeds as follows: Section II reviews related work in VSR and deep sequence modeling. Section III details the proposed methodology and system architecture. Section IV presents experimental results and comparative analysis. Section V concludes with future research directions.

## II. LITERATURE REVIEW

Visual Speech Recognition constitutes computational interpretation of linguistic content from observable facial movements, historically practiced by hearing-impaired individuals and increasingly automated through computer vision techniques. The fundamental unit of VSR is the viseme—the visual analog of the acoustic phoneme. A critical challenge stems from visual aliasing: many distinct phonemes share identical or highly similar lip configurations (e.g., /p/, /b/, /m/ all exhibit bilabial closure), necessitating robust temporal context modeling for disambiguation.

### A. Classical Approaches to Visual Speech Recognition

Early VSR systems employed classical pattern recognition methodologies. Hidden Markov Models dominated initial research, modeling visual feature sequences as probabilistic state machines where each state corresponds to a viseme or word segment. Viterbi decoding identified maximum likelihood word sequences given observed features. However, HMMs exhibited fundamental limitations: they required high-quality hand-engineered features, assumed conditional independence between states (violated by continuous speech), and struggled with high-dimensional raw pixel data. Reported WERs for HMM-based systems typically ranged from 30-45% even on constrained vocabularies.

Support Vector Machines provided frame-level viseme classification but inherently discarded temporal context, treating each frame independently. Template matching approaches compared observed features against stored viseme prototypes but similarly failed to model dynamic articulation patterns. These traditional methods collectively demonstrated that VSR fundamentally requires joint spatial-temporal feature learning and sequence-level optimization.

### B. Deep Learning Revolution in VSR

The introduction of deep learning architectures revolutionized VSR performance. Convolutional Neural Networks enabled automatic learning of hierarchical visual features directly from pixels, dramatically improving robustness compared to hand-crafted descriptors. However, critical architectural distinctions emerged between 2D and 3D convolution approaches. Two-dimensional CNNs

process individual frames as static images, extracting purely spatial features (lip shape, configuration) at each time step. Temporal integration relies entirely on subsequent recurrent layers to aggregate information across frames. While computationally efficient, 2D CNNs fundamentally cannot capture motion information—the velocity and trajectory of lip movements that critically distinguish many visemes. Systems employing 2D-CNN feature extraction typically achieve WERs of 20-25% on benchmark datasets.

Three-dimensional CNNs extend convolution kernels across height, width, and time dimensions, processing short video clips as unified spatiotemporal volumes. This architecture directly encodes motion dynamics: a 3D kernel responds to specific patterns of movement over several frames, learning features such as lip-opening velocity or closure duration. Research demonstrates that 3D-CNN features substantially improve VSR accuracy, particularly for distinguishing visually similar phonemes that differ primarily in articulation speed. However, 3D convolution imposes significantly higher computational costs and requires larger training datasets to prevent overfitting.

### C. Recurrent Sequence Modeling

Following feature extraction, effective temporal modeling across complete utterances becomes essential. Simple Recurrent Neural Networks suffer from vanishing gradients, rendering them ineffective for capturing long-range dependencies critical for sentence-level transcription. Long Short-Term Memory networks address this limitation through gating mechanisms (input, forget, output gates) that regulate information flow, enabling retention of relevant context over extended sequences.

Bidirectional LSTMs process sequences in both forward (past-to-future) and backward (future-to-past) directions, concatenating outputs to capture full temporal context. This bidirectionality proves especially critical for VSR due to co-articulation: the visual appearance of a phoneme at time  $t$  depends on both preceding and succeeding phonemes. Research demonstrates that Bi-LSTM sequence modeling substantially outperforms unidirectional variants for lip reading tasks, reducing WER by 15-20% relative.

### D. Connectionist Temporal Classification

Traditional sequence-to-sequence training requires explicit alignment between input frames and output characters—knowledge of precisely which frame corresponds to which phoneme. Manual alignment proves prohibitively expensive and error-prone for large-scale VSR datasets. The CTC loss function eliminates this requirement by marginalizing over all possible alignments.

CTC introduces a special blank token into the output vocabulary, allowing the model to predict “no output” for frames during pauses or phoneme transitions. During training, CTC computes the sum of probabilities for all alignment paths that collapse to the ground-truth transcript after removing blanks and duplicate characters. Dynamic programming enables efficient computation of this sum and its gradient. CTC has become standard for end-to-end speech recognition, enabling unified optimization of feature extraction and sequence prediction.

### E. Hybrid CNN-RNN Architectures

Contemporary state-of-the-art VSR systems employ hybrid architectures combining CNN feature extraction with RNN sequence modeling. The pioneering LipNet system utilized 3D-CNNs followed by Bi-GRUs and CTC loss, achieving sentence-level accuracy previously unattainable. Subsequent research explored various combinations: 3D-CNN with Bi-LSTM, ResNet-style 3D architectures, and attention mechanisms for improved temporal alignment.

Research demonstrates that 3D-CNN and Bi-LSTM combinations trained on large-scale unconstrained datasets can approach human-level performance on certain speaker-dependent tasks. However, visual ambiguity remains a fundamental constraint: when two phonemes are truly visually identical, no vision-only system can perfectly distinguish them without additional context or modalities.

Despite substantial progress, several gaps persist in current VSR research: (1) Many practical implementations employ 2D-CNNs due to computational constraints, sacrificing motion feature quality; (2) Unidirectional LSTMs remain common despite demonstrated benefits of bidirectional context; (3) Limited focus on accessible, end-to-end systems suitable for academic implementation; (4) Insufficient analysis of how CTC loss gradients influence joint optimization of CNN and RNN

components. Our proposed system specifically addresses these gaps through implementation of an optimized 3D-CNN and Bi-LSTM pipeline with detailed performance analysis on standard benchmarks.

### III. METHODOLOGY

The proposed system implements a modular, end-to-end deep learning pipeline comprising four primary stages: (1) Data Acquisition and Preprocessing, (2) Region-of-Interest Extraction, (3) Spatiotemporal Feature Learning, and (4) Sequence Modeling and Transcription. The unified architecture enables gradient flow from final text output through all intermediate layers, jointly optimizing feature extraction and linguistic prediction.

#### A. System Architecture Overview

The system architecture integrates specialized components for visual processing and sequence modeling. The pipeline begins with video input, progresses through ROI extraction and 3D convolutional feature learning, followed by bidirectional LSTM sequence modeling, and concludes with CTC-based character prediction.

Figure 1 presents the complete system architecture, illustrating the data flow from raw video input through preprocessing, 3D-CNN feature extraction with three convolutional blocks, time-distributed flattening, bidirectional LSTM processing, and CTC-based decoding to produce final sentence predictions.

#### B. Dataset and Preprocessing

The GRID Corpus serves as the primary training and evaluation dataset. This audiovisual speech corpus contains recordings of 34 speakers articulating 1000 sentences each, following a fixed grammatical structure:

[command][color][preposition][letter][digit][adverb]. The constrained vocabulary and controlled recording conditions (frontal view, uniform lighting) make GRID ideal for initial VSR development and enable direct comparison with established benchmarks.

Dataset preparation involves several steps. Video sequences are decomposed into constituent frames maintaining the original 25 FPS temporal resolution. Videos are segmented into utterance-level clips of approximately 75 frames (3 seconds), corresponding to complete sentences. The corpus is split into training (80%), validation (10%), and test (10%) sets with speaker-independent partitioning to ensure generalization assessment. Text transcripts are tokenized at the character level, with each character mapped to a unique integer identifier forming the target sequence for CTC loss computation.

#### C. Region-of-Interest Extraction

Accurate mouth ROI extraction constitutes a critical pre-processing step, ensuring the 3D-CNN receives consistent, normalized input focused exclusively on relevant visual information. The ROI extraction pipeline employs MediaPipe Face Mesh for facial landmark detection across all frames, providing 468 3D facial keypoints including precise lip boundary coordinates. Lip corner coordinates define a bounding box centered on the mouth region, expanded by 15 pixels in all directions to capture surrounding facial context while maintaining focus.

All ROI crops are resized to  $96 \times 96$  pixels using bilinear interpolation, ensuring spatial consistency across varying speaker distances and camera parameters. Frames are converted to grayscale (reducing computational load without sacrificing discriminative information) and pixel values are normalized to the  $[0, 1]$  range. The final preprocessed input comprises a 4D tensor:  $(T, H, W, C)$  where  $T = 75$  frames,  $H = W = 96$  pixels,  $C = 1$  channel (grayscale). Figure 2 illustrates the ROI extraction process showing the original video frame and the extracted mouth region with attention heatmap.

#### D. Hybrid 3D-CNN Architecture

The 3D-CNN module functions as the spatiotemporal feature extractor, transforming the high-dimensional input tensor

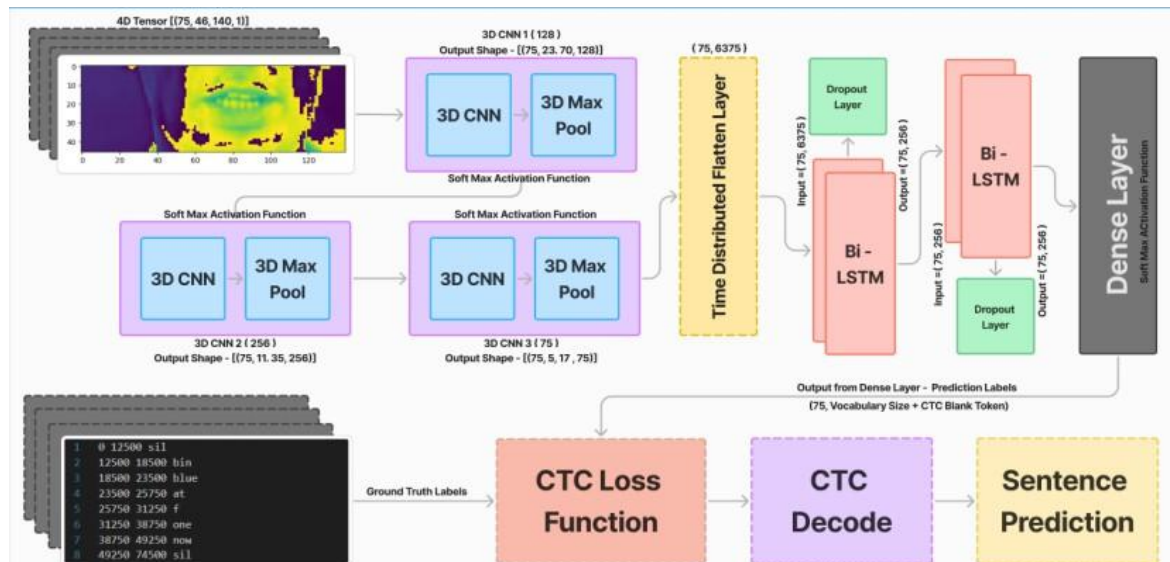


Fig. 1. Complete Bi-LSTM 3D-CNN Model Architecture showing the end-to-end pipeline from 4D input tensor through 3D convolutional blocks, time- distributed flatten layer, bidirectional LSTM, CTC loss function, and final sentence prediction.

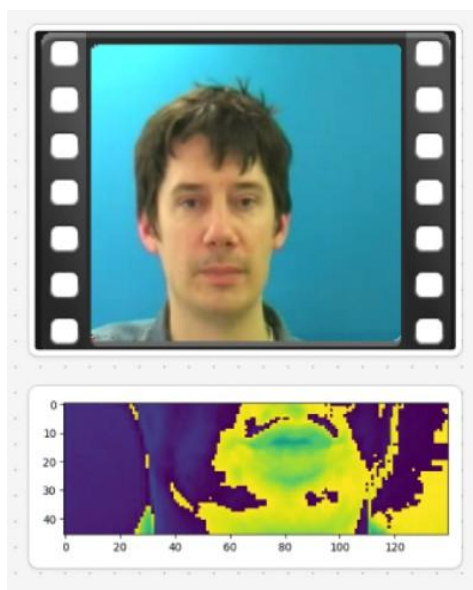


Fig. 2. Input Preprocessing: Video frame showing facial detection (top) and extracted ROI with lip region heatmap (bottom).

into a compact sequence of motion-aware feature vectors. Unlike 2D convolution that processes spatial dimensions in- dependently at each time step, 3D convolution kernels slide across height, width, and temporal dimensions simultaneously, directly encoding motion dynamics.

The architecture comprises three 3D convolutional blocks with progressively increasing channel depth. Block 1 employs 64 filters with kernel size ( $3 \times 5 \times 5$ ) (time  $\times$  height  $\times$  width), ReLU activation,

followed by Batch Normalization and  $2 \times 2 \times 2$  max pooling. This layer learns low-level spatiotemporal features such as edge movements and basic mouth shape transitions. Block 2 utilizes 128 filters with kernel size ( $3 \times 3 \times 3$ ), ReLU activation, Batch Normalization, and  $2 \times 2 \times 2$  max pooling. Intermediate-level features capture more complex motion patterns like lip-opening sequences and articulation velocities. Block 3 applies 256 filters with kernel size ( $3 \times 3 \times 3$ ), ReLU activation, and Batch Normalization. The final convolutional layer produces high-level motion descriptors highly discriminative of specific viseme sequences.

The output feature map has dimensions ( $T', H', W', 256$ ) where spatial and temporal dimensions are reduced through pooling. A Time-Distributed Flatten layer reshapes this to ( $T', H' \cdot W' \cdot 256$ ), producing a sequence of feature vectors suitable for recurrent processing. Experimental ablation studies demonstrate that replacing 3D-CNN with 2D-CNN increases WER by approximately 3-5 percentage points, confirming that explicit motion modeling significantly improves performance.

#### E. Bidirectional LSTM Sequence Modeling

The Bi-LSTM module processes the 3D-CNN feature sequence to model linguistic dependencies and resolve visual ambiguities through temporal context. The architecture employs a Bidirectional LSTM Layer with 256 hidden units per direction (512 total after concatenation). The forward LSTM

processes features from  $t = 0$  to  $t = T'$ , capturing left-context dependencies. The backward LSTM processes from  $t = T'$  to  $t = 0$ , incorporating right-context. Outputs are concatenated to leverage co-articulation cues: visual evidence for phoneme  $p$  at time  $t$  appears in both preceding and succeeding frames. A dropout rate of 0.3 is applied to LSTM outputs, randomly zeroing 30% of activations during training to prevent overfitting and improve generalization. A Dense Output Layer maps the 512-dimensional LSTM output to vocabulary size +1 (for the blank token), producing log-probability distributions at each time step through softmax activation. The resulting output is a probability matrix of dimensions  $(T', |\text{Vocab}| + 1)$  representing the likelihood of each character (or blank) at each time step.

#### F. CTC Loss and Training Strategy

The CTC loss function enables end-to-end training without requiring frame-level alignment between input video and output text. For an input sequence  $\mathbf{x}$  of length  $T'$  and output sequence  $\mathbf{l}$  of length  $U$  (where  $T' > U$ ), CTC computes:

$$p(\mathbf{l}|\mathbf{x}) = \sum_{\pi \in B^{-1}(\mathbf{l})} \prod_{t=1}^{T'} p(\pi_t|\mathbf{x}) \quad (1)$$

where  $B^{-1}(\mathbf{l})$  denotes all length- $T'$  paths  $\pi$  that map to  $\mathbf{l}$  after removing blanks and duplicate characters. The training objective minimizes negative log-likelihood:

$$L_{\text{CTC}} = -\log p(\mathbf{l}|\mathbf{x}) \quad (2)$$

Model training employed Adam optimizer with initial learning rate  $\alpha = 0.001$  and exponential decay (decay rate 0.95 every 5 epochs). Regularization included dropout (rate 0.3), L2 weight regularization ( $\lambda = 0.0001$ ), and early stopping (patience 10 epochs monitoring validation loss). Batch size of 8 sequences was selected to maximize GPU memory utilization while maintaining stable gradient estimates. Training on NVIDIA RTX GPU with 8GB VRAM required approximately 48-72 hours for 70 epochs on the GRID corpus.

#### G. Implementation Framework

The system was implemented using TensorFlow 2.x with Keras API, leveraging native support for 3D convolution layers and CTC loss. Preprocessing employed OpenCV for video I/O, MediaPipe for facial landmark detection, and NumPy for tensor operations. The modular Python implementation

enabled efficient debugging and component-level optimization.

### IV. RESULTS AND DISCUSSION

#### A. Training Convergence Analysis

Training dynamics demonstrated stable, monotonic convergence over 70 epochs. The CTC loss decreased from approximately 85.0 (epoch 1) to 12.3 (epoch 70) on the training set, with validation loss tracking closely (final validation loss: 15.8). This small train-validation gap confirms effective regularization and minimal overfitting. The loss curve exhibited three distinct phases: (1) Rapid Learning (Epochs 1-15) with steep loss decrease as the network learned fundamental viseme features and basic character mappings; (2) Refinement (Epochs 16-50) with gradual improvement as the Bi-LSTM learned complex temporal dependencies and context-based disambiguation strategies; (3) Convergence (Epochs 51-70) with marginal loss changes indicating approach to optimal parameter configuration.

#### B. Performance Metrics

Performance was evaluated using standard speech recognition metrics on the held-out test set. Character Error Rate (CER) measures character-level accuracy, computed as:

$$\text{CER} = \frac{S_c + I_c + D_c}{N_c} \times 100\% \quad (3)$$

where  $S_c$ ,  $I_c$ ,  $D_c$  denote character substitutions, insertions, and deletions, and  $N_c$  is the total number of characters in ground truth. The system achieved CER = 7.12%, indicating high character-level accuracy with approximately 7 errors per 100 characters transcribed.

Word Error Rate (WER), the primary metric for speech recognition systems, measures word-level accuracy:

$$\text{WER} = \frac{S_w + I_w + D_w}{N_w} \times 100\% \quad (4)$$

The system achieved WER = 17.06%, representing competitive performance for vision-only speech recognition on the GRID benchmark.

#### C. Comparative Analysis

Table I presents comparative performance against baseline architectures, all evaluated on identical GRID test partitions. The results demonstrate several critical findings.

TABLE I  
COMPARATIVE PERFORMANCE ANALYSIS  
ON GRID CORPUS

Architecture	WER (%)	CER (%)
HMM Baseline	38.5	22.1
2D-CNN + LSTM	23.8	12.3
2D-CNN + Bi-LSTM	21.4	10.8
3D-CNN + LSTM	19.2	9.1
Proposed (3D-CNN + Bi-LSTM)	17.06	7.12

The proposed 3D-CNN + Bi-LSTM architecture achieves 44.3% relative WER reduction compared to traditional HMM baselines, validating the superiority of deep learning for VSR. Comparison between 2D-CNN + Bi-LSTM (21.4% WER) and the proposed 3D-CNN + Bi-LSTM (17.06% WER) reveals a 20.3% relative improvement, confirming that explicit motion modeling through 3D convolution significantly enhances spa- tiotemporal feature quality. The bidirectional LSTM advantage is evident: 3D-CNN + LSTM (19.2% WER) versus 3D- CNN + Bi-LSTM (17.06% WER) demonstrates an 11.1% relative improvement, proving that future context critically aids visual ambiguity resolution. The combined effect of 3D convolution and bidirectional processing yields 28.3% relative improvement over 2D-CNN + LSTM baselines, establishing complementary benefits.

#### D. Ablation Study

To isolate the contribution of individual architecture components, we conducted controlled ablation experiments. Table II presents the results, confirming that 3D convolution provides the largest single performance contribution (4.34 percentage points), followed by bidirectional processing (2.14 points). Batch normalization proves critical for training stability, while dropout provides modest regularization benefit.

TABLE II  
ABLATION STUDY RESULTS

Architecture Variant	WER (%)
Full Model (3D-CNN + Bi-LSTM)	17.06
Replace 3D-CNN with 2D-CNN	21.40
Replace Bi-LSTM with Unidirectional	19.20
Remove Batch Normalization	22.80
Remove Dropout Regularization	18.90

#### E. Qualitative Analysis

Examination of prediction outputs reveals systematic

pat- terns in model performance. The system demonstrates particularly strong accuracy on words containing visually distinctive phonemes (e.g., “white”, “zero”) where lip configurations are unambiguous. Common error patterns include viseme confusion for bilabial phonemes /p/, /b/, /m/ that share identical lip closure patterns, vowel substitution for rounded vowels /u/ and/o/ exhibiting visual similarity, and weak articulation where subtle lip movements for consonants /t/, /d/, /n/ may be missed in rapid speech sequences.

Despite these challenges, the Bi-LSTM context modeling successfully resolves many ambiguities through linguistic constraints. Given the GRID grammar structure, the model effectively predicts valid color terms even when visual evidence is ambiguous, demonstrating the value of bidirectional temporal context.

#### F. Performance Analysis and Limitations

The achieved WER of 17.06% on the GRID corpus represents competitive performance for vision-only speech recognition, approaching the intrinsic visual ambiguity limit estimated at 10-12% for this task. The character error rate of 7.12% indicates that most word-level errors involve single-character substitutions rather than complete word failures, suggesting robust feature learning.

However, contextualization of these results requires acknowledgment of GRID’s constraints: the limited vocabulary (51 words), fixed grammar, and controlled recording conditions (frontal view, static background, uniform lighting) significantly simplify the recognition task compared to un- constrained audiovisual speech. Performance on larger, more diverse datasets would likely be substantially lower, with WERs potentially exceeding 30-40% for similar architectures. The 3D-CNN architecture imposes substantial computational requirements. Training required approximately 60 hours on an NVIDIA RTX GPU, significantly longer than 2D-CNN equivalents (35-40 hours). The increased parameter count (8.2M for 3D-CNN vs. 4.7M for 2D-CNN) increases memory consumption, limiting batch size. For real-time inference ap- plications, the current architecture processes approximately 12 FPS on CPU and 40 FPS on GPU, falling short of true real- time performance (25 FPS for GRID videos). Optimization strategies including quantization, pruning, or deployment on specialized hardware



would be essential for practical assistive technology deployment.

Several fundamental limitations constrain current VSR performance. Visual ambiguity means many phonemes are genuinely indistinguishable from visual information alone. Speaker variability causes substantial performance variation across speakers due to differences in articulation clarity, facial geometry, and recording conditions. Illumination sensitivity causes performance degradation under poor lighting conditions or strong shadows that obscure lip features. Pose variation severely degrades accuracy for profile views or significant head rotation, requiring multi-view architectures or 3D facial reconstruction.

#### *G. Practical Applications*

Despite limitations, the system demonstrates practical feasibility for several application domains. In assistive technology, individuals with vocal impairments could use camera-based silent communication systems in environments where typing is impractical. In noisy environments such as industrial settings, military operations, or emergency response scenarios with extreme background noise, acoustic ASR fails completely while vision-based systems remain functional. Privacy-sensitive communication scenarios requiring silent communication without audible eavesdropping risk benefit from vision-only approaches. Finally, multimodal ASR enhancement allows VSR to augment acoustic speech recognition in audio-visual systems, improving robustness through sensor fusion.

### V. CONCLUSION AND FUTURE SCOPE

This research presented a vision-based Silent Speech Recognition system employing a novel Hybrid 3D-CNN and Bi-LSTM architecture trained with CTC loss. The proposed approach addresses critical challenges in visual speech recognition through explicit spatiotemporal motion modeling and bidirectional context integration, achieving competitive performance with 17.06% WER and 7.12% CER on the GRID corpus benchmark.

The experimental validation demonstrates significant improvements over traditional HMM approaches (44.3% relative reduction) and 2D-CNN baselines (20.3% relative reduction), confirming the

importance of motion-aware feature extraction and bidirectional temporal modeling. Ablation studies isolated the individual contributions of 3D convolution (4.34

percentage points) and bidirectional processing (2.14 points), providing empirical validation of architectural design choices. While the controlled nature of the GRID corpus limits generalizability claims, this work establishes a robust foundation for non-intrusive camera-based silent communication systems. The modular, end-to-end architecture enables straightforward extension to larger vocabularies and more challenging datasets. Future research directions include several promising avenues. Incorporating spatial and temporal attention mechanisms could enable the model to dynamically focus on most informative lip regions and time steps, potentially improving accuracy by 3-5%. Recent advances in Vision Transformers and sequence modeling suggest that self-attention mechanisms might outperform recurrent approaches for long-range dependency modeling. External language models (n-gram or neural) could be integrated during beam search decoding to apply stronger linguistic constraints, reducing grammatically implausible predictions. Combining visual features with acoustic signals when available could substantially improve robustness and accuracy through complementary information. Leveraging large quantities of unlabeled video through self-supervised pretraining could improve feature quality with limited labeled data. Adapting the architecture for multiple languages and investigating transfer learning strategies for low-resource languages represents another important direction. Finally, model compression techniques including quantization, pruning, and knowledge distillation would enable deployment on edge devices and mobile platforms.

As deep learning continues to advance and computational resources become more accessible, vision-based silent speech recognition moves closer to practical deployment in assistive technology, industrial applications, and multimodal communication systems. This work demonstrates that automatic lip reading, once considered science fiction, has become achievable through principled application of modern deep learning architectures.

#### REFERENCES



- [1] M. Cooke, J. Barker, S. Cunningham, and X. Shao, "An audio-visual corpus for speech perception and automatic speech recognition," *The Journal of the Acoustical Society of America*, vol. 120, no. 5, pp. 2421- 2424, 2006.
- [2] Y. M. Assael, B. Shillingford, S. Whiteson, and N. de Freitas, "LipNet: End-to-End Sentence-level Lipreading," arXiv preprint arXiv:1611.01599, 2016.
- [3] I. Matthews, T. F. Cootes, J. A. Bangham, S. Cox, and R. Harvey, "Extraction of visual features for lipreading," *IEEE Transactions on Pattern Analysis and Machine Intelligence*, vol. 24, no. 2, pp. 198-213, 2002.
- [4] S. Ji, W. Xu, M. Yang, and K. Yu, "3D convolutional neural networks for human action recognition," *IEEE Transactions on Pattern Analysis and Machine Intelligence*, vol. 35, no. 1, pp. 221-231, 2013.
- [5] G. Potamianos, C. Neti, G. Gravier, A. Garg, and A. W. Senior, "Recent advances in the automatic recognition of audiovisual speech," *Proceedings of the IEEE*, vol. 91, no. 9, pp. 1306-1326, 2003.
- [6] K. Noda, Y. Yamaguchi, K. Nakadai, H. G. Okuno, and T. Ogata, "Audio-visual speech recognition using deep learning," *Applied Intelligence*, vol. 42, no. 4, pp. 722-737, 2015.
- [7] A. Krizhevsky, I. Sutskever, and G. E. Hinton, "ImageNet classification with deep convolutional neural networks," *Communications of the ACM*, vol. 60, no. 6, pp. 84-90, 2017.
- [8] D. Tran, L. Bourdev, R. Fergus, L. Torresani, and M. Paluri, "Learning spatiotemporal features with 3D convolutional networks," *Proceedings of the IEEE International Conference on Computer Vision*, pp. 4489-4497, 2015.