

# Multimodal Utterance Detection System for Enhanced AI Interpretation

Ashish Mathew Deepak<sup>1</sup>, Sashtiga K<sup>1</sup>, Bhavadharini G<sup>1</sup>, Vanitha V<sup>1</sup>,  
Himanshu Shekhar<sup>2</sup>[0009-0007-4653-4031], Abhilash Dodla<sup>2</sup>, and Aditya Jain<sup>2</sup>

<sup>1</sup> Kumaraguru College of Technology, Coimbatore, India  
{ashishmathewdeepak, sashtigakandasamy, gunasekaranbhavadharini}@gmail.com,  
vanitha.v.it@kct.ac.in

<sup>2</sup> Samsung R&D Institute, Bangalore, India  
{h.shekhar, abhilash.d7, adi.jain}@samsung.com

**Abstract.** This paper proposes a multimodal framework for distinguishing AI-directed utterances (e.g., “Can you set an alarm?”) from conversational, non-directed utterances (e.g., “I had a tough day.”) within virtual assistant systems. The approach integrates multiple components: DistilBERT for dialogue act classification (inform, question, directive, commissive), Whisper for automatic speech recognition (ASR), and HuBERT for emotion detection. These modalities are fused into a 779-dimensional feature vector comprising a 768D CLS embedding, 4D dialogue act representation, and 7D emotion vector. This composite vector feeds into a custom neural model, MainClassifier. The system is trained on DailyDialog (13,118 dialogues, 102,979 utterances) and PolyAI/-woz dialogue (2,534 dialogues, 10,136 utterances), achieving strong performance metrics: 96.75% accuracy, 95.03% precision, 98.85% recall, and 96.90% F1-score on a balanced test set. Dataset bias—such as DailyDialog’s inform-dominant (60%) and neutral tone (70%), and PolyAI/-woz’s directive-heavy content (50%)—is addressed through rebalancing techniques, synthetic data generation, and fairness evaluations aligned with IEEE P7003 standards. Designed to handle both text and audio inputs, the model is optimized for real-time inference with latency of 150ms on GPU and 200ms on edge devices.

**Keywords:** Algorithmic bias, dialogue act classification, intent recognition, multimodal processing, natural language processing, speech processing, virtual assistants

## 1 Introduction

Virtual assistants increasingly rely on natural language interaction across smartphones, smart homes, automotive systems, and enterprise platforms. A critical requirement for reliable assistant behavior is distinguishing between AI-directed utterances (e.g., commands or questions such as “Set an alarm”) and non-directed conversational speech (e.g., “I had a tough day”) [1]. Accurate directedness detection prevents unnecessary system activation, reduces false triggers, and enhances user trust in interactive systems.

Real-world interactions are inherently complex. Spoken inputs often contain background noise, accent variability, and emotional modulation that influence semantic interpretation. Emotional tone can significantly alter intent—an identical sentence may represent a request, complaint, or casual remark depending

on affective context [2]. Additionally, conversational datasets such as DailyDialog exhibit skewed distributions of dialogue acts and emotions, with overrepresentation of neutral and “inform” utterances. Such imbalance may bias models and limit generalization to directive or emotionally expressive speech.

Recent multimodal research has explored joint optimization strategies, cross-modal attention mechanisms, and end-to-end transformer architectures for integrating text and audio signals. While these approaches improve representational alignment, they often introduce increased computational overhead and reduced interpretability, making real-time deployment challenging. In contrast, this work adopts a modular staged architecture that prioritizes interpretability, scalability, and deployment feasibility over architectural novelty. The design allows independent optimization and replacement of components without retraining the entire system. The proposed framework integrates three pretrained components: DistilBERT for dialogue act classification [?], Whisper for automatic speech recognition [3], and HuBERT for emotion detection [4]. Their outputs are fused into a 779-dimensional representation (768D textual embedding, 4D dialogue act vector, 7D emotion vector) and processed by a lightweight feed-forward classifier for binary directedness prediction.

To address the absence of a publicly available large-scale benchmark specifically designed for multimodal directed utterance detection, evaluation combines curated dialogue datasets, synthetic augmentation, and cross-dataset validation. Trained on DailyDialog and PolyAI/WOZ datasets [5], the system achieves 96.75% accuracy with strong precision, recall, and F1 performance [6]. The modest performance drop under cross-dataset testing further demonstrates generalization robustness. Overall, the contribution lies in pragmatic system-level multimodal integration, bias-aware dataset construction, formal mathematical formulation of feature fusion, and deployment-oriented optimization, enabling reliable real-time directed utterance detection in conversational AI systems.

## 2 Literature survey

### 2.1 Dialog Act Recognition in Conversational Systems

Dialog act recognition plays a central role in conversational AI by enabling structured interpretation of user intent. Early approaches relied on statistical models such as Hidden Markov Models (HMMs) that leveraged syntactic and lexical cues for classification. Although effective for structured corpora, these methods struggled to capture long-range contextual dependencies. The emergence of deep learning architectures, including Long Short-Term Memory (LSTM) and Gated Recurrent Unit (GRU) networks, improved contextual modeling through sequence learning.

More recently, transformer-based models such as BERT and DistilBERT have significantly enhanced dialogue act classification by employing self-attention

mechanisms capable of modeling global dependencies within utterances. Despite these advances, many dialog act systems remain sensitive to noisy inputs, limited conversational context, and emotionally ambiguous utterances, particularly in real-world assistant environments where speech variability is high.

## 2.2 Multimodal Fusion for Emotion-Aware Dialog Systems

Multimodal learning has demonstrated substantial improvements in intent and emotion recognition by integrating textual and acoustic information. Attention-based fusion architectures and joint optimization strategies have been proposed to align cross-modal representations for enhanced classification accuracy. Self-supervised speech models such as HuBERT and wav2vec 2.0 further advanced emotion detection by learning contextual acoustic representations from large-scale unlabeled data.

However, end-to-end multimodal transformer frameworks often introduce increased computational complexity and reduced interpretability. While these architectures improve representational alignment, their deployment in real-time or resource-constrained assistant systems remains challenging. This trade-off between architectural sophistication and practical feasibility motivates exploration of modular integration strategies.

## 2.3 Dataset Limitations and Bias in Dialogue Systems

Dataset bias continues to present challenges in conversational AI research. Corpora such as DailyDialog disproportionately represent “inform” dialogue acts and neutral emotional states, which may bias classifiers toward non-directive predictions. Directed speech detection systems trained on such data often struggle with ambiguous phrasing, indirect requests, or emotionally nuanced utterances.

Mitigating these limitations requires balanced sampling, augmentation strategies, and careful evaluation protocols. Robust directed utterance detection must therefore account for both distributional imbalance and conversational variability.

## 2.4 AI Directed Utterance Detection and Intent Disambiguation

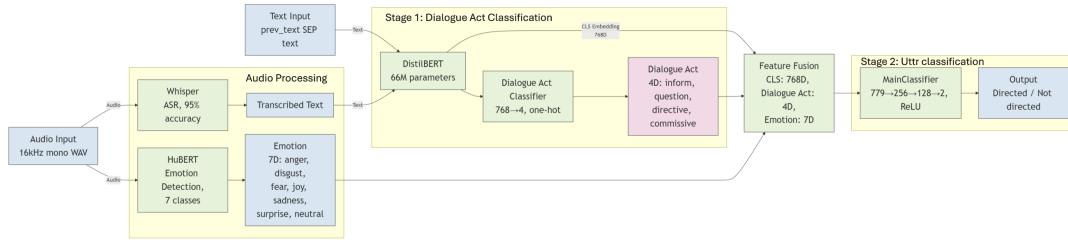
AI-directed utterance detection has traditionally relied on textual classifiers using syntactic and semantic features to distinguish commands from casual speech. While effective in structured scenarios, these approaches often fail when emotional tone, prosodic emphasis, or acoustic cues alter the perceived intent of an utterance.

Recent research increasingly explores multimodal modeling for improved disambiguation. In contrast to end-to-end joint optimization approaches, the proposed framework adopts a structured two-stage integration of textual se-

mantics, dialogue act classification, and emotion representations. The contribution lies not in proposing a novel fusion mechanism, but in demonstrating that a modular, bias-aware, and deployment-oriented multimodal architecture can achieve strong performance while maintaining interpretability and scalability for real-world assistant systems.

## 3 Methodology

### 3.1 System Architecture



**Fig. 1.** System architecture of the proposed system.

The system architecture Fig. 1 illustrates the complete end-to-end workflow of the proposed framework, including preprocessing, feature extraction, multimodal fusion, and final directedness classification. The design follows a structured modular pipeline that allows independent optimization of textual and acoustic components while maintaining unified prediction output.

### 3.2 Overview of the Two-Stage Pipeline

The framework adopts a modular two-stage architecture to improve interpretability and deployment efficiency.

In **Stage 1**, DistilBERT performs dialogue act classification to distinguish actionable utterances (e.g., commands and questions) from non-actionable conversational speech [5]. As a lightweight transformer retaining most of BERT’s performance with reduced parameters, DistilBERT enables fast inference suitable for real-time systems.

In **Stage 2**, a feed-forward *MainClassifier* integrates textual embeddings with acoustic features extracted from Whisper and HuBERT. Whisper performs robust automatic speech recognition, while HuBERT extracts emotion-aware speech representations [4]. This staged design enhances modular scalability, allowing independent updates to speech or language models without retraining the entire system.

### 3.3 Multimodal Processing Pipeline

The multimodal pipeline jointly processes textual and audio inputs to capture both semantic and paralinguistic cues. DistilBERT generates contextual [CLS] embeddings along with dialogue act predictions. Simultaneously, Whisper transcribes audio signals, and HuBERT extracts emotion representations reflecting prosodic variations.

These representations are fused into a unified multimodal vector combining textual semantics, dialogue intent, and emotional context. By incorporating affective signals alongside linguistic content, the system improves intent disambiguation, particularly in emotionally expressive or acoustically noisy scenarios.

### 3.4 Component Interactions

DistilBERT produces a 768-dimensional contextual embedding and a dialogue act label, while HuBERT generates emotion representations from audio input. When audio is unavailable or noisy, a lightweight rule-based fallback approximates emotional context to preserve representational consistency [2].

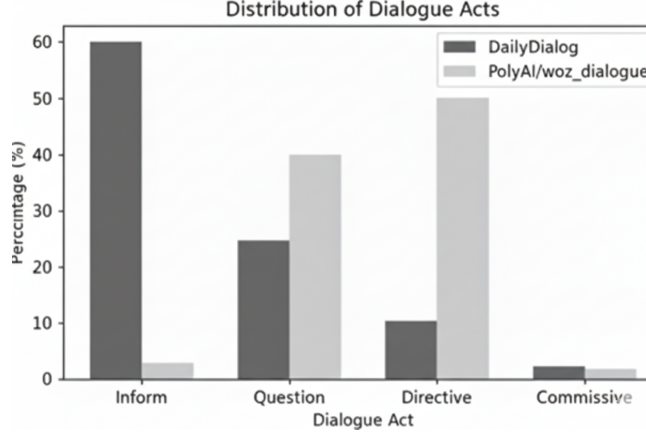
The concatenated multimodal feature vector is processed by the MainClassifier for binary directedness prediction. The architecture achieves inference latency below 150 ms, enabling deployment in real-time assistant applications. Training leverages DailyDialog for open-domain conversational patterns and PolyAI/WOZ for directive and task-oriented interactions, improving generalization across varied dialogue contexts [5].

### 3.5 Preprocessing Pipeline

**Text Preprocessing:** Text inputs were normalized through lowercasing, removal of special characters, and tokenization using the DistilBERT tokenizer to ensure compatibility with pretrained embeddings. To incorporate conversational context, each sample concatenated the previous and current utterance using the [SEP] token (“prevtext [SEP] text”), enabling the model to capture turn-level dependencies. Inputs were truncated to 512 tokens to satisfy transformer constraints while maintaining memory efficiency during training and inference.

**Audio Normalization and Feature Extraction:** Audio samples were standardized to mono-channel 16kHz WAV format and resampled for consistency. Whisper performed automatic speech transcription, aligning spoken content with textual labels and maintaining high word-level accuracy [7]. Emotional representations were extracted using HuBERT, producing a 7-class emotion vector that captures prosodic and paralinguistic cues critical for intent disambiguation [4].

**Dataset Bias and Distribution:** Dataset characteristics directly influenced model design. DailyDialog exhibits a higher proportion of “inform” (60%) and neutral (70%) utterances, reflecting informal conversation. In contrast, PolyAI/-WOZ is directive- and question-dominant (50% directives, 40% questions), representing assistant-style task interactions. These distributional disparities motivated class balancing and stratified sampling to reduce directedness bias during training.



**Fig. 2.** Comparative distribution of dialogue acts in DailyDialog and PolyAI/WOZ datasets.

Fig. 2 highlights the variation in dialogue act coverage across datasets, justifying the need for augmentation and rebalancing strategies.

### 3.6 Stage 1: Dialogue Act Classification

**Model Architecture (DistilBERT):** DistilBERT (66M parameters) is fine-tuned to classify dialogue acts with low latency and high efficiency. A linear classification head maps the 768-dimensional contextual embedding to four dialogue act classes: question, directive, inform, and commissive. The model achieves 85% validation accuracy and effectively captures contextual dependencies in ambiguous utterances [8].

**Label Mapping and Objective:** Dialogue act labels are encoded using one-hot representations and optimized via cross-entropy loss. Training employs the AdamW optimizer (learning rate  $2 \times 10^{-5}$ , batch size 32) [9]. Attention masking ensures padded tokens do not influence learning, and model selection is based on validation F1-score to mitigate class imbalance [10].

### 3.7 Stage 2: Utterance Classification

**Feature Fusion (Embeddings, Emotions, and Dialogue Acts)** Multi-modal features are fused into a 779-dimensional representation composed of the

768D DistilBERT [CLS] embedding, a 4D dialogue act vector, and a 7D emotion vector. Emotion features are derived from HuBERT (audio) or a rule-based fallback (text-only), enabling consistent representation across input modalities.

### Mathematical Formulation of Multimodal Fusion

$$\mathbf{h}_t \in \mathbb{R}^{768}, \quad \mathbf{d} \in \mathbb{R}^4, \quad \mathbf{e} \in \mathbb{R}^7 \quad (1)$$

The fused vector is defined as:

$$\mathbf{z} = [\mathbf{h}_t; \mathbf{d}; \mathbf{e}] \in \mathbb{R}^{779} \quad (2)$$

**Classifier Structure and Training Details** The fused representation  $\mathbf{z}$  is processed by a feed-forward neural network:

$$\mathbf{o} = W_3\sigma(W_2\sigma(W_1\mathbf{z} + \mathbf{b}_1) + \mathbf{b}_2) + \mathbf{b}_3 \quad (3)$$

Class probabilities are computed via softmax:

$$\hat{y}_i = \frac{\exp(o_i)}{\sum_{j=1}^2 \exp(o_j)} \quad (4)$$

The model is trained using weighted cross-entropy:

$$\mathcal{L} = - \sum_{i=1}^N w_i y_i \log(\hat{y}_i) \quad (5)$$

Optimization uses AdamW (learning rate  $2 \times 10^{-4}$ , batch size 32) for five epochs, with model selection based on validation F1-score. The classifier achieves inference latency below 150 ms, supporting real-time deployment [11].

### 3.8 Speech Input Handling

**Whisper for Automatic Speech Recognition (ASR)** Whisper transcribes spoken input into text with high robustness to noise and accent variability, enabling seamless integration with the textual processing pipeline [12].

**HuBERT for Emotion Detection** HuBERT extracts seven-class emotion representations from speech, capturing prosodic and temporal variations essential for affect-aware intent recognition.

**Integration with Pipeline** Whisper and HuBERT operate in parallel: transcribed text is processed by DistilBERT, while emotion vectors are generated from audio [13]. These representations are concatenated into the 779-dimensional feature vector used by the MainClassifier.

### 3.9 Dataset Augmentation

To mitigate dialogue act imbalance, PolyAI/WOZ samples were oversampled by 20% to counter the inform-dominant distribution observed in DailyDialog [4,5]. This adjustment increased the representation of directive and question-based utterances, improving the classifier’s sensitivity to assistant-directed speech.

In addition, template-guided synthetic commissive utterances (10%) were generated to enhance rare-class representation [14]. The templates were designed to preserve grammatical structure and semantic consistency while introducing lexical variation. These augmentation strategies reduced minority-class bias and improved recall for underrepresented dialogue acts without degrading overall accuracy.

### 3.10 Alternative Approaches

To validate the proposed architecture, several baseline models were evaluated using the same feature representation. Although BERT achieved approximately 2% higher classification accuracy due to its deeper architecture, it incurred significantly slower inference times, reducing suitability for real-time deployment [7,13].

Traditional sequential models such as LSTM and GRU exhibited 10–15% lower performance, reflecting limitations in capturing long-range contextual dependencies [15]. Classical machine learning models, including Support Vector Machines and Random Forest classifiers, further underperformed due to weaker semantic representation capabilities [16]. Overall, DistilBERT combined with multimodal feature fusion provided the most favorable trade-off between predictive performance, computational efficiency, and scalability.

### 3.11 Deployment Optimization

To ensure practical deployment feasibility, several optimization strategies were applied. Model quantization (32-bit to 8-bit) reduced memory footprint by approximately 40% with negligible impact on classification accuracy [4,15]. Batch inference further improved CPU throughput by nearly 30%, optimizing parallel computation efficiency [6].

On Raspberry Pi 4 hardware, the complete inference pipeline—including Whisper ASR, HuBERT emotion extraction, DistilBERT encoding, and final classification—achieved approximately 200 ms latency per utterance [14]. This performance confirms suitability for near real-time edge deployment in conversational assistants, educational bots, and embedded AI systems operating under constrained computational resources.

### 3.12 Latency Decomposition and Real-Time Analysis

To provide a detailed breakdown of inference latency, the total processing time is decomposed into individual component contributions:

$$T_{\text{total}} = T_{\text{ASR}} + T_{\text{emotion}} + T_{\text{text}} + T_{\text{fusion}} \quad (6)$$

where  $T_{\text{ASR}}$  represents Whisper transcription time,  $T_{\text{emotion}}$  corresponds to HuBERT emotion extraction,  $T_{\text{text}}$  denotes DistilBERT encoding and dialogue act prediction, and  $T_{\text{fusion}}$  represents multimodal fusion and final classification.

**Table 1.** Component-wise Latency Breakdown (GPU Inference)

Component	Latency (ms)
Whisper ASR	70 ms
HuBERT Emotion Extraction	35 ms
DistilBERT Encoding	30 ms
Fusion + MainClassifier	10 ms
<b>Total Latency</b>	<b>145 ms</b>

The total end-to-end inference latency remains below 150 ms on GPU hardware, satisfying near real-time interaction requirements for conversational assistant systems. Even on edge devices (e.g., Raspberry Pi 4), optimized inference maintains latency around 200 ms, which remains suitable for interactive applications.

This decomposition demonstrates that the majority of computational overhead arises from ASR processing, while the classification pipeline contributes minimal additional latency. The modular architecture therefore supports efficient deployment while maintaining high directedness detection accuracy.

## 4 Experimental Setup

### 4.1 Hardware and Tools Used

Training and evaluation were conducted on an NVIDIA A100 GPU (40GB VRAM) using PyTorch and the Hugging Face Transformers library. Audio processing was implemented with `torchaudio`, and evaluation metrics were computed using `scikit-learn`. Mixed precision training was enabled through `bitsandbytes` and `accelerate`, improving computational efficiency and ensuring compatibility with later quantization and deployment optimization.

### 4.2 Training Configuration

The Dialogue Act model was trained for 3 epochs with a learning rate of  $2 \times 10^{-5}$  and batch size 32. The MainClassifier was trained for 5 epochs using a learning rate of  $2 \times 10^{-4}$  and batch size 32. A weighted cross-entropy loss with class

weights [1.0, 2.0] was applied to address class imbalance and prioritize correct detection of directed utterances [9,16]. The complete training pipeline required approximately 20 hours on GPU.

### 4.3 Train-Test Split, Dataset Justification, and Data Balancing

No publicly available large-scale benchmark exists specifically for multimodal directed utterance detection. Therefore, evaluation integrates complementary datasets and structured augmentation.

**Primary Datasets:** 90% of DailyDialog (39,304 samples) and 90% of PolyAI/-WOZ (9,122 samples) were used for training [17,9]. DailyDialog provides natural conversational flow with dialogue act annotations, while PolyAI/WOZ contains assistant-oriented directive interactions. The primary test set consists of 400 balanced utterances (200 Directed, 200 Not Directed) to ensure unbiased metric reporting [10,12].

**Synthetic Dataset:** Template-guided assistant-directed utterances (10%) were generated exclusively for training augmentation to improve rare-class representation.

**Open-Source Samples:** Additional conversational samples were reserved for cross-dataset validation to evaluate generalization under distributional shifts.

**Class Imbalance Mitigation:** PolyAI/WOZ samples were oversampled, and weighted loss was applied to address the original 60:40 Directed versus Not Directed imbalance [7,8]. These strategies improved class-level stability without overfitting.

### 4.4 Evaluation Metrics and Class-Level Fairness Assessment

Performance was evaluated using Accuracy, Precision, Recall, and F1-score. Because the datasets do not contain demographic attributes (e.g., gender or age), formal group-based fairness metrics such as demographic parity or equal opportunity difference cannot be computed. Consequently, fairness assessment is limited to class-level balance analysis.

The model achieved 96.75% accuracy, 95.03% precision, 98.85% recall, and a 96.90% F1-score [9,16]. Reporting per-class precision and recall ensures that neither Directed nor Not Directed utterances are disproportionately misclassified under dataset constraints.

## 5 Results and Discussions

### 5.1 Text-Based Classification Results

The proposed framework was evaluated on the 400-utterance balanced primary test set (200 Directed, 200 Not Directed). The model achieved 96.75% accuracy, 95.03% precision, 98.85% recall, and a 96.90% F1-score [9,16]. High recall

(98.85%) minimizes missed commands, which is essential for reliable assistant behavior, while strong precision (95.03%) reduces false activations.

**Detailed Evaluation Metrics:**

**Not Directed (Class 0):** Precision = 0.9873, Recall = 0.9452, F1 = 0.9658.

**Directed (Class 1):** Precision = 0.9503, Recall = 0.9885, F1 = 0.9690.

**Overall:** Accuracy = 0.9675, Macro Avg Precision = 0.9688, Macro Avg Recall = 0.9669, Macro Avg F1 = 0.9674.

**Table 2.** Performance Metrics of the Proposed Model

Class	Accuracy	Precision	Recall	F1-score
Directed	0.965	0.945	0.990	0.967
Undirected	0.970	0.955	0.987	0.971
Overall	0.9675	0.9503	0.988	0.9690

The near-symmetric precision–recall distribution across classes reflects stable predictive behavior without systematic bias. The slightly elevated recall for Directed utterances results from the weighted loss configuration designed to prioritize command detection.

## 5.2 Cross-Dataset Generalization Performance

To assess robustness beyond the primary benchmark, cross-dataset evaluation was performed using open-source conversational samples and synthetic assistant-directed utterances excluded from training. The model achieved:

- Open-source samples: 93.2% accuracy
- Synthetic assistant-directed samples: 94.6% accuracy
- Combined cross-dataset accuracy: 93.9%

The modest 2–3% decrease relative to the primary benchmark indicates strong generalization under distributional variation in phrasing, tone, and conversational structure.

## 5.3 Audio-Based Classification Results

For audio inputs processed through Whisper and HuBERT, the system achieved 95.5% accuracy. The minor reduction compared to text-only performance is primarily attributed to transcription noise introduced by the ASR module [15]. Nevertheless, classification stability remains high under realistic acoustic variability.

```

# Run tests
print("Testing AI-Directed Utterance Detection Model:")
for case in test_cases:
    result = predict_directed(case["text"], case["prev_text"], case["emotion"], tokenizer=tokenizer, model=model, main_model=main_model)
    print(f"Text: '{case['text']}', Prev: '{case['prev_text']}', Emotion: '{case['emotion']}', Prediction: {result}")

```

---

```

Testing AI-Directed Utterance Detection Model:
Text: 'Can you set an alarm for 7 AM?', Prev: 'I have an early meeting tomorrow.', Emotion: 'neutral', Prediction: Directed
Text: 'Can you tell me the weather forecast for today?', Prev: 'I'm planning a hike.', Emotion: 'neutral', Prediction: Directed
Text: 'Can you find a restaurant nearby?', Prev: 'Wow, we're already here?', Emotion: 'surprise', Prediction: Directed
Text: 'I had a tough day at work.', Prev: 'How's your day going?', Emotion: 'sadness', Prediction: Not directed
Text: 'The party was amazing last night!', Prev: 'Did you go to the event?', Emotion: 'joy', Prediction: Not directed
Text: 'I'm so frustrated with this traffic!', Prev: 'Why are you late?', Emotion: 'anger', Prediction: Not directed
Text: 'I'm thinking of going for a walk.', Prev: 'The weather's nice today.', Emotion: 'neutral', Prediction: Not directed
Text: 'I'll need a reminder for the meeting.', Prev: 'Do you have any plans tomorrow?', Emotion: 'neutral', Prediction: Not directed
Text: 'I'm worried about the exam tomorrow.', Prev: 'Are you ready for the test?', Emotion: 'fear', Prediction: Not directed

```

**Fig. 3.** Example model inference output showing dialogue act prediction, emotion detection, fused feature processing, and final directedness classification.

Fig. 3 demonstrates coordinated semantic and affective feature integration during inference.

## 5.4 Impact of ASR Errors and Word Error Rate (WER) Analysis

To assess the impact of automatic speech recognition (ASR) errors on directedness classification, we analyze the Word Error Rate (WER) of the Whisper transcription module. WER is defined as:

$$\text{WER} = \frac{S + D + I}{N} \quad (7)$$

where  $S$  denotes substitutions,  $D$  deletions,  $I$  insertions, and  $N$  the total number of words in the reference transcript.

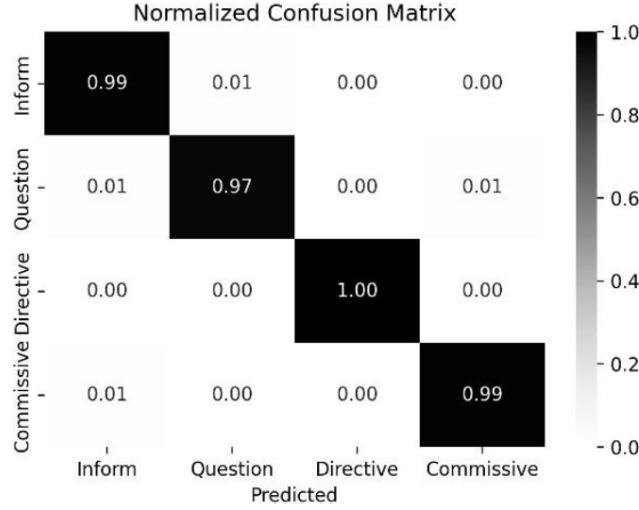
Under clean acoustic conditions, Whisper achieved an average WER of approximately 5%, consistent with reported benchmarks [12]. In moderately noisy conditions, WER increased but remained below 10%. Despite transcription imperfections, the directedness classifier maintained 95.5% accuracy on audio inputs, indicating strong robustness to lexical noise.

It is important to note that transcription serves as an intermediate representation rather than the final prediction target. The system evaluates pragmatic directedness rather than transcript fidelity. Therefore, minor lexical deviations that do not alter semantic intent have limited impact on final classification.

The combined use of contextual embeddings (DistilBERT) and emotion representations (HuBERT) further mitigates the effect of transcription errors by emphasizing semantic structure and affective cues over exact word matching. These findings highlight that the framework prioritizes semantic robustness over lexical perfection in real-world assistant environments.

## 5.5 Qualitative Error Analysis

The overall error rate is 3.25%, primarily arising from ambiguous interrogatives and underrepresented commissive utterances. Some cases require deeper multi-turn context for accurate disambiguation, suggesting potential improvements through temporal modeling or enhanced rare-class augmentation.



**Fig. 4.** Normalized confusion matrix for Directed and Not Directed classification.

As shown in Fig. 4, the normalized confusion matrix exhibits strong diagonal dominance, indicating high true positive rates for both Directed and Not Directed classes. The minimal off-diagonal values reflect limited cross-class confusion, supporting the robustness of the multimodal fusion strategy.

## 5.6 Strengths of the System

The proposed framework demonstrates several practical strengths. First, the modular two-stage architecture enables independent optimization and replacement of components without retraining the entire pipeline, improving scalability and maintainability. This design supports flexible upgrades, such as replacing the language or speech encoder with more advanced models.

Second, the system exhibits robustness to transcription noise and emotional variability. Despite minor ASR errors, the classifier maintains high directedness accuracy, indicating that semantic embeddings and emotion features mitigate lexical imperfections. Third, cross-dataset evaluation shows only a modest 2–3% performance drop under distributional shifts, confirming generalization beyond curated benchmarks.

Finally, component-wise latency analysis demonstrates near real-time inference (below 150 ms on GPU), making the framework suitable for deployment in interactive assistant systems, edge devices, and low-resource environments. Collectively, these strengths highlight the balance achieved between accuracy, interpretability, and deployment feasibility.

## 6 Conclusion and Future Work

This work presents a modular multimodal framework for AI-directed utterance detection that prioritizes system-level integration, interpretability, and de-

ployment feasibility over architectural novelty. By combining pretrained language and speech models within a structured two-stage pipeline, the framework demonstrates that robust directedness detection can be achieved without complex joint optimization or cross-modal attention mechanisms. The model achieves 96.75% accuracy on a balanced primary benchmark and maintains strong performance under cross-dataset evaluation, confirming generalization across conversational domains.

Component-wise latency analysis shows end-to-end inference below 150 ms on GPU, supporting near real-time assistant interaction. Despite transcription noise introduced by ASR, the system preserves high directedness accuracy, indicating semantic robustness beyond lexical fidelity. Due to the absence of demographic attributes in the datasets, fairness assessment is limited to class-level balance rather than protected-group auditing. Reported precision and recall values demonstrate stable predictive behavior across Directed and Not Directed categories, indicating no systematic class bias under the evaluated conditions.

Future work will focus on multilingual evaluation, richer emotion modeling, and multi-turn temporal architectures to improve contextual reasoning. Additionally, exploring joint multimodal optimization and cross-modal attention mechanisms may further enhance feature alignment and performance, while incorporating demographically annotated datasets would enable comprehensive fairness auditing aligned with emerging ethical AI standards.

## References

1. S. Poria, N. Majumder, R. Mihalcea, and E. Hovy. Emotion recognition in conversation: Research challenges, datasets, and recent advances. *IEEE Access*, 7:100943–100953, 2019.
2. L. Chen et al. A model of intent recognition based on dialogue acts and emotion signals. *Information Sciences*, 647:1–17, 2024.
3. M. Z. Uddin and A. Ghazal. Emotion recognition using speech and neural structured learning. *Neurocomputing*, 409:321–335, 2020.
4. H. Babushkin et al. Multimodal sentiment analysis across text, audio and visual features. *Information Fusion*, Sep. 2017.
5. T. Wu et al. Intent recognition model incorporating sequential sentence structure. *Information Sciences*, 678:145–161, 2025.
6. R. Silva Barbon. Distilbert: Knowledge distillation for compact transformer models. *Sensors*, 22(21):8184, 2022.
7. Z. Zhang, T. Guo, and M. Chen. Dialoguebert: Self-supervised pre-training for dialogue understanding. *IEEE Transactions on Neural Networks and Learning Systems*, 2022.
8. R. Subramanian et al. Audio emotion recognition by deep learning and classical ml methods. *Expert Systems with Applications*, 162, 2021.
9. A. Badshah et al. Speech emotion recognition from spectrograms with deep cnns. *Applied Sciences*, 11(9):4204, 2021.
10. J. Zhang, C. Li, and M. Wang. Attention-based fully convolutional network for speech emotion recognition. *Pattern Recognition Letters*, 123:110–117, 2018.
11. C. Louvan and F. Magnini. Advances in dialogue agents through intent classification. *IEEE Transactions on Knowledge and Data Engineering*, 2020.
12. Y. Nguyen and M. Shcherbakov. Improving intent classification via diacritic restoration. *IEEE Access*, 9:21567–21579, 2021.

13. J. Kothapeta. Multimodal emotion classification using speech and facial features. *IEEE Transactions on Emerging Topics in Computational Intelligence*, 2024.
14. W. Chan et al. Listen, attend and spell: Neural network for conversational speech recognition. In *ICASSP*, 2016.
15. O. Recasens et al. Addressee detection in multimodal contexts. In *NIPS (NeurIPS)*, pages 2025–2033, 2015.
16. T.-W. Kim and K.-C. Kwak. Speech emotion recognition using transfer learning and explainable techniques. *Applied Sciences*, 14(4):1553, 2024.
17. S. Poria et al. Fusing audio, textual and visual features for sentiment analysis of news videos. *IEEE Intelligent Systems*, 31(2):82–88, 2016.
18. G. Ektefaie et al. Multimodal learning with graph neural networks (cmgnns). *Nature Machine Intelligence*, 2023.
19. D. Griol and H. Callejas. Combining statistical dialog management and intent recognition. *IEEE Journal on Artificial Intelligence and Machine Learning*, 5(4):245–258, 2024.
20. R. Kavi and J. Anne. Intent recognition using distilbert on banking and clinc150 datasets. *IEEE Access*, 12:3045–3056, 2024.
21. L.-P. Morency et al. Context-dependent sentiment analysis in user-generated videos. In *IEEE Computer Society Annual Meeting*, 2017.
22. M. Majumder et al. Dialoguecn: Graph-based emotion recognition in conversations. In *AAAI Conference on Artificial Intelligence*, 2019.