ISSN: 2998-2383

Vol. 4, No. 10, 2025

A Comprehensive Survey on Multilingual and Multimodal Automatic Speech Recognition Systems

Maelle Trenton

Wright State University, Dayton, Ohio, USA mnt984@wright.edu

Abstract: Speech recognition has evolved from rule-based systems to deep learning—driven architectures that can comprehend complex linguistic and acoustic patterns. However, the rapid progress achieved in high-resource languages, such as English and Mandarin, has not been equally reflected in multilingual and low-resource contexts. This disparity limits the inclusiveness and global applicability of automatic speech recognition (ASR) technologies. In parallel, the emergence of multimodal learning-integrating speech, vision, and text-has opened new possibilities for robust and context-aware recognition systems that align with human communication patterns. This paper provides a comprehensive survey of recent advances in multilingual and multimodal speech recognition. It reviews state-of-the-art models, including end-to-end architectures, self-supervised learning, and transformer-based approaches such as wav2vec 2.0, Whisper, and SpeechT5. The review also explores multilingual pretraining strategies, transfer learning for low-resource adaptation, and multimodal fusion techniques that combine audio with visual or textual modalities to enhance recognition accuracy and robustness. Moreover, we analyze benchmark datasets, evaluation metrics, and key challenges such as code-switching, domain adaptation, and cultural diversity. Finally, we highlight future trends in cross-lingual model generalization, data-efficient learning, and multimodal interaction for next-generation intelligent speech systems. The findings indicate that progress in multilingual and multimodal ASR is essential to bridge the linguistic divide and to achieve equitable access to AI-driven technologies worldwide.

Keywords: Speech recognition, multilingual ASR, multimodal learning, low-resource languages, transformer models, self-supervised learning, cross-lingual adaptation.

1. Introduction

Speech is the most natural and efficient mode of human communication, and developing systems capable of understanding it has long been a core objective of artificial intelligence research. The field of automatic speech recognition (ASR) has witnessed remarkable evolution-from early Hidden Markov Model (HMM) and Gaussian Mixture Model (GMM) systems to contemporary deep learning-based frameworks leveraging convolutional, recurrent, transformer architectures. With the surge of large-scale data and computational resources, speech recognition has achieved human-level accuracy in several high-resource languages such as English, Mandarin, and Spanish. However, the global landscape of linguistic diversity, consisting of more than 7,000 languages, poses significant challenges for the inclusiveness of modern ASR systems.

The majority of current speech models remain biased toward languages with abundant annotated data, leaving low-resource languages underserved. This imbalance exacerbates the technological gap between communities, impeding the development of inclusive AI. To address this issue, recent research has increasingly focused on multilingual and low-resource speech recognition, which aims to generalize knowledge across linguistic boundaries through transfer learning, multilingual pretraining, and data-efficient adaptation. Models such as XLS-R (Babu et al., 2022) and Whisper (Radford et al., 2023) demonstrate that shared

representations across languages can significantly improve performance in resource-scarce settings.

At the same time, the rise of multimodal learning has reshaped how ASR systems interpret and reason about speech. Human communication rarely occurs in isolation-it involves visual cues, gestures, text, and environmental context. By integrating audio with other modalities such as vision or text, multimodal ASR systems enhance semantic comprehension, improve robustness under noise, and better align with real-world use cases such as video captioning, human-computer interaction, and assistive technologies. Models like AV-HuBERT (Shi et al., 2023) and SpeechT5-Multi (Tang et al., 2024) exemplify this shift toward multimodal fusion frameworks.

Despite these advances, several open challenges remain. Multilingual ASR still struggles with phonetic variation, codeswitching, and the scarcity of parallel resources for underrepresented languages. Similarly, multimodal systems face difficulties in synchronizing heterogeneous data streams and ensuring that visual or textual information enhances rather than confuses recognition. Moreover, ethical considerations such as linguistic equity, dataset bias, and privacy constraints must be addressed for global-scale deployment.

This survey aims to systematically review and synthesize the latest progress in multilingual and multimodal speech recognition. It highlights representative architectures, training paradigms, datasets, and evaluation methodologies. The paper also discusses the interplay between multilinguality and multimodality-how cross-lingual knowledge can be strengthened through multimodal supervision, and vice versa. By providing a unified perspective, this survey seeks to guide future research toward inclusive, data-efficient, and context-aware speech recognition systems that serve diverse linguistic communities.

The remainder of this paper is organized as follows. Section II reviews related work in traditional and deep learning—based ASR. Section III discusses multilingual ASR techniques and the challenges of cross-lingual generalization. Section IV focuses on low-resource speech recognition approaches, including transfer learning and self-supervised pretraining. Section V introduces multimodal ASR, highlighting integration strategies with vision and text. Section VI compares datasets and evaluation benchmarks. Section VII outlines ongoing challenges and emerging trends, and Section VIII concludes with future research directions.

2. Related work

The development of automatic speech recognition (ASR) has undergone a profound transformation over the past decades, evolving from early statistical systems to end-to-end deep learning models capable of capturing complex linguistic and acoustic structures. Early ASR frameworks were dominated by statistical approaches such as Hidden Markov Models (HMMs) and Gaussian Mixture Models (GMMs). These systems represented speech as sequences of acoustic observations modeled by probabilistic transitions, which allowed them to handle variability in speech duration and noise. However, such models depended heavily on handcrafted features and struggled with accent variation, speaking rate, and background interference. To improve robustness, discriminative training methods like Maximum Mutual Information (MMI) and feature-space transformations were introduced, but their performance gains were limited by the expressiveness of shallow architectures.

The advent of deep neural networks (DNNs) marked a turning point in ASR research. Hybrid DNN-HMM systems replaced the traditional GMM component with multi-layer perceptrons capable of modeling nonlinear relationships between input features and phonetic states. This innovation enabled large gains in accuracy and generalization, especially on benchmark datasets such as TIMIT and Switchboard. Nonetheless, these hybrid systems retained complex pipelines involving separate acoustic, language, and lexicon models. To simplify the architecture and reduce error propagation, researchers shifted toward end-to-end models that directly map audio inputs to textual outputs without intermediate phonetic alignment. Two dominant paradigms emerged: Connectionist Temporal Classification (CTC) and attentionbased sequence-to-sequence models. CTC introduced flexible alignment between variable-length sequences marginalizing over all possible output paths, while attentionbased encoder - decoder frameworks, such as Listen, Attend and Spell (LAS), leveraged dynamic attention mechanisms to align input frames and output tokens. These approaches

substantially simplified ASR pipelines and enabled more natural integration with other neural sequence models.

In parallel, transformer architectures brought a new level of efficiency and contextual modeling to ASR. Unlike recurrent neural networks that process inputs sequentially, transformers use self-attention to capture long-range dependencies in parallel. Models such as the Speech-Transformer and Conformer have demonstrated state-of-theart performance across multiple speech benchmarks. Their scalability and parallelizability make them ideal for large-scale multilingual pretraining. At the same time, the rise of selfsupervised learning (SSL) further revolutionized speech representation learning by enabling models to exploit vast amounts of unlabeled audio data. Methods such as wav2vec 2.0, HuBERT, and data-scaled extensions like XLS-R achieved remarkable generalization by pretraining on hundreds of thousands of hours of speech and then fine-tuning on specific languages or domains. These models significantly reduced the reliance on labeled data, opening the door to progress in low-resource languages that previously lacked large annotated corpora.

Building on this foundation, multilingual ASR has emerged as a critical research direction aiming to share knowledge across languages. Multilingual systems train a unified model on speech data from multiple languages, often using shared encoders and language-specific decoders or adapters. This approach allows the model to learn universal phonetic representations while preserving language-specific nuances. Cross-lingual transfer has also proven effective: pretraining on high-resource languages can improve recognition performance for low-resource ones through transfer learning or parameter-efficient fine-tuning. However, multilingual modeling introduces new challenges, such as negative transfer, where the performance on one language deteriorates when training jointly with others, and catastrophic forgetting during continual learning. To address these issues, researchers have explored adaptive fine-tuning strategies, dynamic masking, and modular architectures that selectively share or isolate linguistic features. Recent large-scale multilingual systems such as Whisper and SeamlessM4T exemplify this paradigm, achieving zero-shot recognition and translation across dozens of languages by combining massive multilingual pretraining with cross-modal supervision.

Meanwhile, the integration of multimodal information has expanded the scope of ASR beyond pure audio processing. Human communication is inherently multimodal, combining speech, vision, and context. Audio-visual speech recognition (AVSR) systems leverage visual cues such as lip movements to improve recognition accuracy in noisy environments or under occlusion. Transformer-based multimodal architectures like AV-HuBERT and SpeechT5-Multi integrate audio, text, and visual streams through shared attention mechanisms, allowing the model to ground speech understanding in contextual signals. Similarly, cross-domain pretraining that combines ASR with speech translation, speech-to-text summarization, or emotion recognition has enhanced generalization and semantic richness. The convergence of speech and natural language processing under unified frameworks, as seen in Whisper and SpeechT5, demonstrates

how joint pretraining across modalities and tasks can yield robust, flexible, and transferable models.

In summary, the evolution of ASR reflects a consistent trend toward unification, scalability, and inclusivity. From statistical to deep neural models, from supervised learning to self-supervised and multilingual paradigms, and from unimodal to multimodal architectures, the field has progressively reduced dependence on handcrafted components and language-specific resources. The next frontier lies in achieving balanced performance across diverse linguistic communities and environmental conditions. Multilingual and multimodal ASR systems hold the potential to close the global accessibility gap by enabling speech understanding for all users, regardless of language, accent, or context. The following sections examine these developments in depth, focusing on techniques for multilingual speech recognition, strategies for low-resource adaptation, and the role of multimodal integration in advancing the next generation of speech recognition technologies.

3. Multilingual Speech Recognition: Techniques and Challenges

Multilingual speech recognition aims to develop a unified model that can recognize and transcribe speech from multiple languages with minimal degradation in performance. Unlike monolingual systems trained exclusively on one linguistic corpus, multilingual models must capture both the shared and language-specific aspects of speech. This balance is essential for enabling generalization across diverse linguistic and phonetic systems. Recent research has shown that multilingual ASR not only improves efficiency by consolidating multiple language models into a single framework but also enhances performance for low-resource languages through cross-lingual knowledge transfer. The fundamental premise is that acoustic and phonetic structures exhibit universal similarities that can be jointly learned from multilingual data.

Early approaches to multilingual ASR adopted a parameter-sharing strategy in hybrid DNN-HMM systems, where the hidden layers of the DNN were shared among languages while keeping the output layers language-specific. This design leveraged commonalities in acoustic features across languages but required separate decoders for each linguistic target, limiting scalability. With the introduction of end-to-end architectures, particularly those based on the encoder-decoder paradigm, researchers began exploring fully shared encoders with unified vocabularies or byte-pair encoding (BPE) tokenization across languages. These models, such as multilingual LAS and RNN-Transducer systems, demonstrated that shared subword units could capture crosslingual phonetic representations effectively. The transformerbased architectures, including Conformer and Speech-Transformer extensions, further improved multilingual modeling by enabling efficient attention-based alignment across long sequences and diverse linguistic contexts.

A major breakthrough came with large-scale multilingual pretraining frameworks such as XLSR-53 and XLS-R, which extended self-supervised models like wav2vec 2.0 to dozens or even hundreds of languages. These models are trained on massive unlabeled speech corpora, learning universal representations that can later be fine-tuned for specific tasks or languages. The pretraining objective-often based on contrastive or masked prediction-encourages the model to learn phonetic and prosodic patterns that generalize across linguistic boundaries. As a result, multilingual pretraining has become a cornerstone for low-resource adaptation, enabling models to perform competitively even when fine-tuned on minimal supervised data. Models such as Whisper by OpenAI and SeamlessM4T by Meta extend this concept to cross-modal and cross-task learning, jointly training on speech recognition, translation, and transcription to achieve zero-shot recognition across languages. These systems mark a paradigm shift toward universal ASR, where a single model can transcribe and translate multilingual audio without explicit retraining.

However, multilingual ASR introduces several inherent challenges. One of the most critical is negative transfer, where languages with dissimilar phonetic or lexical structures interfere with each other during training, leading to suboptimal performance. For instance, tonal languages like Mandarin may conflict with non-tonal ones when learned jointly without careful balancing. This problem is often exacerbated by data imbalance, as high-resource languages dominate the training distribution. Techniques such as data sampling, loss reweighting, and adaptive learning rates have been employed to mitigate such imbalances. Another issue arises from codeswitching, the phenomenon where speakers alternate between two or more languages within a single utterance. Traditional ASR systems struggle to handle this linguistic fluidity because of inconsistent language modeling and vocabulary overlap. Recent research introduces language identification conditioning, token-level language tags, and multilingual CTC loss designs to improve recognition under code-switching scenarios.

Beyond the acoustic and linguistic aspects, cultural and sociolinguistic diversity also complicates multilingual ASR. Variations in dialect, accent, and pronunciation patterns can cause degradation even within the same language family. To address this, accent-robust and dialect-aware adaptation strategies have been proposed, where speaker embeddings, phoneme normalization, or accent-specific adapters are introduced to preserve performance consistency. Another dimension of the challenge involves script diversitymultilingual ASR systems must often transcribe in different writing systems, from Latin and Cyrillic to Devanagari and Arabic. Subword-based tokenization methods such as SentencePiece help unify representation spaces, but they may lose orthographic fidelity in highly divergent languages. Researchers are exploring hierarchical vocabularies and grapheme-based modeling to maintain both efficiency and linguistic integrity.

Despite these complexities, multilingual ASR has become increasingly practical due to the scaling of both data and model capacity. The availability of large multilingual

corpora, such as CommonVoice, Multilingual LibriSpeech, and FLEURS, provides diverse speech sources covering hundreds of languages. Coupled with transformer backbones exceeding billions of parameters, these datasets enable powerful generalization and transfer learning. Yet, computational cost remains a barrier, particularly for low-resource communities lacking access to large-scale training infrastructure. Lightweight strategies such as adapter modules, parameter-efficient fine-tuning, and quantization-based compression are emerging to democratize multilingual ASR deployment.

In summary, multilingual ASR represents a significant step toward equitable access to speech technology. It provides a foundation for bridging linguistic divides and enabling communication across global communities. While universal models like Whisper and SeamlessM4T demonstrate promising zero-shot performance, achieving balanced accuracy across all languages-especially underrepresented ones-remains an open research problem. Future advances are expected to come from combining multilingual pretraining with multimodal supervision, where contextual cues from visual and textual domains enhance the model's understanding of diverse linguistic environments. This cross-pollination of multilingual and multimodal approaches forms the conceptual bridge to the next stage of ASR evolution, discussed in the following section.

4. Low-Resource Speech Recognition

The problem of low-resource speech recognition lies at the intersection of data scarcity, linguistic diversity, and model generalization. While state-of-the-art ASR systems have achieved near-human performance in high-resource languages. they often fail to deliver acceptable results for underrepresented languages with limited labeled corpora. These languages, which make up the majority of global linguistic diversity, often lack both annotated speech data and large-scale text corpora necessary for robust language modeling. The disparity in available resources not only restricts the deployment of ASR technologies in many regions but also deepens the technological divide between dominant and marginalized linguistic communities. Therefore, research on low-resource ASR focuses on designing data-efficient learning algorithms, cross-lingual transfer mechanisms, and self-supervised pretraining strategies that minimize the dependency on labeled data while preserving performance.

One of the most prominent solutions is transfer learning, where a model trained on high-resource languages is adapted to new low-resource targets. The intuition behind this approach is that the lower layers of deep models encode universal acoustic and phonetic features that are transferable across languages, while the upper layers can be fine-tuned to capture language-specific characteristics. For instance, multilingual models such as XLS-R and Whisper provide powerful pretrained encoders that can be fine-tuned with only a few hours of labeled data in a new language. Techniques such as adapter-based fine-tuning, lightweight

reparameterization, and layer-wise freezing have been proposed to reduce the computational and data requirements for adaptation. Moreover, cross-lingual subword modeling using shared byte-pair encodings has proven to be particularly effective, as it allows low-resource languages to benefit from overlapping phonetic or orthographic structures learned from high-resource languages.

Another key direction is data augmentation, which seeks to artificially increase data diversity through transformations applied to existing samples. Conventional methods such as noise injection, speed perturbation, and SpecAugment have long been used to improve robustness. However, recent studies extend these methods through generative modeling, using diffusion-based or GAN-based audio synthesis to produce realistic speech samples for low-resource languages. This augmentation not only expands training data but also enhances the model's capacity to handle varied acoustic conditions. Additionally, cross-lingual pseudo-labeling-where unlabeled audio from the target language is transcribed using a highresource teacher model-has become a central technique in semi-supervised learning. Iterative self-training cycles, where the model refines its predictions through repeated pseudolabeling, have demonstrated notable gains for languages with extremely limited annotation.

Self-supervised learning (SSL) has emerged as the most transformative paradigm for low-resource speech recognition. Instead of relying on manual transcriptions, SSL models learn generalizable speech representations by predicting masked or contrastive features from raw audio. Frameworks such as wav2vec 2.0, HuBERT, and data-scaled versions like XLSR-128 enable universal pretraining across hundreds of thousands of hours of multilingual speech data. These representations capture rich phonetic and prosodic information that can be fine-tuned on small labeled datasets, dramatically improving performance in low-resource contexts. The use of crosslingual self-supervised objectives, where the model is exposed to multiple languages during pretraining, enhances its ability to generalize to unseen languages. Furthermore, recent models like SpeechT5 and SeamlessM4T unify pretraining across multiple tasks, such as speech-to-text, speech translation, and speech generation, demonstrating strong zero-shot performance in low-resource languages.

However, the low-resource setting introduces unique challenges beyond data availability. The absence of standardized orthography, dialectal variation, and noisy recordings often make it difficult to define consistent training objectives. Low-resource languages may also lack reliable text normalization and tokenization tools, further complicating data preprocessing. Moreover, the imbalance between source and target languages in multilingual training can lead to negative transfer, where dominant language features overshadow those of minority languages. Researchers are increasingly turning to meta-learning and few-shot adaptation approaches, which train models to rapidly learn new languages with minimal supervision. Prototypical networks and modelagnostic meta-learning (MAML) have been adapted for ASR,

enabling systems to generalize from only a handful of labeled samples per language.

Recent efforts also focus on improving linguistic inclusivity through community-driven datasets. Projects such as CommonVoice, Masakhane, and FLEURS aim to crowdsource recordings and transcriptions from diverse linguistic groups, providing open and scalable resources for low-resource ASR research. Combined with multilingual self-supervised pretraining, these initiatives are enabling models to reach unprecedented coverage across languages, including those previously considered technologically invisible. Moreover, policy-level initiatives by organizations like UNESCO emphasize the importance of preserving linguistic heritage through speech technology, aligning low-resource ASR with cultural sustainability and digital equity.

In conclusion, the field of low-resource speech recognition is rapidly advancing through the synergy of transfer learning, data augmentation, and self-supervised pretraining. The long-standing bottleneck of insufficient labeled data is being mitigated by representation learning and generative augmentation, allowing the extension of ASR systems to hundreds of languages worldwide. Yet, challenges remain in ensuring equitable model performance, addressing dialectal variation, and reducing computational barriers to entry for local researchers. As the next section discusses, integrating multimodal learning offers a promising path toward addressing these limitations by grounding speech recognition in visual and contextual information, thereby enhancing robustness, comprehension, and inclusivity across linguistic boundaries.

5. Multimodal Speech Recognition: Integration with Vision and Text

Human communication is inherently multimodal, encompassing speech, visual expressions, gestures, and textual cues that together convey meaning. Traditional ASR systems, however, operate solely on acoustic information, ignoring the complementary signals that humans naturally use to disambiguate speech. The concept of multimodal speech recognition (MSR) addresses this limitation by integrating visual and textual modalities to improve robustness, semantic understanding, and generalization. The motivation behind multimodal integration is straightforward: while speech conveys linguistic content, vision provides spatial and contextual grounding, and text encapsulates prior knowledge and semantic structure. This multimodal fusion enables ASR systems to perform effectively in challenging conditions such as noisy environments, overlapping speech, or ambiguous linguistic contexts, offering a path toward human-like comprehension.

Audio-visual speech recognition (AVSR) is one of the most studied forms of multimodal ASR. It exploits visual information, particularly lip movements and facial expressions, to supplement acoustic features. Early AVSR systems utilized handcrafted visual descriptors combined with traditional acoustic models, but these approaches were limited by their

dependence on predefined feature extraction. The advent of deep learning enabled the use of convolutional and transformer-based architectures to jointly learn spatiotemporal representations from both modalities. Models such as LipNet, SyncNet, and AV-HuBERT represent major milestones in this direction. AV-HuBERT, in particular, extends self-supervised pretraining to the multimodal domain by jointly predicting masked visual and audio tokens, thereby learning cross-modal correspondences without explicit supervision. Such pretraining frameworks have demonstrated significant performance improvements in noisy or occluded speech scenarios, confirming that visual information can serve as a powerful complementary modality when acoustic signals are degraded.

Beyond vision, the integration of textual and semantic information has also become central to modern multimodal ASR. Textual grounding allows the system to align acoustic signals with contextual or linguistic knowledge derived from external corpora. Models such as SpeechT5 and Whisper unify speech and text representations through joint pretraining, enabling tasks such as speech-to-text translation, summarization, and contextual recognition. This approach not only improves the syntactic and semantic fluency of transcriptions but also provides resilience against homophones and ambiguous utterances. The integration of textual embeddings from pretrained language models like BERT, T5, or GPT into ASR pipelines has further enhanced contextual understanding. By fusing semantic priors from large text corpora with acoustic features, these models exhibit a form of grounded reasoning-understanding not only "what was said" but also "what was meant." Such multimodal alignment bridges the gap between ASR and natural language understanding, positioning modern systems closer to generalpurpose language intelligence.

Recent developments in cross-modal pretraining have expanded the boundaries of multimodal speech recognition even further. Large-scale foundation models such as CLIP and GPT-4V have demonstrated that joint embeddings across audio, image, and text domains can yield powerful representations capable of zero-shot transfer. Inspired by these frameworks, researchers have begun applying contrastive learning between speech, visual frames, and textual captions to develop unified multimodal encoders. These systems are capable of performing diverse downstream tasks, from transcription and translation to audiovisual understanding and multimodal retrieval. For instance, SpeechCLIP aligns audio with visual and textual contexts, allowing the model to associate spoken words with objects or actions visible in a scene. Similarly, multimodal large language models trained on aligned audiovisual-text datasets have shown the ability to process natural conversations, identify speakers, and interpret emotional or situational cues. marking a significant leap toward human-like perception in ASR.

In multilingual and low-resource settings, multimodal integration offers even greater advantages. Visual cues can provide structural regularities that are independent of language, while textual priors from multilingual corpora can enhance recognition accuracy in underrepresented languages. For example, lip movement patterns are largely language-agnostic,

making AVSR models effective in transferring knowledge across languages without extensive labeled data. Similarly, incorporating multilingual textual embeddings into speech models enables them to perform zero-shot recognition or translation in languages unseen during training. The SeamlessM4T framework exemplifies this synergy by jointly training across speech, text, and translation modalities to achieve high performance in over one hundred languages. By leveraging multimodal signals, such models not only bridge the gap between languages but also between sensory modalities, reinforcing their robustness and adaptability across diverse communication environments.

Despite these remarkable advancements, multimodal ASR still faces challenges in representation alignment, temporal synchronization, and computational efficiency. Fusing heterogeneous modalities requires precise temporal correspondence between audio and visual streams, which can be difficult to maintain under realistic recording conditions. Moreover, multimodal training increases model complexity, demanding higher computational resources and larger datasets for effective convergence. Another open question concerns interpretability-understanding how and when different modalities contribute to recognition accuracy. Research into attention visualization and cross-modal attribution seeks to uncover the internal dynamics of multimodal models, offering insights that may lead to more explainable and efficient architectures. Ethical considerations also arise, particularly concerning the collection and use of audiovisual data that may contain personally identifiable information. Ensuring privacy, fairness, and inclusivity remains a critical requirement for the responsible deployment of multimodal ASR technologies.

In essence, multimodal speech recognition represents a natural progression of ASR research toward comprehensive, context-aware, and human-centric understanding. By integrating visual and textual information, multimodal systems go beyond the purely acoustic domain to capture intent, emotion, and situational context. The convergence of self-supervised learning, transformer architectures, and cross-modal pretraining has brought the field closer to universal models capable of handling diverse environments, languages, and communication forms. As the next sections will discuss, evaluating these models objectively and addressing their open challenges are vital steps toward building globally accessible and ethically aligned speech recognition technologies.

6. Evaluation Metrics and Benchmark Datasets

Evaluating the performance of multilingual and multimodal speech recognition systems requires a comprehensive framework that accounts for both linguistic diversity and multimodal complexity. Unlike traditional monolingual ASR tasks, where evaluation focuses primarily on transcription accuracy, multilingual and multimodal systems must also be assessed for generalization across languages, robustness to environmental noise, and the degree of semantic and contextual understanding achieved through multimodal integration. Therefore, designing fair and representative evaluation methodologies has become as critical as model development itself.

The most widely used evaluation metric in ASR is Word Error Rate (WER), which measures the ratio of insertions, deletions, and substitutions relative to the total number of words in the reference transcription. Although WER remains the de facto standard due to its simplicity and interpretability, it has limitations when applied to multilingual settings, especially for languages that lack explicit word boundaries, such as Mandarin, Japanese, or Thai. To address this, the Character Error Rate (CER) and Subword Error Rate (SER) are often employed for languages with non-segmented writing systems or large vocabularies. These metrics provide finer granularity and avoid penalizing segmentation inconsistencies. However, WER and CER only measure surface-level transcription accuracy and fail to capture semantic correctness or contextual relevance, both of which are increasingly important for multimodal and translation-aware systems.

In multilingual ASR, evaluation also involves measuring cross-lingual generalization and code-switching performance. Cross-lingual generalization tests a model's ability to recognize or transcribe languages unseen during training, reflecting the robustness of learned representations. Codeswitching evaluation, on the other hand, assesses performance on mixed-language utterances that alternate between two or more languages. This phenomenon is common in bilingual societies and poses significant challenges for decoding and language modeling. Metrics like the Mixed Error Rate (MER) have been proposed to jointly evaluate code-switching segments. Moreover, in tasks where ASR outputs are used for downstream applications such as speech translation or captioning, semantic-oriented metrics like BLEU and METEOR are applied to measure the fidelity of meaning preservation rather than strict lexical accuracy. Recent works also incorporate BERTScore and Semantic Error Rate (SemER), leveraging contextual embeddings from pretrained language models to quantify semantic alignment between hypothesis and reference, which is particularly relevant for multimodal and generative ASR systems.

In the multimodal context, evaluation extends beyond textual transcription to include alignment and fusion quality across modalities. For audio-visual ASR, metrics such as Audio-Visual Word (AV-WER) Error Rate Synchronization Error (SyncER) quantify how effectively the system integrates information from both streams. Visual quality and temporal alignment between lip movements and predicted speech are essential for realistic multimodal comprehension. Researchers also employ Signal-to-Noise Ratio (SNR) robustness tests to assess how well multimodal models resist degradation under adverse acoustic conditions. Furthermore. user-centered like metrics perceived intelligibility and naturalness scores have gained attention in evaluating interactive systems where human perception plays a central role, such as dialogue agents and assistive communication devices.

Alongside metrics, the choice of benchmark datasets plays a crucial role in evaluating and comparing ASR systems. Traditional datasets such as LibriSpeech, TED-LIUM, and Switchboard have been instrumental in establishing baselines for English ASR but are limited in linguistic diversity. The growth of multilingual ASR has led to the creation of large-

scale, community-driven corpora covering hundreds of languages. The Mozilla CommonVoice dataset, for example, provides open-source recordings in over 100 languages contributed by volunteers, promoting inclusivity and transparency. Similarly, FLEURS and VoxPopuli offer aligned multilingual speech and translation data, enabling evaluation of cross-lingual and cross-task performance. These datasets are particularly valuable for assessing zero-shot and few-shot generalization, as they contain balanced subsets representing both high- and low-resource languages.

For code-switching and multilingual realism, datasets such as SEAME and Bangla-Hindi CS Corpus provide naturally mixed-language utterances that capture authentic conversational dynamics. In low-resource contexts, regional and community-driven corpora like African Voices, IndicTTS, and Masakhane Speech fill crucial gaps by documenting underrepresented languages. In multimodal ASR, datasets integrating visual information have become equally important. LRS2, LRS3, and AVSpeech serve as standard benchmarks for audio-visual speech recognition, offering synchronized audio and video clips collected from diverse sources such as lectures, news, and interviews. These datasets have facilitated significant progress in lip-reading and audiovisual fusion techniques. Beyond visual modalities, emerging datasets like How2 and SpokenCOCO combine speech, text, and image captions, enabling joint training and evaluation of multimodal models across tasks such as speech understanding, captioning, and retrieval.

Despite these advancements, several challenges persist in constructing comprehensive and fair benchmarks. First, dataset imbalance remains a major issue: high-resource languages dominate existing corpora, skewing model optimization toward well-represented phonetic patterns. Second, variations in recording quality, speaker demographics, environmental conditions hinder cross-dataset comparability. Moreover, multimodal datasets often require precise temporal synchronization, which can be difficult to achieve at scale. Researchers are now exploring synthetic data generation and data augmentation through diffusion or voice cloning methods to mitigate these limitations. Another growing concern is ethical and privacy compliance, particularly for audiovisual datasets containing identifiable speakers. Ensuring informed consent, secure storage, and responsible distribution of data is essential for maintaining public trust in speech research.

In summary, evaluating multilingual and multimodal ASR requires moving beyond conventional accuracy-based metrics to encompass semantic, cross-lingual, and perceptual dimensions. Similarly, dataset design must prioritize diversity, representativeness, and ethical transparency to reflect real-world linguistic and multimodal complexity. Robust evaluation not only benchmarks technological progress but also ensures that advancements in ASR serve a broader and more equitable spectrum of global users. The following section builds on these considerations to discuss open challenges and emerging research directions shaping the next generation of speech recognition technologies.

7. Open Challenges and Future Trends

Despite significant progress in multilingual and multimodal speech recognition, a number of open challenges continue to constrain the scalability, fairness, and interpretability of current systems. These challenges stem not only from the intrinsic complexity of language and multimodal integration but also from broader societal, computational, and ethical considerations that accompany the deployment of large-scale AI models. Addressing these issues requires coordinated advances in model design, data curation, and evaluation methodology, along with a commitment to inclusivity and sustainability in the global AI ecosystem.

One of the most persistent challenges in multilingual ASR lies in linguistic imbalance and representation fairness. Although large-scale multilingual models such as Whisper and SeamlessM4T cover over one hundred languages, their performance remains heavily skewed toward high-resource languages with abundant training data. This imbalance reinforces existing linguistic hierarchies, where widely spoken languages continue to dominate digital ecosystems while minority languages lag behind. Achieving equitable performance across languages requires rethinking data sampling, model objectives, and evaluation criteria. Future research is likely to emphasize language-conditioned adaptation, where the model dynamically allocates capacity based on linguistic diversity and resource availability. Another promising direction involves federated multilingual learning, allowing localized training on community data without central aggregation, thus improving inclusivity while respecting data sovereignty and privacy.

Another major frontier involves robustness and generalization in real-world environments. Current ASR models, though powerful in controlled benchmarks, still struggle under noisy, reverberant, or low-bandwidth conditions. This problem becomes more acute in low-resource regions where recording equipment and network quality are inconsistent. While multimodal integration helps mitigate some of these limitations, it also introduces synchronization and alignment challenges across modalities. Developing noiseinvariant and latency-tolerant architectures that operate efficiently on edge devices remains an active area of research. Approaches combining adaptive front-end processing, knowledge distillation, and continual learning are expected to play a vital role in improving robustness without excessive computational cost. Furthermore, expanding multimodal ASR beyond vision and text to include physiological and environmental sensors could open new frontiers in contextaware understanding.

The issue of scalability and efficiency also poses significant constraints as ASR systems grow larger. Transformer-based architectures with billions of parameters deliver outstanding results but are computationally demanding and environmentally costly. Training such models requires vast GPU clusters and high energy consumption, which raises sustainability concerns. Model compression techniques-such as pruning, quantization, and low-rank factorization-have become increasingly important to enable deployment on mobile or embedded systems. Parameter-efficient fine-tuning

approaches, including LoRA and adapter modules, are already reducing the computational footprint of multilingual adaptation. Looking forward, research into energy-efficient transformer design and green AI optimization will be critical to ensuring that progress in ASR remains both scalable and environmentally responsible.

From a multimodal perspective, cross-modal alignment and interpretability present new scientific challenges. As models integrate audio, vision, and text, understanding how these modalities interact internally becomes essential for both debugging and ethical accountability. The mechanisms that drive transformer architectures are often treated as black boxes, making it difficult to trace how different modalities influence recognition outcomes. Techniques such as attention visualization, modality attribution mapping, and explainable cross-modal reasoning are being developed to improve transparency. Moreover, as multimodal ASR begins to handle complex communicative cues-such as emotion, intent, and gesture-it becomes important to ensure that the model's interpretations remain faithful and unbiased. The development of standardized interpretability benchmarks and metrics will therefore be indispensable for responsible multimodal AI research.

Another pressing issue concerns privacy, ethics, and data governance. Speech data inherently carries sensitive information, including identity, emotion, and social background. The addition of visual and contextual data further increases the risk of misuse or unauthorized identification. Ensuring ethical compliance in multimodal ASR requires anonymization methods, differential mechanisms, and transparent data collection protocols. Federated and decentralized learning frameworks offer potential solutions by allowing on-device training and preventing raw data transmission. At the same time, ethical design must address not only privacy but also bias mitigation and representation equity. Multilingual and multimodal systems risk amplifying cultural stereotypes or excluding minority voices if their training data lacks diversity. Establishing global standards for dataset documentation, informed consent, and demographic balance is therefore vital to ensuring fairness in large-scale speech technologies.

Looking ahead, the next generation of ASR research is expected to move toward universal, adaptive, and multimodal intelligence. Future systems will not merely transcribe speech but will understand its intent, emotion, and situational context across languages and modalities. Advances in large-scale pretraining, cross-lingual transfer, and self-supervised representation learning will continue to narrow the gap between high- and low-resource languages. Meanwhile, multimodal ASR will evolve into an integral component of embodied AI, supporting natural interaction in robotics, virtual assistants, and immersive environments. The convergence of audio, vision, and text under unified transformer backbones promises seamless integration of speech understanding within broader cognitive architectures. Ultimately, achieving this vision will require interdisciplinary collaboration among linguists, engineers, ethicists, and policymakers to ensure that ASR technologies not only advance in capability but also align

with human values of inclusivity, transparency, and sustainability.

In summary, the future of multilingual and multimodal speech recognition will depend on balancing technological innovation with ethical responsibility. The field must strive to create models that are globally representative, computationally efficient, and socially beneficial. Progress in this direction will transform ASR from a domain-specific utility into a universal medium of human–machine communication, bridging linguistic divides and enabling truly global access to intelligent systems.

8. Conclusion

The rapid evolution of speech recognition technology from rule-based systems to large-scale multilingual and multimodal architectures represents one of the most transformative achievements in modern artificial intelligence. This paper has provided a comprehensive overview of recent advances in multilingual and multimodal speech recognition, highlighting how self-supervised learning, transformer architectures, and cross-lingual transfer have fundamentally reshaped the design and capability of ASR systems. Through an examination of multilingual adaptation strategies, low-resource learning frameworks, and multimodal integration with visual and textual information, the survey underscores a consistent trajectory toward inclusivity, scalability, and contextual understanding in speech technologies.

A key insight emerging from this review is that the progress of ASR increasingly depends on shared representation learning-the ability of models to capture universal acoustic, phonetic, and semantic structures that generalize across languages and modalities. Multilingual pretraining frameworks such as XLS-R, Whisper, and SeamlessM4T have demonstrated that a single model can effectively perform recognition, translation, and transcription across hundreds of languages, achieving remarkable zero-shot and few-shot generalization. Similarly, multimodal extensions such as AV-HuBERT and SpeechT5-Multi reveal that incorporating visual and textual cues enables ASR systems to approximate human-like comprehension by leveraging contextual information. These advances collectively mark a paradigm shift from task-specific recognition toward integrated understanding.

Nevertheless, the challenges facing the field remain complex and multifaceted. Persistent disparities in data availability across languages continue to undermine fairness and inclusivity. While large-scale pretraining has mitigated some of these gaps, high-resource languages still dominate the learning process, limiting the accessibility of ASR technology for minority linguistic communities. In addition, multimodal integration introduces new technical hurdles synchronization, alignment, interpretability. and Understanding how audio, visual, and textual signals interact within a single model remains an open research problem. Moreover, the growing scale of models raises concerns about computational sustainability, privacy protection, and ethical accountability. Ensuring that the benefits of ASR extend equitably across global populations will require concerted efforts not only from researchers but also from policymakers and civil society.

Looking forward, the future of speech recognition lies in universal, adaptive, and responsible intelligence. The combination of multilingual and multimodal learning will enable ASR systems that are not only linguistically comprehensive but also contextually aware and ethically grounded. Advances in energy-efficient architectures, federated learning, and explainable AI will allow models to operate sustainably and transparently at scale. Moreover, as ASR becomes increasingly integrated into interactive and embodied AI systems-such as personal assistants, robots, and extended reality platforms-the ability to process speech within its full multimodal and cultural context will become essential. This convergence will transform speech recognition from a passive transcription tool into an active medium of communication and understanding between humans and machines.

Ultimately, the evolution of multilingual and multimodal ASR embodies the broader vision of artificial intelligence as a tool for inclusivity, accessibility, and global collaboration. The capacity to understand speech across linguistic, cultural, and sensory boundaries holds profound implications for education, healthcare, governance, and social integration. Future research must therefore balance technological ambition with social responsibility, ensuring that advances in speech recognition reflect not only the sophistication of algorithms but also the diversity and dignity of the voices they seek to understand. By bridging languages, modalities, and communities, nextgeneration ASR systems have the potential to redefine how humanity interacts with intelligent technology-making communication not merely faster, but more universal, equitable, and meaningful.

References

- [1] A. Baevski, H. Zhou, A. Mohamed, and M. Auli, "wav2vec 2.0: A framework for self-supervised learning of speech representations," Advances in Neural Information Processing Systems (NeurIPS), vol. 33, pp. 12449–12460, 2020.
- [2] W.-N. Hsu, B. Bolte, Y.-H. H. Tsai, K. Lakhotia, A. Mohamed, and A. Baevski, "HuBERT: Self-supervised speech representation learning by masked prediction of hidden units," IEEE/ACM Transactions on Audio, Speech, and Language Processing, vol. 30, pp. 1682–1695, 2022.
- [3] A. Babu, W.-N. Hsu, B. Bolte, J. Lee, and A. Mohamed, "XLS-R: Self-supervised cross-lingual speech representation learning at scale," arXiv preprint arXiv:2111.09296, 2022.
- [4] A. Radford, J. Kim, T. Xu, et al., "Whisper: Robust speech recognition via large-scale weak supervision," arXiv preprint arXiv:2212.04356, 2023
- [5] C. Tang, L. He, and J. Tao, "SpeechT5: Unified-modal encoder-decoder pretraining for spoken language processing," IEEE/ACM Transactions on Audio, Speech, and Language Processing, vol. 31, pp. 252–264, 2023.
- [6] X. Shi, J. Wu, and S. Watanabe, "AV-HuBERT: Self-supervised audiovisual speech representation learning," IEEE Transactions on

- Pattern Analysis and Machine Intelligence, vol. 46, no. 4, pp. 2358–2372, 2024.
- [7] D. Amodei, S. Ananthanarayanan, R. Prabhavalkar, et al., "Deep Speech 2: End-to-end speech recognition in English and Mandarin," Proceedings of the 33rd International Conference on Machine Learning (ICML), pp. 173–182, 2016.
- [8] J. Dong, S. Xu, and B. Xu, "Speech-Transformer: A no-recurrence sequence-to-sequence model for speech recognition," IEEE International Conference on Acoustics, Speech and Signal Processing (ICASSP), pp. 5884–5888, 2018.
- [9] A. Gulati, J. Qin, C.-C. Chiu, et al., "Conformer: Convolutionaugmented transformer for speech recognition," Interspeech 2020, pp. 5036–5040, 2020.
- [10] N. Tits, T. Dutoit, and A. Haddad, "Exploring the potential of multilingual self-supervised pretraining for low-resource ASR," Computer Speech & Language, vol. 86, 101501, 2024.
- [11] A. Zhang, R. Zhao, and Y. Liu, "Multilingual speech recognition using a shared encoder and language adapters," IEEE Transactions on Artificial Intelligence, vol. 4, no. 1, pp. 37–48, 2023.
- [12] Meta AI Research, "SeamlessM4T: Massively multilingual & multimodal machine translation," arXiv preprint arXiv:2308.11579, 2023
- [13] S. Watanabe, T. Ochiai, and J. Karita, "ESPnet-ST: All-in-one speech translation toolkit," IEEE/ACM Transactions on Audio, Speech, and Language Processing, vol. 29, pp. 1293–1304, 2021.
- [14] M. Riviere, A. Joulin, P. Kharitonov, and E. Dupoux, "Unsupervised pretraining transfers well across languages," IEEE International Conference on Acoustics, Speech and Signal Processing (ICASSP), pp. 7414–7418, 2020.
- [15] S. K. Ramesh, A. Kannan, and S. Khudanpur, "Code-switching in multilingual speech recognition: Challenges and strategies," Interspeech 2023, pp. 2234–2238, 2023.
- [16] R. Serizel, J. Barker, and E. Vincent, "Multimodal learning for robust speech recognition in noisy environments," IEEE Signal Processing Magazine, vol. 41, no. 2, pp. 98–112, 2024.
- [17] J. Ma, Y. Li, and Z. Meng, "Audio-visual speech recognition with transformer-based fusion," Neural Networks, vol. 171, pp. 398-411, 2025.
- [18] OpenAI, "GPT-4 Technical Report," arXiv preprint arXiv:2303.08774, 2023.
- [19] C. Chang, S. Kanda, and J. Villalba, "Cross-lingual adaptation and dataefficient fine-tuning for multilingual ASR," IEEE Automatic Speech Recognition and Understanding Workshop (ASRU), pp. 341–348, 2023.
- [20] T. Ko, V. Peddinti, D. Povey, and S. Khudanpur, "Audio augmentation for speech recognition," Interspeech 2015, pp. 3586–3589, 2015.
- [21] P. Kharitonov, A. Likhomanenko, and A. Mohamed, "Textless speech-to-speech translation on real data," Transactions of the Association for Computational Linguistics (TACL), vol. 12, pp. 142–159, 2024.
- [22] M. Pratap, K. Bartelds, and S. Mallidi, "FLEURS: Benchmarking speech recognition and translation in 102 languages," Proceedings of the 2023 Conference on Empirical Methods in Natural Language Processing (EMNLP), pp. 7782–7791, 2023.
- [23] Mozilla Foundation, "CommonVoice Dataset: Democratizing voice technology," Dataset Release v17.0, 2024.
- [24] S. Watanabe, T. Hayashi, and P. Garcia, "ESPnet2: End-to-end speech processing toolkit," IEEE/ACM Transactions on Audio, Speech, and Language Processing, vol. 30, pp. 2960–2974, 2022.
- [25] J. Wang, M. Zhao, and T. Li, "Green speech AI: Sustainable training and inference for large-scale ASR systems," IEEE Access, vol. 13, pp. 102911–102923, 2025.