

Summary of Automatic Audio Recognition Papers

Abdulkader Sardini

Sohaib Belaroussi

Henry Odongo

Syuja Akmal

April 22, 2025

1 Kader's part

1.1 Arabic Paper

try

Key findings: The study reveals that...

2 A Historical Perspective of Speech Recognition

- This 2014 paper, authored by Xuedong Huang, James Baker, and Raj Reddy, provides a historical overview of automatic speech recognition (ASR) research, tracing its evolution from the early days of limited capabilities in 1976 to the current era of sophisticated voice assistants like Siri and Google Assistant. The authors highlight the key breakthroughs that have driven ASR's progress, including the development of hidden Markov models (HMMs), statistical modeling techniques, and the advent of deep neural networks (DNNs).

- The paper emphasizes the importance of large datasets and computing power in advancing ASR, noting that Moore's law has played a crucial role in enabling the development of increasingly complex and accurate systems. However, the authors also acknowledge the limitations of current ASR systems, particularly in handling noisy or accented speech, and highlight the need for further research in areas like data efficiency, robustness, and generalization.

- The authors identify six main challenges that must be addressed to move ASR to the next level, including the need for more data, improved computing infrastructure, better handling of uncertainties, and more robust speaker-independent and adaptive systems. They also discuss the importance of incorporating prosody (intonation, rhythm, and stress) into ASR models, as this crucial aspect of human speech has been largely ignored in the past. –

3 Trends and Developments in Automatic Speech Recognition Research

- This 2023 paper, authored by Douglas O'Shaughnessy, delves into the intricacies of automatic speech recognition (ASR) research, focusing on the unique challenges posed by the complex nature of human speech. The author contrasts the traditional approach of using hidden Markov models (HMMs) with the more recent trend of employing deep neural networks (DNNs), highlighting the advantages and limitations of each approach. The paper emphasizes the importance of understanding the acoustic-phonetic properties of speech, as well as the limitations of current ASR systems in handling noisy or accented speech.

- O'Shaughnessy explores various aspects of speech analysis, including spectral analysis, Mel-frequency cepstral coefficients (MFCCs), and formant tracking, and discusses the trade-offs between accuracy and computational efficiency. He also examines different types of supervised and unsupervised learning methods used in ASR, along with the challenges of data scarcity and speaker variability.

- The paper concludes by suggesting potential avenues for future research in ASR, including the development of more robust and efficient systems that can better handle noisy and accented speech, incorporate prosody, and exploit the unique characteristics of human speech. The author emphasizes the need for a deeper understanding of the underlying principles of speech production and perception to guide the development of more effective ASR systems.

4 Suhaib's part [replace text with content]

4.1 Test

Speech Recognition by Machine: A Review

This article offers a comprehensive survey of the technological developments and foundational concepts in Automatic Speech Recognition (ASR) spanning six decades of research. It explores three major approaches in the field: the Acoustic-Phonetic approach, which decodes speech based on phonetic units; the Pattern Recognition approach, which applies statistical models such as Hidden Markov Models (HMMs) and techniques like Dynamic Time Warping (DTW); and the Artificial Intelligence approach, which involves expert systems and neural networks for modeling and adapting to speech patterns.

The article categorizes ASR systems by the types of speech they process—such as isolated words, connected speech, continuous, and spontaneous speech—and outlines practical applications in telecommunications, education, healthcare, and military domains. It highlights performance factors like vocabulary size, noise handling, speaker variability, and system adaptability. Essential techniques such as Mel Frequency Cepstral Coefficients (MFCC), Linear Predictive Coding (LPC), and Support Vector Machines (SVMs) are also emphasized for their importance in feature extraction and classification.

The paper also includes a historical overview, tracing progress from early analog devices like “Radio Rex” in the 1920s to modern pattern-based and neural systems. It emphasizes how innovations in hardware and algorithms have driven advances in ASR. The review concludes by identifying current challenges in spontaneous speech recognition, robust performance in diverse environments, and the need for system personalization.

The History of Speech Recognition to the Year 2030

This article by Awni Hannun reflects on the evolution of ASR technology from 2010 to 2020 and anticipates future developments through 2030. Major breakthroughs during the last decade—including deep learning, large-scale annotated datasets, and GPU acceleration—have enabled dramatic reductions in Word Error Rates (WER), surpassing human transcription performance in benchmark tasks. Innovations like Kaldi, LibriSpeech, Deep Speech models, and streaming on-device systems have redefined the ASR landscape.

Looking forward, Hannun predicts a shift in research focus from reducing WER to enhancing system usability and integration with downstream applications. He emphasizes the growing importance of self- and semi-supervised learning, lightweight model design, and on-device inference. These approaches offer benefits such as improved privacy, lower latency, and consistent performance in offline settings.

The article concludes by highlighting the need for personalized ASR systems that adapt to individual users’ accents, speech patterns, and environments. Hannun also warns that increasing centralization of ASR research in large technology companies may hinder academic progress. Nevertheless, he remains optimistic about ASR’s future in enabling accessible, intelligent, and context-aware speech technologies across various industries.