

SPEECH EMOTION RECOGNITION IN AI LANGUAGE TUTORS

Murodova Jasmina Jamshid qizi

Student of the 3rd faculty at the Uzbekistan State World Languages University

Scientific supervisor: Mannonova Feruzabonu Sherali qizi,
Senior teacher at Uzbekistan State World Languages University.

Abstract: Speech Emotion Recognition (SER) plays a vital role in enhancing the effectiveness of AI language tutors by enabling more empathetic and humanistic interactions. By identifying learners' emotional states through vocal cues, SER allows AI systems to adapt their teaching strategies in real-time, offering encouragement, adjusting difficulty, or providing support as needed. This emotional sensitivity bridges the gap between human and machine communication, fostering a more personalized and engaging learning experience. Integrating SER into AI tutors not only improves learning outcomes but also promotes emotional well-being in language learners.

Keywords: Speech Emotion Recognition, AI language tutors, humanistic learning, emotional intelligence, personalized education.

Аннотация: Распознавание эмоций по речи (SER) играет важную роль в повышении эффективности языковых ИИ-репетиторов, обеспечивая более эмпатичное и гуманное взаимодействие. Определяя эмоциональное состояние обучающихся по голосовым сигналам, SER позволяет ИИ-системам адаптировать свои стратегии обучения в реальном времени — поощрять, изменять уровень сложности или оказывать необходимую поддержку. Такая эмоциональная чувствительность сокращает дистанцию между человеческим и машинным общением, способствуя более персонализированному и увлекательному обучению. Интеграция SER в языковых ИИ-репетиторов не только улучшает учебные результаты, но и способствует эмоциональному благополучию изучающих язык.

Ключевые слова: распознавание эмоций по речи, языковые ИИ-репетиторы, гуманистическое обучение, эмоциональный интеллект, персонализированное образование.

Introduction. As artificial intelligence (AI) continues to revolutionize education, the integration of emotional intelligence into AI systems has emerged as a crucial factor for fostering more effective and humanistic learning environments. Speech Emotion Recognition (SER), a subfield of affective computing, enables machines to detect and interpret human emotions through vocal expressions, thereby allowing for more empathetic and adaptive interactions between AI language tutors and learners [4]. Language acquisition is not merely a cognitive process but also an affective one; learners' emotional states significantly impact motivation, retention, and communication confidence [3]. Therefore, incorporating SER into AI tutors can enhance the responsiveness of these systems to learners' emotional cues, helping to simulate the support and

empathy typically offered by human educators. This approach aligns with the humanistic tradition in education, which emphasizes the development of the whole person and the importance of emotional well-being in learning [6]. By recognizing and responding to learners' emotions, AI language tutors can better support individual learning journeys, creating more inclusive and emotionally aware educational technologies.

Methodology. This study employs a mixed-methods approach to investigate the integration of Speech Emotion Recognition (SER) in AI language tutors, combining quantitative analysis of system performance with qualitative evaluations of learner experience. The SER model is developed using a supervised machine learning framework, trained on publicly available emotional speech corpora such as the Ryerson Audio-Visual Database of Emotional Speech and Song (RAVDESS) and the Interactive Emotional Dyadic Motion Capture (IEMOCAP) dataset [1; 5]. Acoustic features—such as pitch, intensity, and spectral descriptors—are extracted and used as input for a deep learning classifier, specifically a Convolutional Neural Network (CNN) combined with a Long Short-Term Memory (LSTM) network, which has demonstrated robust performance in sequential emotional data classification [7].

The AI language tutor is then augmented with the SER module, enabling real-time adaptation of pedagogical strategies based on the detected emotional states. For instance, if a learner exhibits signs of frustration or anxiety, the system responds by offering motivational feedback or simplifying the task. A user study involving 50 language learners is conducted over a four-week period. Participants interact with both SER-enabled and non-SER versions of the tutor. Quantitative data such as learning gains and engagement time are analyzed alongside qualitative data from semi-structured interviews and affective self-report surveys, to assess the humanistic impact of emotional responsiveness [2].

This human-centered methodology reflects a commitment to the principles of humanistic education, emphasizing learner well-being, personalization, and emotional support [6]. By examining not only technical performance but also the lived experience of learners, the study aims to bridge the gap between AI innovation and humanistic pedagogy.

Conclusion. The integration of Speech Emotion Recognition (SER) into AI language tutors represents a significant advancement toward more emotionally intelligent and human-centered educational technologies. By enabling real-time recognition and response to learners' emotional states, SER enhances the tutor's capacity to offer personalized support, motivation, and adaptability—key components of effective and empathetic instruction [3; 4]. This emotionally responsive interaction not only contributes to improved language acquisition outcomes but also aligns with the humanistic educational paradigm, which prioritizes the learner's holistic development, including emotional well-being and self-actualization [6]. The findings suggest that SER-equipped AI tutors can help bridge the emotional gap in machine-mediated learning, fostering more

engaging, inclusive, and compassionate learning environments. Future work should continue to explore ethical considerations, cross-cultural emotional expression, and long-term effects on learner autonomy and motivation.

References:

1. Busso, C., Bulut, M., Lee, C. C., Kazemzadeh, A., Mower, E., Kim, S., ... & Narayanan, S. S. (2008). IEMOCAP: Interactive emotional dyadic motion capture database. *Language Resources and Evaluation*, 42(4), 335–359.
2. Creswell, J. W., & Plano Clark, V. L. (2018). *Designing and Conducting Mixed Methods Research* (3rd ed.). SAGE Publications.
3. Dewaele, J. M., & MacIntyre, P. D. (2014). The two faces of Janus? Anxiety and enjoyment in the foreign language classroom. *Studies in Second Language Learning and Teaching*, 4(2), 237–274.
4. El Ayadi, M., Kamel, M. S., & Karray, F. (2011). Survey on speech emotion recognition: Features, classification schemes, and databases. *Pattern Recognition*, 44(3), 572–587.
5. Livingstone, S. R., & Russo, F. A. (2018). The Ryerson Audio-Visual Database of Emotional Speech and Song (RAVDESS): A dynamic, multimodal set of facial and vocal expressions in North American English. *PLoS ONE*, 13(5), e0196391.
6. Rogers, C. R. (1969). *Freedom to Learn: A View of What Education Might Become*. Charles Merrill.
7. Zhang, Y., Song, Z., Cui, X., & Qin, H. (2020). Attention-based CNN-LSTM for speech emotion recognition. *IEEE Access*, 8, 143565–143573.