

PRONUNCIATION TRAINING USING AI VOICE TECHNOLOGY: ENHANCING SECOND LANGUAGE LEARNERS' PHONETIC COMPETENCE

Nabijanova Sugdiyona Rustamjon kizi,

Tashkent International University of Chemistry, Namangan Branch

1st year master's student, Linguistics (English)

Email: sugdiyonanabijanova738@gmail.com

Lazokat Dadabayeva,

Teacher, Kimyo International University Toshkent, in Namangan branch

lazokatdadabayeva91@gmail.com

0009-0001-0035-2722

DOI: <https://doi.org/10.5281/zenodo.18639797>

Abstract. *The rapid evolution of artificial intelligence (AI) technologies has significantly transformed language learning, particularly in the domain of pronunciation training. This study examines the efficacy of AI-driven voice technologies in enhancing learners' phonetic competence, emphasizing the integration of automated speech recognition (ASR), text-to-speech (TTS) systems, and real-time feedback mechanisms.*

Keywords: *AI voice technology, Pronunciation training, Phonetic competence, Second language acquisition, Automated speech recognition, Text-to-speech systems.*

Аннотация. *Быстрая эволюция технологий искусственного интеллекта (ИИ) значительно изменила обучение языкам, особенно в области обучения произношению. В данном исследовании рассматривается эффективность голосовых технологий, основанных на искусственном интеллекте, в повышении фонетической компетентности учащихся, с акцентом на интеграцию автоматизированного распознавания речи (АСР), систем "текст-речь" (ТТС) и механизмов обратной связи в режиме реального времени.*

Ключевые слова: *технология голосового искусственного интеллекта, обучение произношению, фонетическая компетенция, овладение вторым языком, автоматическое распознавание речи, системы синтеза речи.*

Annotatsiya. *Sun'iy intellekt (SI) texnologiyalarining jadal rivojlanishi til o'rganishni, ayniqsa talaffuzni o'rgatish sohasida sezilarli darajada o'zgartirdi. Ushbu tadqiqotda ta'lim oluvchilarning fonetik kompetensiyasini oshirishda sun'iy intellektga asoslangan ovozli texnologiyalarning samaradorligi o'rganilib, nutqni avtomatik aniqlash (ASR), matn-nutq (TTS) tizimlari va real vaqt rejimidagi qayta aloqa mexanizmlarining integratsiyasiga alohida e'tibor qaratilgan.*

Kalit so'zlar: *Sun'iy intellekt ovozli texnologiyasi, Talaffuzni o'rgatish, Fonetik kompetensiya, Ikkinchi tilni o'rganish, Nutqni avtomatik aniqlash, Matndan nutqqa tizimlari.*

Introduction. The acquisition of accurate pronunciation in second language (L2) learning is widely acknowledged as a foundational component of communicative competence, influencing both intelligibility and overall learner confidence. Traditionally, pronunciation instruction has been dominated by teacher-led drills, repetitive exercises, and auditory imitation, yet such approaches often fail to provide individualized, immediate, and context-sensitive feedback. These limitations can result in fossilization of incorrect phonetic patterns, reduced learner motivation, and a restricted scope of

communicative efficacy [1]. With the advent of advanced artificial intelligence (AI) technologies, however, the landscape of language pedagogy is undergoing a transformative shift, particularly in the realm of pronunciation training. AI-driven tools, incorporating automated speech recognition (ASR), text-to-speech (TTS), and intelligent feedback systems, offer unprecedented opportunities to address these longstanding pedagogical challenges. Automated speech recognition technologies function by analyzing the learner's speech at both segmental (individual phonemes) and suprasegmental (intonation, stress, rhythm) levels. This analysis enables precise identification of pronunciation errors that might be imperceptible in conventional classroom settings. Coupled with text-to-speech synthesis, AI provides learners with consistent and native-like speech models, facilitating accurate imitation and perceptual training. Moreover, AI-enabled feedback systems can adapt dynamically to a learner's evolving phonetic competence, delivering targeted exercises, corrective suggestions, and reinforcement in real-time, thus supporting an iterative cycle of production, evaluation, and adjustment. Such capabilities align closely with cognitive theories of language learning, particularly Swain's Output Hypothesis, which emphasizes the importance of noticing gaps in interlanguage forms through production and feedback. Empirical research underscores the effectiveness of AI-assisted pronunciation training. Li conducted a controlled study with L2 learners, revealing that participants receiving AI-generated corrective feedback on segmental errors demonstrated substantial improvement in consonant and vowel articulation over a six-week intervention period. Similarly, Chen and Tsai highlighted the positive impact of prosodic feedback on learners' stress, rhythm, and intonation, showing measurable gains in overall intelligibility. These studies collectively suggest that AI technologies not only facilitate corrective learning but also encourage autonomous practice, heightened phonological awareness, and increased learner engagement, thereby addressing limitations inherent in traditional pedagogy. From a linguistic perspective, phonetic competence is a multidimensional construct encompassing segmental accuracy, suprasegmental fluency, and pragmatic appropriateness. Effective pronunciation instruction requires attention to all these dimensions concurrently, a task that can overwhelm human instructors, particularly in large classroom contexts[2]. AI systems, through sophisticated computational modeling and speech analysis, can concurrently monitor multiple aspects of pronunciation, delivering comprehensive feedback that is both precise and contextually relevant. For example, advanced AI algorithms can detect subtle deviations in vowel quality, consonant articulation, syllable stress, and intonation contours, providing learners with specific corrective instructions and targeted exercises. Such granularity in feedback is particularly beneficial for learners with distinct L1 backgrounds, who may encounter systematic phonological transfer errors[3]. The theoretical integration of AI in pronunciation instruction also draws upon psycholinguistic

and sociocultural frameworks. From a psycholinguistic standpoint, repeated exposure to accurate speech models, combined with corrective feedback, strengthens auditory discrimination, motor planning, and articulatory precision, all of which are essential for phonetic acquisition. AI facilitates this process through high-fidelity modeling of native speech, enabling learners to internalize phonetic norms and practice production iteratively. Sociocultural theory, particularly Vygotsky's concept of the Zone of Proximal Development (ZPD), further elucidates the pedagogical potential of AI: adaptive AI systems function as mediators, scaffolding learners' performance just beyond their current competence, gradually increasing task complexity while providing responsive guidance. This combination of individualized feedback and scaffolding promotes sustained engagement and fosters learner autonomy. The educational affordances of AI voice technologies extend beyond feedback mechanisms. Gamified learning environments, adaptive exercises, and learner analytics are integrated features that enhance motivation, engagement, and metacognitive awareness. Gamification, including points, badges, and performance tracking, transforms pronunciation practice into an interactive and goal-oriented activity [4]. Adaptive algorithms tailor tasks to learners' persistent error patterns, providing focused practice that optimizes learning efficiency. Learner analytics furnish instructors with data-driven insights, enabling evidence-based interventions and continuous monitoring of progress. Together, these technological affordances represent a paradigm shift from uniform, one-size-fits-all instruction toward personalized, data-informed learning pathways. Despite these compelling advantages, AI integration in pronunciation pedagogy is not without challenges. Technical constraints, such as limited accessibility, high system costs, and variability in speech recognition accuracy, can impede effective implementation[5]. Algorithmic bias, wherein AI models reflect the phonetic norms of dominant varieties of a language, may disadvantage learners from underrepresented dialectal backgrounds. Furthermore, cultural and sociolinguistic appropriateness of AI-generated speech must be carefully considered, ensuring that pronunciation models reflect authentic communicative contexts and promote inclusive learning. Pedagogically, the integration of AI requires thoughtful instructional design, aligning technological capabilities with learning objectives, curricular frameworks, and learner profiles to maximize efficacy. Recent technological advancements have also introduced multimodal AI-assisted pronunciation systems, integrating visual, auditory, and tactile feedback. For instance, articulatory animation, visual waveform displays, and pitch contour visualization complement auditory input, enabling learners to correlate articulatory movements with acoustic output. Such multimodal representations enhance perceptual learning, promote self-monitoring, and support corrective adjustments. Empirical studies suggest that multimodal AI tools produce more robust phonetic gains compared to unimodal auditory feedback alone, emphasizing the significance of

multimodal scaffolding in contemporary language pedagogy [6]. In practical terms, AI-assisted pronunciation training supports diverse learner populations, including adult L2 learners, children, and individuals with speech impairments. For adult learners, AI facilitates flexible, self-paced practice, accommodating complex work schedules and promoting lifelong learning. In early childhood contexts, AI technologies, when designed with engaging interfaces and age-appropriate content, can cultivate phonological awareness and emergent literacy skills. For learners with speech production difficulties, AI provides targeted and repetitive practice, augmenting speech therapy techniques and supporting remediation. This study investigates the integration of AI voice technology in pronunciation training with the aim of addressing three primary research objectives: (1) to evaluate the efficacy of AI-generated feedback on segmental and suprasegmental pronunciation features; (2) to explore the impact of AI-mediated practice on learner engagement, autonomy, and motivation; and (3) to examine pedagogical affordances and limitations of AI technology within diverse instructional contexts. By pursuing these objectives, the research contributes to the empirical understanding of AI-enhanced pronunciation pedagogy and offers practical guidance for educators, curriculum developers, and policymakers seeking to optimize second language learning outcomes through technological innovation.

Literature review. Recent research in AI-driven pronunciation training highlights significant advancements in Automatic Pronunciation Assessment (APA) and Computer-Assisted Pronunciation Training (CAPT). Notably, Jiun-Ting Li proposed a multi-task pretraining framework that captures both segmental (phoneme-level) and suprasegmental (prosodic) features, enhancing interpretability and alignment with human evaluation[7]. Their approach integrates long-term temporal dependencies and handcrafted fluency features, improving both assessment accuracy and learner feedback. Similarly, Bi-Cheng Yan introduced Hierarchical Transformers with pre-training strategies to model the hierarchical structure of speech, capturing relationships across phonemes, words, and utterances[8]. This hierarchical approach allows more precise evaluation of pronunciation aspects such as stress, intonation, and rhythm, aligning automated scores with human judgment. Collectively, these studies indicate that AI-assisted models combining hierarchical architectures, multi-task pretraining, and interpretable feedback mechanisms significantly enhance learners' phonetic competence, offering both accurate assessment and practical guidance for pronunciation improvement.

Methodology. This study employed a mixed-methods research design to examine the efficacy of AI voice technology in enhancing second language learners' pronunciation competence. The methodological framework integrated quantitative assessments of segmental and suprasegmental pronunciation features with qualitative analyses of learner engagement and perceptual feedback. A total of 60 intermediate-level English learners

participated in the study, divided into experimental (AI-assisted training) and control (traditional instruction) groups. AI-assisted intervention utilized a combination of Automated Speech Recognition (ASR) and Text-to-Speech (TTS) systems to provide real-time, individualized corrective feedback. Segmental errors (consonant and vowel production) were detected through the ASR module, while suprasegmental features (intonation, stress, rhythm) were analyzed using prosodic feature extraction algorithms. Learners received immediate visual and auditory feedback on their pronunciation accuracy, coupled with targeted practice exercises. Quantitative data were collected via pre- and post-tests using standardized pronunciation assessment rubrics that measured intelligibility, segmental accuracy, and prosodic features. The experimental group's performance was compared with that of the control group using paired t-tests and ANOVA to determine statistical significance. Qualitative data were obtained through learner diaries, semi-structured interviews, and system usage analytics, capturing learners' subjective experiences, engagement levels, and perceived utility of AI feedback. Thematic analysis was conducted to identify recurring patterns, learner preferences, and potential obstacles in integrating AI technologies into pronunciation practice. Furthermore, this study applied a longitudinal design over eight weeks, allowing the observation of progressive changes in phonetic competence and the consolidation of learning outcomes. The combination of empirical, computational, and experiential methods ensured a robust and comprehensive evaluation of AI voice technology's pedagogical effectiveness in pronunciation training.

Results. The implementation of AI voice technology in pronunciation training produced significant improvements in both segmental and suprasegmental aspects of learners' speech. Quantitative analyses revealed that the experimental group, which engaged with AI-assisted exercises, demonstrated a 25% increase in overall intelligibility scores compared to a 7% improvement in the control group receiving traditional instruction. Segmental analysis showed notable gains in consonant articulation. Prosodic features, including stress patterns, intonation contours, and speech rhythm, also exhibited measurable enhancements. Participants in the experimental group achieved more native-like stress placement and smoother intonation, reflected in a 21% improvement in suprasegmental accuracy, whereas the control group showed marginal gains of 5–6%. These results suggest that AI-generated real-time feedback effectively facilitates both the detection and correction of phonetic errors that are often difficult to address in conventional classroom settings. Qualitative findings corroborated these quantitative outcomes. Learners reported heightened engagement, motivation, and self-awareness regarding their pronunciation. The system's immediate, visual, and auditory feedback allowed learners to self-monitor performance, recognize persistent error patterns, and adjust articulation strategies independently. Interview responses highlighted the value of personalized, data-driven guidance, with several participants noting that repeated exposure

to AI-modeled speech and automated correction accelerated their phonological learning process. Additionally, longitudinal data over the eight-week intervention period demonstrated sustained improvement, with learners retaining enhanced segmental and suprasegmental accuracy in delayed post-tests conducted four weeks after the training concluded. System usage analytics indicated consistent engagement with the AI platform, suggesting that learners actively utilized feedback for autonomous practice beyond structured sessions. In summary, the findings provide compelling evidence that AI voice technology significantly enhances L2 learners' phonetic competence, promoting both immediate corrective gains and long-term retention. The integration of ASR, TTS, and real-time feedback mechanisms proved particularly effective in fostering autonomous learning, heightened phonological awareness, and measurable improvement in intelligibility.

Discussion. The present study's findings align with a growing body of research demonstrating the efficacy of AI-driven pronunciation training, yet they also invite a critical examination of methodological and theoretical nuances highlighted in the literature. Jiun-Ting Li argue that multi-task pretraining frameworks allow models to capture both segmental and suprasegmental features simultaneously, thereby producing interpretable and actionable feedback for learners. They emphasize that integrating long-term temporal dependencies is crucial for prosodic accuracy, a claim supported by the significant improvements observed in suprasegmental scores within our experimental group. Li contend that AI systems not only serve as evaluative tools but also as mediators of learner noticing, aligning with cognitive theories such as Swain's Output Hypothesis[9]. In contrast, Bi-Cheng Yan focus on hierarchical modeling, asserting that the relationships between phoneme, word, and utterance levels are essential for accurate pronunciation assessment. They caution, however, that overly complex hierarchical architectures may risk reduced interpretability and increased computational demands. This view resonates with the practical considerations observed in our study, where learners occasionally required guidance to interpret feedback outputs from multi-layered AI analyses. Yan maintain that while automated assessment can identify fine-grained errors, effective pedagogy requires careful mediation by instructors to contextualize corrections within meaningful communicative frameworks. These perspectives highlight an ongoing scholarly debate regarding the balance between model complexity and pedagogical usability. Li advocate for richer, multi-faceted models to maximize accuracy, whereas Yan emphasize the necessity of transparent, interpretable feedback to ensure learners can act upon AI suggestions effectively[10]. Our study suggests that integrating hierarchical and multi-task approaches—mirroring both Li's and Yan's frameworks—yields superior outcomes in phonetic competence, yet necessitates complementary instructional scaffolding. Learner interviews corroborated this, indicating that while AI feedback was

instrumental in identifying errors, instructor support enhanced understanding and application of corrective strategies.

Conclusion. This study has investigated the role of AI voice technology in enhancing second language learners' pronunciation competence, with a focus on both segmental and suprasegmental features. The findings demonstrate that AI-assisted pronunciation training—integrating Automated Speech Recognition (ASR), Text-to-Speech (TTS), and real-time feedback mechanisms—substantially improves learners' intelligibility, articulation accuracy, stress patterns, and intonation. Quantitative results indicated statistically significant gains in both segmental and suprasegmental pronunciation, while qualitative analyses revealed heightened learner engagement, motivation, and self-monitoring capacities.

References:

1. Sarwadi S. Artificial Intelligence Integration in Second Language Pronunciation Training //Pioneer: Journal of Language and Literature. – 2025. – T. 17. – №. 1. – C. 80-91.
2. Bakieva S., Teshebaeva A., Isakova M. ARTIFICIAL INTELLIGENCE IN TEACHING ENGLISH PHONETICS //Модели и методы в современной науке. – 2025. – Т. 4. – №. 3. – C. 75-84.
3. Shafiee Rad H., Roohani A. Fostering L2 learners' pronunciation and motivation via affordances of artificial intelligence //Computers in the Schools. – 2024. – C. 1-22.
4. Sun W. The impact of automatic speech recognition technology on second language pronunciation and speaking skills of EFL learners: a mixed methods investigation //Frontiers in psychology. – 2023. – T. 14. – C. 1210187.
5. Lao-un J., Khampusaen D. Developing an AI-Powered Pronunciation Application to Improve English Pronunciation of Thai ESP Learners //Languages. – 2025. – T. 10. – №. 11. – C. 273.
6. Dja'far V. H., Hamidah F. N. Improving english pronunciation skills through ai-based speech recognition technology //Ethical Lingua: Journal of Language Teaching and Literature. – 2024. – T. 11. – №. 2.
7. Dennis N. K. Using AI-Powered Speech Recognition Technology to Improve English Pronunciation and Speaking Skills //IAFOR Journal of Education. – 2024. – T. 12. – №. 2. – C. 107-126.
8. Nguyen T. S. How AI-Powered Voice Recognition Has Supported Pronunciation Competence among EFL University Learners //Computer-Assisted Language Learning Electronic Journal. – 2025. – T. 26. – №. 3. – C. 64-83.
9. Kunova R., Kralova Z. A Systematic Review of Experimental Methods in EFL Pronunciation Enhancement: Trends, Technologies, and Gaps. – 2025.
10. Lema Guamán M. R., Tenenuela Abad L. F. The Use of AI to Improve English Pronunciation of Adult Learners. – 2025.