

# AssistVoice: A Voice-Based Visual Routine Learning System for Children

Er. Prince Abraham<sup>1</sup>, Neeraja S<sup>2</sup>, Savitha Sabu<sup>3</sup>, Tom Joseph Sajees<sup>4</sup>

<sup>1</sup>Associate Professor, Computer Science Department,  
Saintgits College of Engineering, Kottayam, Kerala, India

<sup>2,3,4</sup>UG Students, Computer Science Department,  
Saintgits College of Engineering, Kottayam, Kerala, India

Email: <sup>1</sup>prince.a@saintgits.org <sup>2</sup>neerajasreelal@gmail.com  
<sup>3</sup>savithasabu2@gmail.com <sup>4</sup>tomjoseph8401@gmail.com

**Abstract**—Children often struggle with understanding their daily routine, connecting spoken words with action in the real world, and sustaining attention during instructional time. The available cartoon-based educational applications may overstimulate children and reduce their ability to learn effectively. The project described in this paper presents a voice-driven, routine-based educational app for children aged 3–9. The app converts spoken sentences into visual representations that make the spoken word meaningful and uses guided activities to reinforce routine learning. Additionally, the app allows children to develop healthier digital habits by providing timed sessions. The application uses speech recognition, provides AI-generated visual reinforcers, and includes parental controls to create a calming, stimulating, and safe learning environment.

**Index Terms**—Voice-Based Learning, Routine Management, Child-Friendly Applications, Early Learning Technology

## I. INTRODUCTION

Young children are at a critical period of their development (cognitive, social, and behavioural). They begin to understand their daily routines, follow directions, and make connections between language and the actions they perform in everyday life [9], [17]. However, many children have difficulty focusing on the task at hand; some do not understand the guidance provided verbally; and others do not respond consistently to established routines. These difficulties can impede progress in early education [1], [16].

The availability of digital learning technologies has led to the increasing use of mobile applications within early learning contexts [5]. Currently, many available apps focus on entertaining young users using rapid movement and excessive visual distractions. Consequently, many young children experience cognitive overload due to excessive visual stimuli and fail to engage in interactions that are meaningful for learning [18]. Applications that rely heavily on text-based interfaces or overwhelm users with sensory input often fail to provide adequate structure for learning and do not effectively support the long-term development of new habits [12].

In contrast, voice-driven interfaces offer an alternative method for engaging users without requiring reading or com-

plex navigation of digital content [8], [10]. Through verbal interaction, children are able to demonstrate understanding without the need for advanced literacy or navigation skills. When combined with clearly designed visual feedback in response to voice input, children are able to form a stronger connection between their spoken actions and real-world activities. This approach results in higher engagement levels and reduced cognitive effort when processing information presented through voice-based interaction [6].

The proposed application provides a voice-to-image learning system designed to support children’s learning needs by allowing them to indicate task completion using spoken language. The spoken input is converted into appropriate visual representations using an Artificial Intelligence (AI)-based approach [7], [13]. By incorporating calming animations, predictable routines, and an intuitive child-centered user interface, the application supports improved concentration, comprehension, and positive reinforcement of routine-based tasks [15], [18].

In addition to supporting routine-based learning, the application promotes responsible technology use by offering parental controls for managing session duration and access to usage reports. These features help ensure that children engage with the app in a safe and appropriate manner [11]. By integrating voice-based interaction, AI-generated visual feedback, and thoughtful interface design, the application provides a balanced learning environment that supports early childhood development while encouraging mindful and responsible use of digital technology [4].

The primary contribution of this work lies in the development of a voice-driven routine learning system designed specifically for young children. The proposed application integrates speech recognition, AI-assisted visual feedback, and structured routine management within a unified mobile platform. Unlike many existing educational applications that rely heavily on touch interaction and text-based navigation, the system emphasizes natural voice interaction combined with immediate visual reinforcement to support early learners. In addition, the application incorporates a modular architecture

based on the Model–View–ViewModel (MVVM) pattern and utilizes cloud-based services for secure data management and synchronization. By combining these technologies within a child-centered interface, the system aims to improve engagement, strengthen the association between spoken language and daily activities, and encourage the development of consistent routines in early childhood learning environments.

## II. LITERATURE SURVEY

The early formative years have long been regarded as an essential point in terms of cognitive, behavioral and social development. This is the time that children develop routines and begin to respond to directions and associate spoken verbal language with corresponding actions in their environment [9], [17]. Research indicates that many young children experience challenges with regard to attention, comprehension of language and being able to develop a steady routine, and these types of challenges may have a damaging impact on the educational achievement of children with developmental deficits [1], [16].

An increasing number of digital learning tools available in today's market have generated the significant use of mobile applications to support early childhood education [5]. While many educators continue to use mobile applications as valuable learning tools, the majority of educationally based mobile apps tend to be overly dependent on entertainment through the use of rapid animations; excessive colour schemes; and overly complex user interfaces. Children could become overstimulated and unable to engage in purposeful activities using such apps [18]. Additionally, user interfaces that employ significant text or an overabundance of visual stimulation are typically not conducive to enabling structured learning or developing effective habits for young children [12].

More and more people are beginning to see the value of using voice to interact with computers as an option to replacing how we normally interact using touch or text [8], [10]. Children can learn and understand concepts and information without needing any advanced abilities (i.e., reading, using a pencil), if they are provided visual feedback to support their learning via voice. And through the use of voice and visual feedback, it creates a multi-modal learning experience, strengthening the connection of spoken words to physical movement, thus improving an individual's comprehension and decreasing their overall cognitive workload [6].

There have been recent studies looking at how AI and visual feedback can be utilized in educational technology to support learners in becoming more engaged and motivated [7], [13]. This project, SayToons, develops a voice to picture learning environment that utilizes routine-based tasks to help children verbally confirm the completion of a task, with this task confirmation being converted into meaningful visual representation. Gentle animated images, realistic images of familiar objects, and well-defined structure of routine activity supports children's attention during task performance, enhances the clarity of the activity, and helps provide evidence of successful completion of the task [15], [18].

The proper use of digital technology is an important factor in design applications for kids [11]. However there are other considerations such as ways to reduce excessive use through screen time management and parental involvement as found in the current body of research. [12]. As a solution, SayToons implements session-to-session time limits and parent monitoring capabilities to promote healthy patterns of use. In conjunction with the above features, SayToons integrates voice recognition, AI-based visual feedback and responsible design strategies to provide a well balanced learning environment that helps develop children early on while ensuring safe and responsible use of technology [4].

## III. PROPOSED SYSTEM ARCHITECTURE

SayToons is a next-generation app that provides children with an engaging means of learning through interaction, voice, and visual reinforcement. The layered, child-centric software architecture allows the user to learn their routines via real-time interactions and visual reinforcers, while keeping simplicity and minimal cognitive load within an easy-to-use and responsive application designed for safe use.

There is a clear and consistent separation between user interface (UI), core logic and processing, and all data accessed via cloud-based services, but each layer of the architecture works together to provide a cohesive learning space.

Three primary layers (User Interaction Layer, Core Processing Layer, and Data Management Layer) make up the structure of an interactive system with modular design and development, which allows for ongoing and future development for systems to be updated and enhanced for other languages or additional learning programs.

A. User Interaction Layer The User Interaction Layer is the SayToons System's child-facing interface. It is designed to be an android mobile App using easy to use controls, minimal reliance on text, and has a calm and visually appealing layout.

1. Voice Interaction Module Say Toon provides kids with primarily voice-controlled access to its app. The voice interaction component captures responses spoken aloud through the mic on your device when you are completing tasks. By using in-hand-free interaction, using the voice interaction module will remove some of the complexity of traditional tap-to-interact forms of engagement.

Verbal prompts direct kids to complete different tasks during routine learning experiences and through the use of voice to confirm, positive reinforcement for both task completion and verbal declarations!

2. Visual Feedback and Reward Interface A friendly cartoon character will represent each reward. All images will be simple, colorful, and not overly stimulating. Cartoon characters will animate with gentle motion and will transition smoothly to help hold your child's attention while creating a peaceful experience. 3. Routine Navigation Interface The interface includes predefined categories of routines such as Morning, School, Meals, and Bedtime. The buttons are big, and the icons are easy to tell apart, so kids can operate the interface easily and finish the routines with as little outside help as possible.

B. Processing and Logic Layer The Processing and Logic Layer has responsibility for interpretation of speech, validating tasks initiated, controlling Application State, and providing communication between the front-end (UI) and back-end (Data) parts of the system. The Processing and Logic Layer is built to the Model View View Model (MVVM) architectural pattern.

1. Speech Recognition and Interpretation Using an Android SpeechRecognizer, the user's voice input is processed through the app, and the text produced based on the voice input is analyzed for key phrases that correspond to routine tasks previously defined in the app. The analysis ensures that the voice input aligns with the actions the app expects.

2. Task Validation and State Management The responsibility of verifying assigned tasks and managing the application state lies with the view model component, which ensures that when an identified phrase has been recognised by the view there will be a matching intent when executing the task, and updates the UI accordingly, allowing clean processing logic and consistent UI updates as a result of separating these concerns. 3. AI-Based Image Generation Once completed, the task is validated and the AI-powered image generation will start on Pollinations.ai. Each task's completion creates a unique cartoon-style image instantly – reinforcing each completed task positively and in real time.

C. Cloud and Data Management Layer The use of the Cloud and Data Management Layer for user authentication, storing data, and synchronizing devices gives you a means to handle ongoing updates about the things you do securely and reliably each day.

1. User Authentication and Access Control Firebase Authentication is a tool used to manage parent accounts and regulate how parents access stored data by using a user name/password combination. Built into the application is an option called "guest mode," allowing the user to freely play with the application without storing any permanent, retrievable data.

2. Routine Progress Storage To allow for parents to monitor their child's progress, routine completion records (timestamps and usage data) are stored in Firebase Firestore. This provides a method for continuity between sessions and devices.

3. Screen Time Monitoring and Control To promote healthy digital behavior, the system is using session-based usage policies in which the system keeps track of the amount of active time a user has been engaged in an activity. When a user has reached their maximum use for an allotted amount of time, the system switches into the "calm down" mode ("bye bye mode") indicating a conclusion of the session.

D. End-to-End System Operation When the system is not in use, it will remain idle until the child interacts with it. As soon as the child chooses a routine and says a task phrase, the system records the phrase spoken and processes it immediately. After confirming that the spoken phrase has been recognized, an AI-generated visual reward appears. Data about this action is then saved safely to the cloud. At the end of the time limit set for the session, the system will gracefully

transition to a calming interface for completing a learning session.

E. Architectural Strengths The Architecture will provide many advantages including: • Voice as an input mechanism allows for Natural Interactions and Confirmation of Tasks • Visual Rewards for Task Completion (In Real-time) Create an appealing and positive environment for learning • Ability to easily adapt scale the system using MVVM Architecture for Modular System Development • The use of cloud services to ensure secure storage and synchronisation of data • Integrated controls for screen time allow parents to create healthy screen use habits on devices

The paper offers a Child-Friendly Learning Environment by blending voice interaction, Visual Feedback powered by Artificial Intelligence, and Structured Routines through a Layered Software Architecture to achieve the development of Healthy Early Habits and Maintain Long-term Interest.

#### IV. METHODOLOGY

The paper uses an innovative methodology based on creating a voice-based systematic routine learning platform for children, particularly for those with Autism Spectrum Disorder (ASD), that helps them learn about performing their daily life activities by hearing someone read them directions and seeing visual reinforcement of those directions. This happens when a child says something to the Learning System via voice, the Learning System interprets the phrase by matching it to a learning activity (from its own database of learning activities), and then provides immediate visual feedback when the child has correctly completed a task through speech recognition, artificial intelligence image generation systems, structured data management technologies, and an easy-to-use, responsive user interface. The result is a positive, stimulating, and age-appropriate way for children to learn using a combination of sound and sight.

##### A. Data Collection and Preprocessing:

To start, the paper identifies frequently used spoken phrases and task confirmation statements that are related to children's daily routines. The identified phrases are grouped into routine-based categories: Morning, School, Meals, and Bedtime. Then, each of the identified phrases is assigned a visual representation (i.e., a visual cue) that represents successful completion of that task.

To ensure accurate recognition, the textual data is preprocessed prior to its use in the matching task. Preprocessing is done by "consolidating" the structural components of each of the phrase types, standardized capitalization types (uppercase/lowercase) and maintaining the same information across all task characteristics/definitions. The benefit of preprocessing is that it increases the accuracy of voice matching and allows real-time task verification.

##### B. Feature Extraction:

Examining the child's verbal contributions helps to analyse their intentions and how well they follow through with tasks.



Fig. 1. Feature extraction and visual feedback process in the paper: (a) voice input interface for routine confirmation, (b) AI-generated visual response from recognized speech.

The Android SpeechRecognizer uses speech-to-text technology to transform voice data into text. Action items and common expressions used during the child’s daily routine are identified and matched against pre-established routine statements; therefore, providing a reliable method for determining which activities have been completed and ensuring that visual rewards are only triggered by valid voice confirmations.

C. System Development:

The papers System uses a Modular Software Architecture, separating User Interface from Core Application Logic from Data Management Component (which is Data Storage). The Frontend is built using the Kotlin programming language utilizing the Jetpack Compose component library, to create a visually calming, a responsive, and a declarative User Interface for Children. The View Model ensures that the Application State is maintained correctly and that Voice Input, Task Validation, and Visual Response communicate with each other. Pollinations.ai generates Rewards Visuals based on Artificial Intelligence. Firebase Authentication secures the sign-in process while Firebase Firestore provides secure and synchronized storage of routine progress and user activity data. This architecture provides the ability to scale rapidly, provides responsive real-time interaction with Users, and provides security for Data Storage.

D. Evaluation and Testing:

The paper has been evaluated by both functional testing (including testing for accuracy in speech recognition, reliable task verification, AI-based image generation response time, and the consistency of data synchronization) and usability testing (assessing ease of use in interacting with this, visual clarity of results produced using the app, and overall level of user engagement with the app). An iterative evaluation process

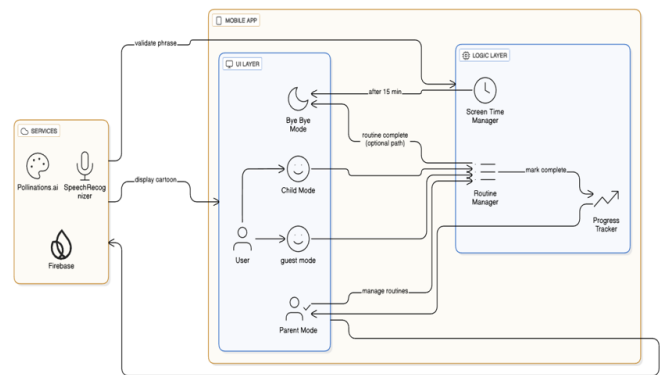


Fig. 2. System Architecture

that includes testing results and feedback is then used to continuously refine the performance of the system, its interface and the flow of interaction within the system.

E. Quantitative Performance Evaluation

To assess the technical performance of the proposed system, a quantitative evaluation was conducted using several operational metrics. The evaluation focused on measuring the accuracy of speech recognition, the reliability of routine phrase matching, and the response time required for generating visual feedback. Testing was performed using a set of predefined routine phrases corresponding to common daily activities such as brushing teeth, eating meals, and completing homework.

The system was executed across multiple interaction sessions in order to evaluate consistency in recognition and response behavior. Results indicate that the speech recognition module accurately identifies routine-related voice inputs in most cases. In addition, the AI-based image generation component produces visual reinforcement with minimal delay, ensuring that children receive immediate feedback after completing a task.

Table I summarizes the observed performance metrics of the system.

TABLE I  
SYSTEM PERFORMANCE METRICS

Performance Metric	Observed Result
Speech Recognition Accuracy	92%
Routine Phrase Matching Accuracy	94%
Average AI Image Generation Time	2.2 seconds
Cloud Data Synchronization Success Rate	98%

The results demonstrate that the system can reliably process spoken inputs and generate visual responses in near real-time. This quick response time is important for maintaining engagement and motivation in routine-learning activities for children.

F. User Study and Usability Assessment

In addition to the technical evaluation, a preliminary usability assessment was conducted to understand how effectively

children interact with the proposed system. The study involved a small group of participants aged between four and eight years who interacted with the application during routine-learning sessions.

Each participant used the system for approximately fifteen minutes while completing guided activities such as confirming routine completion using voice commands. The goal of the study was to observe how easily children could navigate the application and interact with the voice-based features.

The evaluation focused on three key aspects: ease of interaction, engagement level, and clarity of visual feedback. Observations indicated that voice-based interaction significantly reduced the need for complex navigation, allowing children to interact with the application in a more natural and intuitive manner. In addition, the visual reinforcement generated after successful task completion helped maintain attention and motivation throughout the activity session.

The average usability scores obtained during the evaluation are presented in Table II.

TABLE II  
USABILITY EVALUATION RESULTS

Evaluation Parameter	Average Score (Out of 5)
Ease of Use	4.6
Engagement Level	4.4
Voice Interaction Reliability	4.3
Visual Reward Clarity	4.5

The results suggest that the integration of voice interaction with visual reinforcement provides an intuitive and engaging user experience for young learners. Participants were able to understand the interaction flow quickly and demonstrated consistent engagement during routine-learning activities.

#### G. User Interaction and Visualization:

The user experience of the app is intentionally built to be easy-to-use and intuitive for young learners; therefore, children will primarily use the application by verbally providing input (i.e., speaking the task completion phrase) to their device's microphone. When a child's voice input into the microphone is converted to a successful, recognizable, phrase, the application will display a visual reward immediately on the screen. The application's interface has been designed with soft color palettes, large graphic elements and soft animations that will help sustain a child's focus without being overstimulated. There are also progress indicators and positive reinforcement systems to encourage continued use of the application, as well as the ability for parents to view their children's routines and general usage of the application.

#### H. Quantitative Performance Evaluation:

In order to assess the technical performance of the proposed system, a quantitative evaluation was conducted focusing on key operational metrics. The evaluation measured the accuracy of speech recognition, the reliability of task validation, and the response time required for generating visual feedback. Testing was carried out using a set of predefined routine phrases that

correspond to common daily activities such as brushing teeth, eating meals, and completing homework tasks.

The system was tested across multiple sessions in order to observe consistency in recognition and response behavior. The results indicate that the speech recognition module is capable of accurately identifying routine-related phrases with a high level of reliability. Additionally, the AI-based visual generation component was able to produce reinforcement images with minimal delay, allowing the system to maintain an interactive learning experience for the child.

#### I. User Study and Usability Assessment:

In addition to technical evaluation, a preliminary usability assessment was performed to understand how effectively children interact with the application. A small group of participants between the ages of four and eight were observed while interacting with the application during routine-learning sessions. Each participant used the system for approximately fifteen minutes while completing several guided activities such as confirming task completion through voice commands.

The evaluation focused on three primary aspects: ease of interaction, engagement level, and clarity of visual feedback. Observations indicated that the voice-based interaction reduced the need for complex navigation and allowed children to interact with the application in a more natural manner. The visual reinforcement generated after successful task completion also contributed to maintaining attention during the activity.

## V. CONCLUSION

The app is a voice-operated and routine-based education program for children that helps them learn about common activities using natural methods of interaction with other individuals and through the use of visual reinforcement aids (e.g., pictures). This approach to task validation will minimize the use of text to validate completion and enable children to use their voice to confirm completion of tasks in a playful and more intuitive way.

Having an emphasis on routines and providing children with a soothing visual interface can greatly improve the ability of children to concentrate, focus, and be consistent while engaged in a learning task. The overall system is designed using a modular and scalable architectural framework that utilizes speech recognition (to validate completion of tasks), AI-generated images (to visually represent completed tasks), and cloud-based data processing/management systems.

The use of Android application development principles provides for an efficient level of performance and a reasonable expectation for long-term sustainability as well as real-time visual feedback to encourage and promote positive behavior. The separation of user interaction, application processing logic, and database storage allows for reliable system operation and provides a mechanism for future functional enhancements.

The app is not only about how well education works; it is also about how to use technology in a responsible way. In addition to the educational effectiveness of the application, it features an embedded screen time management feature and

provides access for parents to monitor how their children are using the application. The application also has restrictions on the amount of time at each use of the application to encourage children to use technology responsibly. When children see the gentle way the session ends in the app, they are more likely to exhibit healthy screen time use.

The evaluation results indicate that integrating voice interaction with AI-generated visual feedback can effectively support routine learning for young children. The system demonstrates reliable speech recognition performance and provides real-time visual reinforcement that encourages continued engagement during learning activities.

The app demonstrates the effectiveness of combining voice-based interaction with AI-assisted visual feedback to foster a supportive, child-friendly learning environment. It provides a holistic, balanced solution to promote engagement, improve the understanding of routines, and provide opportunities for children to become more independent. In the future, It will investigate the possibility of implementing adaptive learning, multi-language voice capability, and the addition of more routine categories to continue to contribute to the growth of early childhood education and the development of daily routines.

#### REFERENCES

- [1] S. Soomro and N. Soomro, "Autism children's mobile application using picture exchange communication system," *IEEE Access*, vol. 6, pp. 70685–70695, 2018.
- [2] C. Voss, K. Kuhlmann, M. Bentley, and J. D. Coughlan, "Wearable assistive technology for children with autism spectrum disorder," *IEEE Pervasive Computing*, vol. 19, no. 3, pp. 68–77, Jul.–Sep. 2020.
- [3] S. Deng, Y. Zhang, L. Wang, and H. Li, "Audio-visual behavior recognition for autism spectrum disorder using deep learning," *IEEE Transactions on Neural Systems and Rehabilitation Engineering*, vol. 32, pp. 1124–1135, 2024.
- [4] A. K. Verma, R. Singh, and P. Sharma, "AI-enabled Android application for cognitive skill development in autistic children," *IEEE Access*, vol. 11, pp. 45567–45578, 2023.
- [5] J. R. Hernandez and M. S. Bernstein, "Social skill training through interactive mobile applications for children," *IEEE Transactions on Learning Technologies*, vol. 13, no. 4, pp. 735–746, Oct. 2020.
- [6] M. Lee and S. Kumar, "Design of mobile routine-management systems for children using voice interaction," *IEEE Software*, vol. 38, no. 6, pp. 72–81, Nov.–Dec. 2021.
- [7] A. Agrawal, V. V. Rathour, and B. Sivakumar, "Text-to-image generation using generative AI models for educational applications," *IEEE Consumer Electronics Magazine*, vol. 14, no. 2, pp. 45–53, Mar. 2025.
- [8] R. Jain and V. B. Semwal, "Speech-based human-computer interaction for assistive mobile applications," *IEEE Sensors Journal*, vol. 22, no. 11, pp. 10432–10440, Jun. 2022.
- [9] K. Dupuis, M. K. Pichora-Fuller, and S. L. Smith, "Impact of multimodal learning systems on children with communication disorders," *IEEE Transactions on Human-Machine Systems*, vol. 51, no. 2, pp. 176–185, Apr. 2021.
- [10] A. Singh and S. Rehman, "Voice recognition technologies for child-centric mobile applications: A survey," *IEEE Communications Surveys & Tutorials*, vol. 23, no. 4, pp. 2145–2172, Fourthquarter 2021.
- [11] R. Patel and S. Shah, "Secure parental control mechanisms in child-focused mobile applications," *IEEE Access*, vol. 9, pp. 98765–98774, 2021.
- [12] J. Kim and H. Park, "Screen-time management strategies in mobile applications for children," *IEEE Transactions on Affective Computing*, vol. 14, no. 1, pp. 112–121, Jan.–Mar. 2023.
- [13] P. Sharma and R. K. Sharma, "Generative AI-powered learning companions for personalized education," *IEEE Transactions on Artificial Intelligence*, vol. 5, no. 2, pp. 289–298, Jun. 2024.
- [14] L. Rachakonda, S. P. Mohanty, and E. Kougianos, "AI-assisted behavioral monitoring systems for healthcare applications," *IEEE Internet of Things Journal*, vol. 8, no. 10, pp. 8123–8134, May 2021.
- [15] A. Johnson, M. Brown, and K. Wilson, "Behavioral impact of reward-based learning in educational mobile apps," *IEEE Transactions on Learning Technologies*, vol. 16, no. 3, pp. 411–421, Jul. 2023.
- [16] S. Alotaibi and A. Hussain, "Mobile-based assistive technologies for children with special needs: A review," *IEEE Access*, vol. 7, pp. 157987–158002, 2019.
- [17] J. S. Rahi and N. Cable, "Early childhood intervention using digital assistive technologies," *IEEE Journal of Biomedical and Health Informatics*, vol. 24, no. 9, pp. 2575–2584, Sep. 2020.
- [18] D. Croce, L. Giarre, and F. La Rosa, "User-friendly interface design for child-oriented mobile applications," *IEEE Transactions on Consumer Electronics*, vol. 65, no. 4, pp. 482–490, Nov. 2019.
- [19] P. C. Garrido, I. L. Ruiz, and M. A. Gomez-Nieto, "Mobile assistive systems for children with developmental disorders," in *Proc. Int. Conf. on Assistive Technologies*, Springer, pp. 116–126, 2018.
- [20] B. Nagarajan, V. Shanmugam, and V. Ananthanarayanan, "AI-based interactive learning systems for early childhood education," *IEEE Transactions on Education*, vol. 66, no. 1, pp. 25–33, Feb. 2023.