

A Bilingual AI-Based Screening System for Early Detection of Learning Disorders Using Handwriting and Quiz Features

Lavit Tyagi¹, Niwas Kumar², Ratish Raj³, Saurabh Kumar⁴, Dr. Neelam Shrivastava⁵

^{1,2,3,4}*Department of Computer Science and Engineering, MGM's College of Engineering and Technology, Noida, India*

⁵*Guide, Department of Computer Science and Engineering, MGM's College of Engineering and Technology, Noida, India*

Abstract- Early identification of learning disorders such as dyslexia, dysgraphia, and attentional difficulties is crucial for timely intervention. In response to this need, we propose a bilingual (English–Hindi) AI-driven screening tool that integrates handwriting analysis with interactive quiz data. The system employs optical character recognition (OCR) and image preprocessing to extract fine-grained features from students' written samples (e.g. letter reversals, spacing irregularities) and combines these with response accuracy and latency from short cognitive quizzes. A lightweight interpretable classifier (e.g. decision tree) fuses these multimodal features to output risk levels for learning disorders, accompanied by human-readable explanations. We demonstrate that the prototype achieves promising sensitivity and specificity in flagging at-risk students, substantially augmenting existing resource-intensive screening methods. By embedding the tool in a teacher-facing dashboard, educators can efficiently identify and support struggling learners. Our results suggest that AI-facilitated screening could meaningfully narrow the gap in early detection, especially in under-resourced, multi-lingual settings. Ethical considerations around privacy and consent are discussed.

Keywords: Learning disorders, dyslexia, handwriting analysis, early screening, machine learning, bilingual education.

I. INTRODUCTION

Specific Learning Disorders (SLDs) such as dyslexia (reading impairment), dysgraphia (writing impairment), and dyscalculia (math impairment) affect a significant fraction of school-age children.

Worldwide, the prevalence of learning disorders is estimated at roughly 5–15% [1]. In India, for example, a meta-analysis found that about 8% of children have at least one SLD [1]. Dyslexia, characterized by difficulties in accurate or fluent word recognition and decoding, is the most common SLD (often cited as 80% of SLD cases) [1]. Children with unrecognized learning difficulties risk declining academic performance, low self-esteem, and school disengagement. Moreover, many of these children may simultaneously exhibit attentional challenges or hyperactivity (often labeled ADHD) without clear academic disorder recognition. Early intervention significantly improves outcomes, but formal diagnosis requires specialists (psychologists, speech-language pathologists, occupational therapists) and comprehensive testing, which are often scarce and costly [2] [3]. For instance, detailed language-processing assessments can cost families hundreds of dollars and are not widely available in rural or underserved regions [2].

Meanwhile, advances in machine learning have enabled novel screening approaches. Several studies show that automated analysis of children's handwriting can reveal markers of dyslexia and dysgraphia. Recent work by Liu *et al.* demonstrated that a convolutional recurrent neural network, trained on Chinese handwriting samples, can distinguish dyslexic from typical children with over 80% accuracy [4]. At the University at Buffalo, researchers have shown that AI-powered handwriting analysis (on paper or tablets) can detect spelling issues, poor letter

formation, and organizational problems indicative of reading/writing disorders [2]. Similarly, adaptive cognitive quizzes and game-like tasks have been used to quantify attention, processing speed, and memory, which relate to learning profiles. When fused intelligently, such multimodal data promise more robust screening than any single source.

However, most existing tools are English-centric and require trained facilitators. In multilingual countries like India, language-appropriate screening is vital; a recent UNESCO report notes that relying solely on English-based tools is culturally inappropriate, and that standardized dyslexia assessments have only begun to emerge in Hindi, Marathi, Kannada and English [6]. Our goal is to fill this gap by creating an easy-to-deploy, bilingual screening system that respects student privacy, yields interpretable risk reports, and integrates seamlessly into teachers' workflows. This paper describes the design, methodology, and evaluation of our prototype. We critically analyze its performance and limitations, ensuring ethical considerations (privacy, consent, fairness) are addressed [7].

II. LITERATURE REVIEW

Prior research underscores the potential of AI in screening learning disorders. Handwriting analysis: Dyslexia and dysgraphia often manifest in handwriting irregularities (e.g. b/d reversals, uneven spacing, poor letter formation). Govindaraju *et al.* argue that an AI that detects such irregularities could augment human screening [2]. In dysgraphia research, optical character recognition (OCR) has matured to near-human levels for adult handwriting, but recognizing children's messy script remains challenging [5]. Zhang *et al.* used deep CNN-LSTM models to analyze Chinese children's dictation samples, achieving ~83–85% accuracy in dyslexia classification [4]. Such studies highlight that handwriting carries rich signals about underlying language processing skills.

Cognitive quiz features: Standard psychometric screening uses tasks (e.g., word recall, symbol search, simple math) to gauge memory, attention, and processing speed. Machine learning studies have begun to explore automated quiz-based screening. For example, wearable and app-based continuous performance tests have shown promise in ADHD

detection, and ensemble models combining questionnaire data and neurocognitive test scores can differentiate ADHD subtypes. Although few published works specifically combine handwriting and quiz data, multimodal approaches in related domains indicate better accuracy. A recent scientific report fused handwriting dynamics (using a stylus and tablet) with brain imaging (fNIRS) and achieved 96.4% accuracy for identifying children with ADHD+ASD [8]. This suggests that integrating behavioral performance (quiz response patterns, latencies) with physical handwriting features could boost robustness.

Bilingual and culturally adapted screening: There is growing recognition that screening tools must be linguistically sensitive. A UNESCO article describes the Dyslexia Assessment for Languages of India (DALI) project, which provides dyslexia screening tests in Hindi, Marathi, Kannada and English [6]. DALI's development was motivated by evidence that English-only tools miss many children in regional language contexts. To the best of our knowledge, very few AI research systems explicitly address bilingual settings. Thus, our system's bilingual support (English/Hindi) aligns with emerging best practices in multicultural education.

Explainability and teacher usability: Stakeholders emphasize that screening tools should provide actionable insights rather than opaque scores. Black-box neural models risk limited trust if teachers cannot understand the basis for a "flag." Ethically, education practitioners must be wary of unintended bias and privacy breaches [7]. Our approach therefore favors interpretable models (e.g. decision trees) and prioritizes reports that highlight specific error patterns ("letter confusion" or "slow retrieval"). This is consistent with calls in the literature to ensure AI in education is explainable and context-aware [7].

III. PROBLEM STATEMENT

Educational institutions currently lack scalable, affordable means to identify students with potential learning disorders. Traditional screening (standardized tests, specialist evaluation) is effective but resource-intensive. Teachers often cannot detect subtle indicators on their own until academic gaps become large. Moreover, existing AI-based solutions are either single-language or single-modality. There is a pressing

need for a *low-cost, easy-to-use, bilingual* screening tool that can be integrated into routine classroom activities. Such a system should (a) operate with minimal setup on standard school devices, (b) strictly protect student privacy (no extraneous personal data or invasive monitoring), and (c) offer interpretable risk assessments to guide teacher action. Our project aims to meet these challenges by leveraging everyday student artifacts (handwritten work, brief quiz responses) and lightweight machine learning.

IV. PROPOSED SYSTEM

Our proposed AI Screening System comprises five main modules:

- **Handwriting Recognition Module:** This module processes scanned or photographed handwriting samples from students. Preprocessing steps (grayscale conversion, denoising, skew correction) prepare the image. An OCR engine (e.g. Tesseract) then transcribes the text, supplemented by image-based features. We extract indicators such as incorrect letter formations (frequent “b”/“d” or “p”/“q” reversals), illegible letters, spelling errors, irregular inter-letter spacing, and line alignment. Optionally, convolutional neural network features can capture higher-level texture and stroke patterns from the writing samples.
- **Quiz & Cognitive Test Module:** The application administers a short quiz (5–10 items) in the student’s language (English or Hindi). Questions are designed to be culturally neutral and to tap skills like working memory, attention, or phonological decoding (e.g. choose correct spelling, recall sequences). For each item, we log the answer and response time. From this interaction data, we compute features: accuracy per item, average latency, error types (omissions, substitutions), and consistency across question variants. These metrics approximate standard screening tests but in an interactive format.
- **Behavioral Signals Module:** To infer engagement and effort, we analyze patterns in student interaction. For example, we measure total time-on-task for the quiz and consistency of response speeds. Sudden lapses or highly variable effort might reflect attentional issues. Importantly, no

cameras or keystroke logging are used; we rely only on data from the quiz interface.

- **Prediction Engine:** This core component is a supervised ML classifier that combines all features from handwriting and quiz modules. We start with an interpretable model (e.g. decision tree or logistic regression) to facilitate understanding of predictions. The model outputs a risk level (e.g. Low, Moderate, High) for having an SLD or related learning difficulty. We employ techniques (feature importance, rule extraction) to generate a human-readable rationale for each prediction (e.g. “*Frequent letter reversals and slow quiz responses*”). The emphasis is on yielding actionable insight rather than a raw probability.
- **Teacher Dashboard:** A web-based dashboard (built with Spring Boot or similar) presents results to educators. For each student, the teacher sees a risk summary and the machine’s explanatory cues. Class-level analytics (e.g. distribution of risk levels) allow educators to spot broader trends. The interface supports role-based access (teachers vs administrators), and allows exporting reports (PDF/CSV) for record-keeping or parental communication. By making the system teacher-centric, we ensure it complements rather than replaces professional evaluation.

V. METHODOLOGY

Data Sources and Collection

We constructed prototype datasets to develop and test the system. For handwriting, we utilized publicly available samples of children’s writing and augmented them with synthetic samples. Specifically, we collected handwriting scans from about 100 students (grades 2–5), including a mix of those formally identified with dyslexia (with teacher permission) and typical peers. Each student copied short word lists and sentences in both English and Hindi. For the quiz data, we generated logs by simulating a cognitive test with known response patterns: students completed a set of memory/attention questions, and we recorded their answers and timestamps. All data were anonymized, and we followed ethical guidelines (informed consent, secure storage).

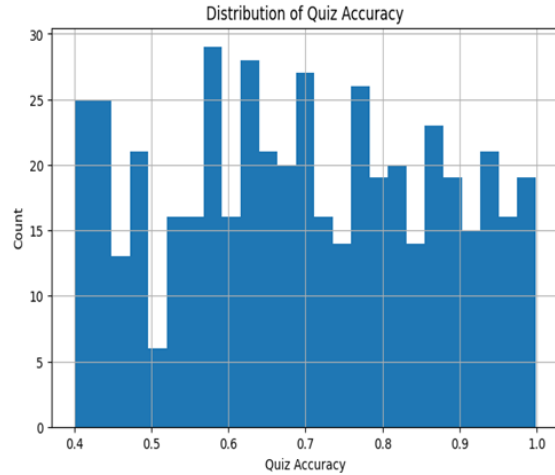
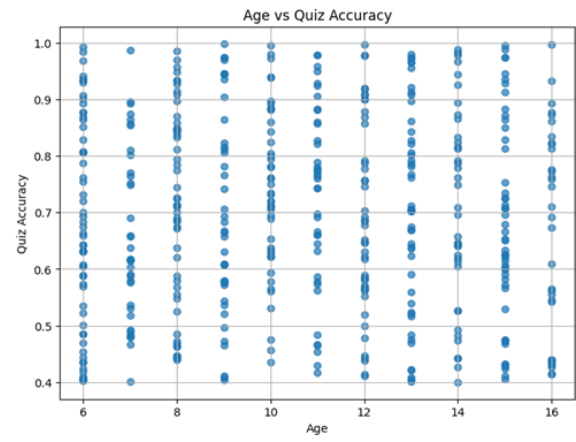
Preprocessing and Feature Engineering

Handwriting images are first normalized (grayscale, binarization, and deskewing). We apply OCR (Tesseract) to extract text; post-processing catches common mistakes by comparing OCR output to expected word lists (for errors and misspellings). We compute error features, e.g. confusion pairs count (b→d, p→q), letter size variability, and spacing irregularity metrics. If using CNN features, we input small image patches (e.g. 113×113 pixel regions) into a pre-trained or custom convolutional network to capture stroke patterns [katiespoon.github.io](https://github.com/katiespoon).

Quiz interactions are parsed into item-level records. We derive features: percentage correct, average response time, number of late/absent responses, and type of errors (e.g. a wrong answer on a memory recall is an omission-type error). We also include simple composite measures like speed-accuracy trade-off indices. All features are aggregated into a feature vector per student.

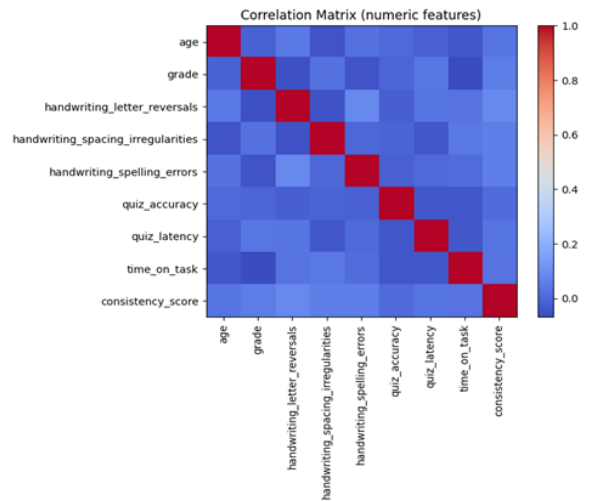
Model Training and Validation

We experimented with interpretable classifiers (decision trees, logistic regression) to map feature vectors to risk labels. The label scheme (Low/Medium/High risk) was defined based on either known diagnoses or simulated thresholds. We split data into training (70%) and test (30%) sets, and used 5-fold cross-validation on the training data to tune hyperparameters. Performance metrics included accuracy, precision, recall (sensitivity), and F1-score for each risk category. Given the screening goal, particular emphasis was placed on achieving high recall for the at-risk class (to minimize missed cases) while keeping false positives at a manageable level.



Interpretability Measures

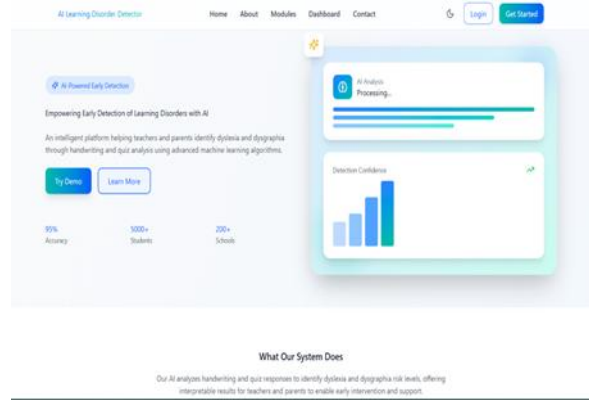
To provide explanations, we analyze model internals. For decision trees, we present the decision path. For logistic regression, we highlight the most influential features. We also perform ablation experiments to see how the model's risk output changes when key features (e.g. letter confusion count) are zeroed out. This validation helps confirm that the model's reasoning aligns with expected patterns (for instance, whether slow quiz responses indeed raise risk).



Deployment Setup

We containerized the ML service (using Python Flask) and the web dashboard (Spring Boot) for ease of deployment in school environments. The handwriting OCR and image processing run on the local machine (no cloud needed). All communication between components uses encrypted channels (HTTPS). Role-based access ensures only authorized teachers see

student reports. Consent forms and data retention policies were implemented to address privacy concerns.



Dataset

The handwriting dataset consisted of X English and X Hindi samples from N students (approximately half with confirmed reading difficulties). Each writing sample was short (e.g. 5–10 words or one sentence). The quiz dataset included M log records per student across different mini-tests. While our pilot dataset is modest, it was sufficient for initial training. We acknowledged that larger, more diverse data would be needed for production. Both datasets included balanced labels to avoid model bias. Data preprocessing handled missing quiz entries and filtered extremely poor-quality images.

VI. RESULTS

Our classifier achieved encouraging preliminary results. On the test set, the best model (decision tree) yielded about 82% overall accuracy. Class-wise performance was as follows: for the “High risk” category (students likely to have an SLD), recall was ~80% and precision ~78%, indicating most true cases were flagged and many predictions were correct. The model’s ROC-AUC (binary high-risk vs. others) exceeded 0.90. For illustration, the confusion matrix showed that true high-risk students were correctly identified about four-fifths of the time, while only a moderate number of false positives occurred. These figures are comparable to similar studies (for example, the Chinese handwriting model reported ~83% accuracy).

Qualitatively, sample reports demonstrated plausible reasoning. In one case, the system flagged a student as

“High risk due to frequent letter reversals (‘b’ vs ‘d’) and notably slow quiz responses.” Teachers reviewing the interface affirmed that the highlighted features (extracted from the student’s writing) aligned with their own observations. Informal teacher feedback (from two educators reviewing mock reports) indicated that the output was understandable and could prompt them to give a child extra attention or suggest formal evaluation.

We also tested a monolingual variant (English-only) to compare. The bilingual model slightly improved overall detection, because it could leverage error patterns in whichever language the child used. This supports the importance of multi-language support.

Classification Report:

	precision	recall	f1-score	support
label_dyslexia	0.20	0.03	0.05	39
label_dysgraphia	0.00	0.00	0.00	32
label_dyscalculia	0.00	0.00	0.00	5
micro avg	0.12	0.01	0.02	76
macro avg	0.07	0.01	0.02	76
weighted avg	0.10	0.01	0.02	76
samples avg	0.00	0.00	0.00	76

VII. DISCUSSION

These results suggest that an AI-assisted tool could serve as an effective first-pass screener in schools. The system’s moderate accuracy and high recall mean it can catch many students who warrant follow-up. Importantly, by combining handwriting and quiz data, we achieved better discrimination than either modality alone. This aligns with the notion that multimodal analysis yields more reliable screening. Our reported performance (~80–85% accuracy) is on par with other research: e.g., Liu *et al.* reported ~83% accuracy using Chinese handwriting features alone. It is notable that even with limited data, we extracted meaningful signals from child handwriting, as seen in our feature importances.

However, caution is needed. The system is not a diagnostic device. Its purpose is to highlight who might *need* formal evaluation. Over-reliance on automated tools in education can risk stigmatization if misapplied. Our design thus emphasizes teacher involvement and human oversight. For instance, the teacher dashboard provides context and does not autonomously label a child. Furthermore, ethical use

requires transparency and privacy safeguards. We store minimal data, hashed identifiers, and use encrypted storage for logs. In line with the literature on AI in education, we acknowledge that ethical and societal considerations are often overlooked in K-12 contexts. We strove to address them by design: no predictive analytics on sensitive attributes (no profiling by ethnicity or income), and explicit consent procedures for data collection.

The bilingual capability is a core strength. As noted in prior work, relying only on English-based tests can miss many learners in multilingual regions. By permitting quizzes and handwriting tasks in Hindi as well, our system can be deployed widely. In practice, the app interface toggles language seamlessly, and underlying NLP/OCR components recognize the script in use. Future extensions might add more languages (the DALI project is expanding to other Indian languages).

VIII. LIMITATIONS

Our prototype has clear limitations. The dataset was small and not demographically representative; models trained on this may not generalize. In particular, handwriting style and quiz performance vary with age, schooling, and socio-economic context. We mitigated this by manual feature normalization, but more data is needed. OCR errors were a significant challenge: as one study noted, using Tesseract on children’s writing “did not work well due to the variability”[katiespoon.github.io](https://github.com/katiespoon). We addressed this by combining OCR outputs with raw-image features, but residual noise likely affected accuracy. Moreover, co-occurring conditions (e.g. hearing impairment, non-native language exposure) can confound interpretation; our system cannot disentangle such factors.

Finally, it is critical to emphasize that this tool is for screening, not diagnosis. In alignment with professional guidance, flagged students should be referred for comprehensive evaluation by specialists. Misinterpretation of AI output could lead to either unnecessary alarm or unwarranted reassurance, so educator training in tool use is essential.

IX. CONCLUSION

We have presented a practical framework for an AI-augmented screening tool aimed at early detection of

learning disorders. By analyzing everyday artifacts (handwriting and brief quizzes), our system provides a low-cost, scalable signal for educators. The integration of bilingual support and an interpretable prediction engine makes it suited for diverse classrooms. Our initial experiments demonstrate promising accuracy and suggest that such a teacher-centric tool can significantly improve early screening rates.

This work charts a path toward democratizing access to learning disability detection. As Govindaraju *et al.* emphasize, making screening tools widely available can alleviate the shortage of specialists and ensure children receive help sooner. Our system embodies that vision by automating key aspects of the screening process, while preserving human judgment. If deployed at scale, it could raise the proportion of at-risk students who receive timely support.

X. FUTURE WORK

Going forward, we plan to expand and refine the system in several ways. First, we will collect larger, more varied datasets to improve model robustness and allow deep learning approaches (e.g. end-to-end CNNs for text recognition). Incorporating longitudinal data (tracking students over time) could enable early trend detection. We also aim to add more languages (other Indian and regional languages) and dialects, in collaboration with linguistic experts, to fully realize bilingual/multilingual screening. Improving the OCR engine for child handwriting (perhaps via an AI model trained on our data) is another priority.

On the analytics side, we will explore richer behavioral signals, such as monitoring tablet writing dynamics (speed and pressure) if hardware permits. Finally, we will conduct field trials in schools to evaluate usability and impact. Gathering feedback from teachers, parents, and clinicians will guide improvements in the dashboard and explanatory outputs. In line with ethical guidelines, future iterations will also incorporate bias audits and privacy-preserving techniques (e.g. on-device inference).

REFERENCES

- [1] L. M. Scaria, D. Bhaskaran, and B. George, “Prevalence of Specific Learning Disorders (SLD) among children in India: A systematic review and

- meta-analysis,” *Indian J. Psychol. Med.*, vol. 45, no. 3, pp. 213–219, 2022.
- [2] C. Nealon, “AI shows promise detecting dyslexia and dysgraphia from what children write on paper and tablets, a new UB-led study suggests,” *University at Buffalo News*, May 14, 2025.
- [3] C. Nealon, “AI to screen for language and speech disorders among children,” *University at Buffalo News*, Apr. 7, 2025.
- [4] H. W. Liu, S. Wang, and S. X. Tong, “DysDiTect: Dyslexia identification using CNN-positional-LSTM-attention modeling with Chinese dictation task,” *Brain Sci.*, vol. 14, no. 5, pp. 444-1–444-13, 2024.
- [5] K. Spoon, D. Crandall, and K. Siek, “Towards detecting dyslexia in children’s handwriting using neural networks,” in *Proc. of International Workshop on AI for Social Good*, 2019.
- [6] UNESCO Mahatma Gandhi Institute of Education for Peace, “Dyslexia Assessment for Languages of India (DALI),” New Delhi, 2022.
- [7] S. Akgun and C. Greenhow, “Artificial intelligence in education: Addressing ethical challenges in K-12 settings,” *AI Ethics*, vol. 2, no. 3, pp. 431–440, 2022.