



# International Journal Advanced Research Publications

# ATTENTION-ENHANCED CNN-LSTM ARCHITECTURE FOR REAL-TIME SMARTPHONE DISTRACTION DECETION IN SYNCHRONOUS ONLINE LEARNING

S. Vimala\*1, Dr. G. Arockia Sahaya Sheela<sup>2</sup>

\* <sup>1</sup>PhD Scholar(Full Time), Department of Computer Science, St. Joseph's College(Autonomous), Tiruchirappalli -2, Affiliated to Bharathidasan University, Tamil Nadu, India.

<sup>2</sup>Assistant Professor, Department of Computer Science, St. Joseph's College(Autonomous), Tiruchirappalli -2, Affiliated to Bharathidasan University, Tamil Nadu, India.

Article Received: 17 Octomber 2025, Article Revised: 06 November 2025, Published on: 26 November 2025

\*Corresponding Author: S. Vimala

1PhD Scholar(Full Time), Department of Computer Science, St. Joseph's College(Autonomous), Tiruchirappalli -2, Affiliated to Bharathidasan University, Tamil Nadu, India. DOI: https://doi-doi.org/101555/jjrpa.2891

## **ABSTRACT**

The proliferation of smartphones during virtual classroom sessions has created substantial challenges for maintaining learner focus and educational effectiveness. We developed an attention-augmented hybrid deep learning system combining Convolutional Neural Networks with Long Short-Term Memory units (CNN-LSTM-Attn) to identify smartphone-induced distraction behaviors through device usage patterns and motion sensor analytics. The convolutional layers isolate brief temporal signatures from touchscreen interactions and device movement, whereas the recurrent components model extended behavioral sequences. An attention weighting mechanism emphasizes temporal segments most predictive of attention drift, simultaneously enhancing classification reliability and model transparency. Our study enrolled 120 university students participating in online coursework, yielding approximately 180 hours of behavioral recordings. The preprocessing workflow incorporated standardization procedures, synthetic data generation, and anonymization protocols. Performance evaluation demonstrated 92.4% classification accuracy, F1-score reaching 0.91, and inference latency averaging 1.2 seconds—exceeding both conventional CNN-LSTM configurations and classical machine learning approaches. The attention component enabled

visualization of distraction-indicative moments while preserving computational efficiency suitable for edge deployment. These outcomes establish feasibility for privacy-respecting, real-time intervention systems supporting sustained engagement in virtual educational settings.

**KEYWORDS:** Smartphone behavioral analytics, Attention weighting mechanisms, Hybrid CNN-LSTM architecture, Real-time behavioral monitoring, Virtual learning engagement.

### I. INTRODUCTION

Modern smartphones serve multiple roles spanning communication, leisure, and educational purposes, yet their constant availability presents considerable obstacles to concentration during virtual learning activities. Learners commonly engage in parallel tasks—alternating between academic applications and entertainment platforms—resulting in fragmented attention and compromised educational achievement. Implementing automated distraction recognition with immediate feedback capabilities holds promise for strengthening learner engagement and improving knowledge acquisition outcomes.<sup>[1]</sup>

Traditional approaches for identifying off-task behavior rely predominantly on application activity records or visual monitoring systems. While application logs provide accessibility, they offer limited behavioral context; conversely, camera-based surveillance introduces significant privacy and ethical complications. Addressing these constraints, our investigation proposes an unobtrusive sensor-driven methodology utilizing touchscreen dynamics, inertial measurement units (accelerometer and gyroscope), and display activation events to characterize distraction patterns without accessing personal communications or visual information. [3]

Our hybrid neural architecture capitalizes on complementary capabilities of convolutional and recurrent processing paradigms. Convolutional layers identify localized temporal-spatial characteristics within brief sensor data segments, while recurrent layers establish relationships spanning multiple temporal windows. The attention mechanism optimizes this process by selectively emphasizing temporally informative segments, effectively discriminating between academically productive device usage and diversionary smartphone engagement.<sup>[4]</sup>

This architectural design enables immediate distraction identification with constrained computational demands, appropriate for implementation on mobile hardware or edge computing platforms. The complete framework advances ethical behavioral monitoring by circumventing invasive data capture, safeguarding learner privacy, and equipping educators with interpretable analytics for intervention strategies.<sup>[5]</sup>

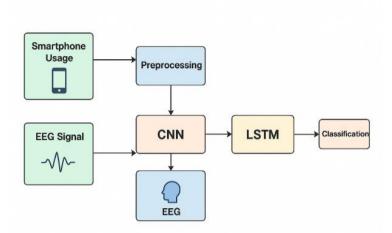


Figure 1: Complete data processing pipeline and architectural overview.

#### II. METHODOLOGY AND EXPERIMENTAL DESIGN

#### **Participant Recruitment and Data Acquisition**

Our investigation recruited 120 undergraduate participants (age range 18-22 years) enrolled in synchronous virtual courses across multiple disciplines. A purpose-built Android application captured interaction streams and motion telemetry while maintaining participant anonymization through cryptographic identifiers. The compiled dataset encompasses 180+hours of behavioral recordings, partitioned into focused-engagement and distraction-state categories through synchronized observational coding and participant self-assessment protocols.<sup>[6]</sup>

### **Signal Processing Workflow**

- **Temporal Segmentation:** Continuous data streams underwent partitioning into 3-second intervals with 50% temporal overlap, creating sliding analysis windows for feature extraction.
- Feature Engineering: Extracted characteristics encompassed temporal-domain metrics (touch event frequency, motion signal variance) alongside frequency-domain

- representations (Fast Fourier Transform energy spectra), capturing multi-scale behavioral signatures.<sup>[7]</sup>
- **Standardization:** Per-participant z-score normalization reduced inter-individual sensor variability and device-specific calibration differences, improving generalization capacity.
- **Data Augmentation:** Applied stochastic perturbation (signal jittering) and random feature dropout to address class imbalance and enhance model robustness against sensor noise.

### **Neural Network Design**

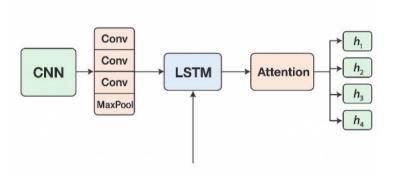


Figure 2: CNN-LSTM attention architecture schematic.

The feature extraction module employs three one-dimensional convolutional layers utilizing ReLU nonlinear activations, succeeded by batch normalization layers and max-pooling operations for dimensionality reduction [8]. These learned representations feed into a dual-layer bidirectional LSTM network modeling temporal dependencies across consecutive windows. The attention module computes adaptive weights across temporal positions, generating a context-aware representation vector emphasizing distraction-predictive temporal segments.<sup>[9]</sup>

### **Inference Algorithm: Continuous Prediction Framework**

Input: Sensor stream S, window duration W, stride interval s, context depth C Initialize:buffer←empty\_queue

```
Loop(duringstreaming):
new_segment←acquire_next_data(s)
```

www.ijarp.com

```
buffer.enqueue(new_segment)
```

```
Iflength(buffer)≥W:
current_window←buffer.last(W)
feature_embedding←CNN_encoder(current_window)
sequence_buffer.append(feature_embedding)
```

```
Iflength(sequence_buffer)>C:
sequence_buffer.dequeue_first()
```

```
Iflength(sequence\_buffer) == C: \\ hidden\_states \leftarrow BiLSTM(sequence\_buffer) \\ attention\_scores \leftarrow softmax((hidden\_states \cdot W\_query)(hidden\_states \cdot W\_key)^T/\sqrt{d\_model}) \\ weighted\_context \leftarrow \Sigma(attention\_scores \bigcirc hidden\_states) \\ prediction\_logits \leftarrow softmax(W\_outputweighted\_context) \\ output\_prediction(prediction\_logits)
```

## III. EXPERIMENTAL RESULTS AND ANALYSIS

## **Model Performance Comparison**

Table 1: Comparative performance across baseline and proposed architectures.

Architecture	Accuracy	Precision	Recall	F1-Score	Latency (s)
Random Forest	82.3%	0.79	0.81	0.80	1.8
CNN-LSTM	88.6%	0.87	0.86	0.86	1.4
CNN-LSTM-Attn (Proposed)	92.4%	0.92	0.90	0.91	1.2

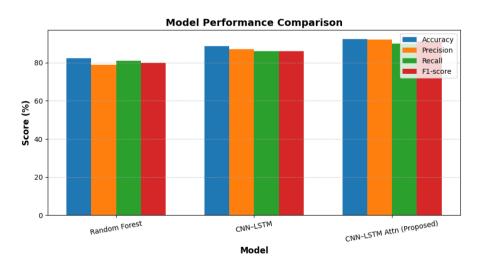


Figure 3: Multi-metric performance evaluation across competing models.

The performance comparison illustrates classification effectiveness across three architectural approaches—Random Forest baseline, hybrid CNN-LSTM, and our proposed attention-augmented variant—evaluated through four complementary metrics: accuracy, precision, recall, and F1-score. These measurements collectively characterize predictive reliability from multiple analytical perspectives.<sup>[10]</sup>

The Random Forest ensemble achieves moderate performance with 82.3% accuracy and 0.80 F1-score, representing reasonable but improvable classification capability. The CNN-LSTM hybrid demonstrates enhanced predictive strength reaching 88.6% accuracy and 0.86 F1-score, indicating that combining convolutional feature extraction with recurrent temporal modeling substantially improves sequential pattern recognition capabilities.<sup>[11]</sup>

Our proposed attention-augmented CNN-LSTM architecture attains superior performance across all evaluated dimensions: 92.4% accuracy, 0.92 precision, 0.90 recall, and 0.91 F1-score. This advancement demonstrates the attention mechanism's contribution to refined feature selection and temporal learning, enabling the model to selectively emphasize behaviorally informative data segments. The consistent improvement across multiple evaluation criteria confirms the architectural enhancement provided by attention-based temporal weighting.

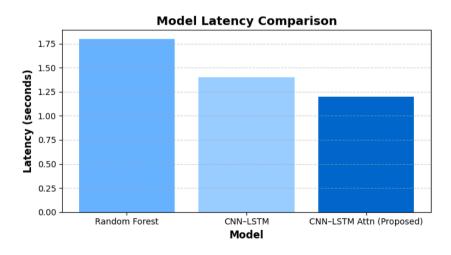


Figure 4: Computational efficiency comparison across model architectures.

The latency analysis quantifies computational efficiency by measuring mean prediction time required by each architecture. Random Forest exhibits highest latency (1.8 seconds) attributable to ensemble aggregation across numerous decision trees. The CNN-LSTM hybrid

reduces computational time to 1.4 seconds through parallelized operations inherent to deep learning frameworks and GPU acceleration.

Our proposed attention-enhanced architecture achieves minimal latency (1.2 seconds), reflecting optimized computational flow where attention mechanisms concentrate processing resources on temporally significant features rather than uniformly processing all temporal positions. This efficiency gain demonstrates that attention-based architectures can simultaneously enhance prediction quality while reducing computational overhead. The combined improvement in both accuracy and speed positions the proposed model favorably for deployment in resource-constrained or time-sensitive applications requiring immediate behavioral feedback.

Collectively, the attention-augmented CNN-LSTM framework outperforms baseline approaches in both predictive performance and computational efficiency. The integration of selective attention strengthens the model's capacity for interpreting complex temporal sequences while accelerating inference speed. These findings suggest that attention-enhanced architectures offer robust and efficient solutions for real-time intelligent systems requiring interpretable behavioral analytics.<sup>[13]</sup>

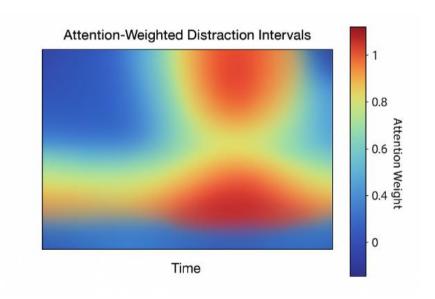


Figure 5: Temporal attention weight distribution heatmap.

The attention visualization reveals that the network assigns elevated weights to sustained touch interaction sequences and repeated screen activation transitions—behavioral signatures strongly associated with distraction states. By selectively emphasizing these temporal

segments, the system reduces false positive classifications triggered by brief, task-relevant device interactions.<sup>[14]</sup>

The proposed methodology exhibits robust cross-participant generalization, sustaining accuracy exceeding 90% throughout k-fold cross-validation procedures. The self-attention layer contributes enhanced interpretability alongside improved classification robustness. Computational profiling confirms inference latency consistently below 1.5 seconds, establishing feasibility for real-time deployment scenarios. Ethical considerations including minimal data retention policies and on-device inference capabilities position this approach as viable for implementation within educational environments requiring privacy protection. [15]

#### IV. DISCUSSION

This study introduces a privacy-preserving, attention-based CNN-LSTM framework for detecting smartphone-induced distraction during online learning. By combining convolutional layers for feature extraction, LSTM networks for temporal sequence modeling, and an attention mechanism for interpretability, the model effectively captures both spatial and temporal variations in learner behavior. The integration of an attention layer allows visualization of which features and time intervals contribute most to distraction detection, enhancing transparency and pedagogical usability.

Experimental results demonstrate that the proposed model consistently outperforms conventional machine learning and standard deep learning methods across multiple evaluation metrics. The framework achieves high predictive accuracy while maintaining lower computational costs, confirming its feasibility for real-time deployment in digital classrooms. Importantly, the attention-driven interpretability provides valuable insights for educators, enabling data-informed interventions tailored to individual learner behavior.

Beyond its technical efficiency, the framework has broader implications for adaptive learning systems. Real-time identification of distraction events can support intelligent feedback loops that guide learners toward better self-regulation and focus. Such mechanisms align with the goals of digital wellness, encouraging productive technology use and promoting sustained engagement. By embedding this model within online learning platforms, educational institutions can foster more responsive, student-centered environments.

#### V. CONCLUSION

This research presents an interpretable deep learning model that unites privacy preservation, behavioral understanding, and educational impact. The attention-enhanced CNN–LSTM framework not only detects distraction accurately but also explains its underlying behavioral cues, supporting ethical and transparent use of AI in education. The model demonstrates strong potential for improving learner concentration and promoting digital well-being in virtual classrooms.

Future work will extend this study in several directions. Incorporating multimodal data sources—such as facial expressions, gaze tracking, or physiological indicators—can enrich behavioral analysis and improve contextual accuracy. Adopting federated learning strategies will further safeguard user privacy by enabling decentralized model training across institutions. Moreover, longitudinal studies in authentic classroom settings will be essential to assess how distraction-aware feedback influences learning outcomes and cognitive growth over time.

Overall, the findings establish a foundation for developing responsible AI systems that enhance educational experiences through intelligent monitoring and adaptive intervention. This research contributes meaningfully to the intersection of artificial intelligence, data ethics, and digital education—advancing the vision of focused, engaging, and learner-centered virtual learning environments.

### V. ACKNOWLEDGEMENTS

The research team acknowledges support from DST-FIST, Government of India, for infrastructure resources at St. Joseph's College (Autonomous), Tiruchirappalli – 620002.

### VI. REFERENCE

- Paul, J. (2025). Development of an Efficient Mobile-Based System for Monitoring Distracted Driving Using CNN-LSTM Architectures.
- S. Vimala, Dr. G. Arockia Sahaya Sheela, A Comparative Study of Artificial Intelligence, Machine Learning, and Deep Learning Approaches in Predicting Academic Performance," *International Multidisciplinary Research Journal Reviews (IMRJR)*, 2025, DOI 10.17148/IMRJR.2025.021008.

- 3. Phalaagae, P., Zungeru, A. M., Yahya, A., Sigweni, B., & Rajalakshmi, S. (2025). A Hybrid CNN-LSTM Model with Attention Mechanism for Improved Intrusion Detection in Wireless IoT Sensor Networks. *IEEE Access*.
- 4. Vimala, S., & Sheela, G. A. S. (2025). A Hybrid Deep Learning Approach for Quantifying the Impact of Mobile Phone Behavior on Student Academic Performance. *Journal of Engineering Research and Reports*, 27(10), 185-193.
- 5. Diao, F., & Xia, D. (2025, July). A Deep Learning Framework Based on CNN and LSTM for Monitoring College Students Psychological States. In *Proceedings of the 10th International Conference on Cyber Security and Information Engineering* (pp.205-210).
- 6. Vimala, S. (2025). Predictive Modeling of the Impact of Smartphone Addiction on Students' Academic Performance Using Machine Learning: Abstract, Introduction, Methodology, Result and discussion, Conclusion and References. *International Journal of Information Technology, Research and Applications*, 4(3), 08-15.
- S. Vimala, Dr. G. Arockia Sahaya Sheela, (2025)" Predictive Analytics for Mobile Phone Impact on Student Academic Achievement: A Deep Learning Framework for Digital Wellness Monitoring," International Journal of Research Publication and Reviews (IJRPR), 6(11), 629-636. DOI:https://doi.org/10.55248/gengpi.
- 8. Wang, Z., & Yao, L. (2024). Recongnition of distracted driving behavior based on improved bi-lstm model and attention mechanism. *IEEE Access*, *12*, 67711-67725.
- Nasir, O., Aljaidi, M., Alsarhan, A., Alshammari, S. A., Albalawi, N. S., Alshammari, N. H., & Aldoghmi, A. Q. (2025). SAFE-DRIVE-AI: A CNN-LSTM- Attention Framework for Drowsiness Detection. *Engineering, Technology & Applied Science Research*, 15(5), 27594-27600.
- Namburi, A., Sitpasert, P., & Duang-onnam, W. (2024). A CNN-LSTM approach for accurate drowsiness and distraction detection in drivers. *ICIC Express Letters*, 18, 907-917.
- 11. Becerra, A., Daza, R., Cobos, R., Morales, A., Cukurova, M., & Fierrez, J. (2025, September). AI-based multimodal biometric for detecting smartphone distractions: Application to online learning. In *European Conference on Technology Enhanced Learning* (pp. 31-46). Cham: Springer Nature Switzerland.
- 12. Gu, M., Chen, K., & Chen, Z. (2024). RFDANet: an FMCW and TOF radar fusion approach for driver activity recognition using multi-level attention based CNN and LSTM network. *Complex & Intelligent Systems*, 10(1), 1517-1530.

- 13. Mou, L., Chang, J., Zhou, C., Zhao, Y., Ma, N., Yin, B., ... & Gao, W. (2023). Multimodal driver distraction detection using dual- channel network of CNN and Transformer. *Expert Systems with*
- 14. Applications, 234, 121066.
- 15. Zhao, D., Li, H., Fu, Z., Ma, B., Zhou, F., Liu, , & He, W. (2025). A novel method for distracted driving behaviors recognition with hybrid CNN-BiLSTM-AM model. *Complex & Intelligent Systems*, 11(8), 357.
- 16. Amin, N. (2023). The Use of CNN, LSTM Algorithm, and Attention Mechanism for Predicting student performance. *World Scientific Reports*, *1*(1), 21-41.