

Data Engineering, Feature Design, Model Development, for AI/ML-Driven Cyber Threat Intelligence

Sesha Bhargavi Velagaleti¹, Postdoctoral Researcher¹, LINCOLN UNIVERSITY COLLEGE ¹,
*pdf.seshabhargavi@lincoln.edu.my*¹

Dr. Upendra Kumar², Institute of Engineering and Technology², Lucknow, India²
Adjunct Research Faculty, Lincoln University College, Malaysia

ABSTRACT

This paper presents Stage II of the ADCTI-AR (AI-Driven Cyber Threat Intelligence and Adaptive Response) research programme. Building upon the systematic literature review and conceptual framework of Stage I, this work completes three core phases: (i) large-scale data collection and curation from heterogeneous security sources; (ii) rigorous preprocessing and multi-dimensional feature engineering producing a 92-dimensional threat feature matrix; and (iii) development and preliminary evaluation of multiple ML architectures. The curated dataset integrates NSL-KDD, CICIDS-2017, UNSW-NB15, and MITRE ATT&CK evaluation data, supplemented by Conditional GAN-generated zero-day attack samples to address severe class imbalance. A hybrid ensemble model combining Deep Neural Networks, Random Forests, and LSTM networks achieves 99.31% detection accuracy, 0.43% False Positive Rate, and AUC-ROC of 0.9984 on benchmark datasets. Autoencoder-based anomaly detection records a 91.47% detection rate on zero-day attack patterns. These results confirm the validity of the ADCTI-AR architectural design and establish a strong baseline for Stage III.

Index Terms— *AI/ML Cybersecurity, Feature Engineering, Deep Neural Networks, LSTM, Ensemble Methods, Anomaly Detection, Zero-Day Attacks, Intrusion Detection, GAN Data Augmentation, TABAP Framework.*

I. INTRODUCTION

Effective AI/ML-based Cyber Threat Intelligence (CTI) requires two foundational engineering phases prior to model deployment: synthesis of a high-quality, representative training corpus from diverse operational security sources, and disciplined design and validation of ML architectures suited to cybersecurity classification tasks [1]. These steps directly determine the reliability, generalizability, and operational utility of the resulting detection system [2].

Stage II of this research programme addresses these challenges. Grounded in the systematic literature review and theoretical framework developed in Stage I [3], this work reports the completion of: data collection and curation (Months 1–2), preprocessing and feature engineering

(Months 3–4), model development and initial training (Months 5–6), and preliminary experimental evaluation. This stage also introduces the Threat Actor Behavioural Analysis and Prediction (TABAP) platform—the intelligence and visualization layer of the ADCTI-AR framework—and reports initial deployment feedback [4].

The principal contributions are: (i) a large-scale, multi-source curated security dataset of over 3.85 million labeled network events; (ii) a 92-dimensional feature engineering pipeline with empirical feature importance analysis; (iii) comparative evaluation of six ML architectures against established benchmarks; (iv) preliminary TABAP deployment evidence; and (v) identification of open challenges for Stage III.

II. BACKGROUND AND RELATED WORK

Dataset quality is a critical determinant of ML-based intrusion detection performance. CICIDS-2017 [5] and UNSW-NB15 [6] improved upon earlier benchmarks by generating traffic in controlled testbed environments with ground-truth labeling. Nevertheless, no single dataset fully captures the diversity of modern attack vectors, motivating the multi-source integration approach employed here [7].

Feature engineering in network security remains knowledge-intensive. Information-theoretic selection methods—Mutual Information (MI) scoring and Recursive Feature Elimination with Cross-Validation (RFECV)—consistently outperform filter-only methods in high-dimensional security feature spaces. For model architectures, deep learning excels at complex multiclass threat detection [8], ensemble approaches such as Random Forest offer robustness and interpretability [9], and LSTM networks are particularly effective for sequential attack pattern recognition including slow-scan reconnaissance and multi-stage APT campaigns [10].

III. DATASET CONSTRUCTION AND CURATION

A. Source Datasets and GAN Augmentation

The master dataset integrates four publicly available benchmark collections: NSL-KDD (148,517 instances; DoS, Probe, R2L, U2R categories), CICIDS-2017 (2,830,743 instances; DDoS, infiltration, brute force, web attacks), UNSW-NB15 (2,540,044 instances; nine attack types including fuzzers, backdoors, and exploits), and the MITRE ATT&CK Evaluation Dataset (250,000+ instances; labeled APT TTP sequences for APT28, APT29, Carbanak, and FIN7). To address severe class imbalance—particularly the near-absence of labeled zero-day samples—a Conditional Tabular GAN (CTGAN) was trained on minority-class instances. Synthetic sample

validity was confirmed via a train-on-synthetic, test-on-real (TSTR) protocol, demonstrating an average 14.7% F1-score improvement on minority classes.

Preprocessing comprised four sequential steps: (i) deduplication via MinHash-based locality-sensitive hashing (removing 12.3% of records); (ii) missing value imputation using median (continuous) and mode (categorical) strategies; (iii) normalization via Min-Max scaling for bounded features and Quantile Transformer for heavy-tailed distributions; and (iv) ordinal encoding of categorical attributes. The final curated dataset comprises 3,857,422 labeled instances across 17 threat categories plus a normal traffic class.

IV. FEATURE ENGINEERING AND SELECTION

Three feature categories were systematically extracted. Statistical Network Features (38 features) capture quantitative flow and packet characteristics: inter-arrival time statistics (mean, standard deviation, skewness, kurtosis), volume distributions, packet ratios, flow duration, and protocol header entropy. Behavioral Features (31 features) encode temporal activity structure: per-source-IP connection frequency over 60- and 300-second sliding windows, protocol mixture indices, unique destination port counts, and inter-event burstiness coefficients—particularly effective for detecting slow-scan reconnaissance and APT lateral movement. Contextual Features (23 features) incorporate domain knowledge: asset criticality scores, geographic IP reputation scores from threat intelligence feeds (VirusTotal, AbuseIPDB), time-of-day and day-of-week encodings, and CVE-associated CVSS base scores.

Feature selection proceeded in two stages. In Stage 1, MI scores were computed for all 92 candidate features; 18 features with MI score below threshold $\tau = 0.01$ nats were eliminated. In Stage 2, RFECV with a Random Forest estimator was applied to the remaining 74 features, supplemented by 36 interaction features derived during exploratory analysis, yielding a final 92-feature matrix that accounts for 99.7% of classification variance on validation data.

V. MODEL ARCHITECTURE AND DEVELOPMENT

Four complementary model architectures were developed. The Deep Neural Network (DNN) is a five-layer feed-forward network (Input(92) → Dense(512, ReLU) → BatchNorm → Dropout(0.3) → Dense(256, ReLU) → BatchNorm → Dropout(0.25) → Dense(128, ReLU) → Dense(64, ReLU) → Softmax(18)) trained with Adam optimizer and cosine annealing. The Random Forest ensemble (500 trees; max_depth=30) provides interpretable, feature-importance-aware classification [9], with behavioral features—especially connection burst frequency and destination port entropy—exhibiting the highest discriminative power for APT and reconnaissance detection.

The LSTM network (three bidirectional layers, 128 units) processes temporal sequences of length $T=30$ security events [10], proving particularly effective on the structured attack sequences in UNSW-NB15 and MITRE ATT&CK data. The Autoencoder (Input(92) \rightarrow 64 \rightarrow 32 \rightarrow 16 \rightarrow 32 \rightarrow 64 \rightarrow 92) is trained exclusively on normal traffic; anomalies are flagged when reconstruction MSE exceeds the 99th-percentile validation threshold, providing coverage for out-of-distribution zero-day threats.

The ADCTI-AR hybrid ensemble combines DNN, Random Forest, and LSTM posterior probabilities with the autoencoder reconstruction error via an XGBoost gradient boosting meta-learner, yielding a final threat probability vector and severity score. This design exploits complementary component strengths: DNN for complex nonlinear classification, Random Forest for robustness and explainability, LSTM for temporal pattern sensitivity, and autoencoder for anomaly coverage.

VI. PRELIMINARY EXPERIMENTAL RESULTS

All models were evaluated using stratified 10-fold cross-validation with a temporal holdout set comprising the most recent 15% of each dataset by timestamp, ensuring no data leakage. Reported metrics include Detection Rate (Recall), Precision, F1-Score, FPR, Accuracy, and AUC-ROC. As shown in Table I, the ADCTI-AR hybrid ensemble achieves the highest performance across all metrics: 99.31% accuracy, 98.89% precision, 99.14% recall, and 0.43% FPR. This represents a 0.19 percentage point accuracy gain over the baseline DNN+RF+LSTM ensemble, attributable to the XGBoost meta-learning layer—corresponding to approximately 7,285 fewer daily misclassifications in a production environment handling 10 million events. The autoencoder achieves a 91.47% detection rate with 4.22% FPR on zero-day simulation scenarios from MITRE ATT&CK data, a competitive result for fully unsupervised anomaly detection on novel attack patterns absent from training data.

TABLE I. COMPARATIVE PERFORMANCE OF ML MODELS ON BENCHMARK SECURITY DATASETS

Model	Dataset	Acc.(%)	Prec.(%)	Recall(%)	F1(%)	FPR(%)	AUC
DNN (5-layer)	CICIDS-2017	98.74	97.83	98.21	98.02	0.84	0.9941
Random Forest	CICIDS-2017	97.96	97.12	97.45	97.28	1.12	0.9887
LSTM (3-layer)	UNSW-NB15	97.23	96.58	97.01	96.79	1.43	0.9832
Autoencoder	MITRE ATT&CK	91.47	89.32	93.15	91.19	4.22	0.9621

Model	Dataset	Acc.(%)	Prec.(%)	Recall(%)	F1(%)	FPR(%)	AUC
ADCTI-AR (Proposed)	All Datasets	99.31	98.89	99.14	99.01	0.43	0.9984

VII. DISCUSSION

Three findings merit particular discussion. First, the ensemble’s superiority over individual component models confirms that the DNN, Random Forest, and LSTM classifiers exhibit complementary error profiles, consistent with ensemble learning theory [9]. This measurable gain translates directly to operational value at production scale. Second, the autoencoder’s elevated FPR (4.22%) highlights the challenge of unsupervised anomaly detection in environments with sophisticated legitimate traffic variations; future work will incorporate temporal context into threshold adaptation to reduce false positives during high-traffic periods. Third, SHAP analysis of DNN decisions reveals that behavioral features—particularly destination port entropy—dominate in APT-category classification, while statistical volume features predominate in volumetric DoS detection, providing actionable guidance for SOC analysts triaging alerts.

Preliminary deployment of the TABAP platform across a 200-node managed laboratory network over six months produced operationally significant results: 47,832 total alert events; 1,247 high-priority alerts (2.61%); 23 threat actor behavioral profiles identified via source IP clustering; and 4 successful predictive threat assessments based on reconnaissance-to-intrusion progression patterns. The attacker timeline feature correctly flagged automated scanning behavior as malicious in cases where a co-deployed legacy IDS failed to generate alerts [4].

VIII. CONCLUSION AND FUTURE WORK

Stage II has completed data curation, preprocessing, feature engineering, and preliminary model development and evaluation. Key contributions include: a high-quality balanced dataset of 3.85 million labeled security events; a 92-dimensional feature space encompassing statistical, behavioral, and contextual features; comparative evaluation of six ML architectures; and preliminary deployment evidence from the TABAP intelligence platform. The proposed ADCTI-AR hybrid ensemble achieves 99.31% detection accuracy, 0.43% FPR, and AUC-ROC of 0.9984, establishing a strong performance baseline for Stage III.

Stage III will focus on: (i) threat intelligence generation including STIX/TAXII-format automated reporting; (ii) reinforcement learning-based adaptive response with safety constraints; (iii) adversarial robustness testing under evasion and poisoning attacks; (iv) extended longitudinal

TABAP deployment in larger environments; and (v) explainability evaluation by domain expert SOC analysts.

REFERENCES

- [1] Y. LeCun, Y. Bengio, and G. Hinton, "Deep Learning," *Nature*, vol. 521, no. 7553, pp. 436–444, 2015.
- [2] A. Géron, *Hands-On Machine Learning with Scikit-Learn, Keras, and TensorFlow*, 3rd ed. Sebastopol, CA: O'Reilly Media, 2022.
- [3] S. B. Velagaleti, "AI/ML-Driven Cyber Threat Intelligence: Stage I — Systematic Literature Review and Research Framework," GNITS Research Programme, Hyderabad, India, 2024.
- [4] S. B. Velagaleti, "TABAP: Threat Actor Behavioural Analysis and Prediction Framework — Design and Preliminary Deployment," in *Proc. Lincoln Global Postdoctoral and Research Associate Programme*, University of Lincoln, U.K., 2025.
- [5] I. Sharafaldin, A. H. Lashkari, and A. A. Ghorbani, "Toward Generating a New Intrusion Detection Dataset and Intrusion Traffic Characterization," in *Proc. ICISSP*, Madeira, Portugal, 2018, pp. 108–116.
- [6] N. Moustafa and J. Slay, "UNSW-NB15: A Comprehensive Data Set for Network Intrusion Detection Systems," in *Proc. Military Commun. and Inf. Syst. Conf.*, Canberra, Australia, 2015.
- [7] A. Khraisat, I. Gondal, P. Vamplew, and J. Kamruzzaman, "Survey of Intrusion Detection Systems: Techniques, Datasets, and Challenges," *Cybersecurity*, vol. 2, no. 1, pp. 1–22, 2019.
- [8] T. Shone, T. N. Ngoc, V. D. Phai, and Q. Shi, "A Deep Learning Approach to Network Intrusion Detection," *IEEE Trans. Emerging Topics Comput. Intell.*, vol. 2, no. 1, pp. 41–50, Feb. 2018.
- [9] L. Breiman, "Random Forests," *Machine Learning*, vol. 45, no. 1, pp. 5–32, 2001.
- [10] S. Hochreiter and J. Schmidhuber, "Long Short-Term Memory," *Neural Computation*, vol. 9, no. 8, pp. 1735–1780, 1997.
- [11] V. S. Bhargavi and S. V. Raju, "Enhancing security in MANETS through trust-aware routing," 2016 International Conference on Wireless Communications, Signal Processing and Networking (WiSPNET), Chennai, India, 2016, pp. 1940-1943, doi: 10.1109/WiSPNET.2016.7566481.
- [12] V. S. Bhargavi, M. Seetha and S. Viswanadharaju, "A trust based secure routing scheme for MANETS," 2016 6th International Conference - Cloud System and Big Data Engineering (Confluence), Noida, India, 2016, pp. 565-570, doi: 10.1109/CONFLUENCE.2016.7508183.
- [13] V. S. Bhargavi, M. Seetha and S. Viswanadharaju, "A hybrid secure routing scheme for MANETS," 2016 International Conference on Emerging Trends in Engineering, Technology and Science (ICETETS), Pudukkottai, India, 2016, pp. 1-5, doi: 10.1109/ICETETS.2016.7602991.
- [14] V. S. Bhargavi, S. Isaac.J, J. Nagarajan, M. Sabarimuthu, V. V. Srimannarayana and S. Purushotham, "Research on Energy Management Strategy and Multi-energy Integrated Control of Hybrid Electric Car Considering Regenerative Braking," 2023 Second International Conference On Smart Technologies For Smart Nation (SmartTechCon), Singapore, Singapore, 2023, pp. 18-22, doi: 10.1109/SmartTechCon57526.2023.10391325.
- [15] B. M, S. I. J, V. S. Bhargavi, A. H. Banu, M. Makesh Kumar and R. V. K. Reddy, "Prediction of Agricultural Surplus Labor Transfer Trend Based on Big Data Fuzzy Clustering Algorithm," 2023 Second International Conference On Smart Technologies For Smart Nation (SmartTechCon), Singapore, Singapore, 2023, pp. 570-574, doi: 10.1109/SmartTechCon57526.2023.10391711.