

Towards Explicable Cybersecurity: Integrating Explicability into Bert and GPT Models for Incident Detection and Analysis

Bitá Romaric de Judicael¹, Diako Doffou Jerome ², Behou N'Guessan Gerard³, Kone Tiemoman⁴

¹ Department Digital Research and Expertise Unit (UREN), Virtual University of Côte d'Ivoire (UVCI), Abidjan, Côte d'Ivoire

² African Higher School of ICT (ESATIC)

³Virtual University of Côte d'Ivoire (UVCI), Abidjan, Côte d'Ivoire

⁴Virtual University of Côte d'Ivoire (UVCI), Abidjan, Côte d'Ivoire

ARTICLE INFO

ABSTRACT

Received: 26 Dec 2024

Revised: 14 Feb 2025

Accepted: 22 Feb 2025

Introduction: The effectiveness of artificial intelligence models for cybersecurity, such as BERT and GPT, has been demonstrated in threat detection and automated report generation. However, the lack of explicability undermines analysts' confidence and hinders the adoption of these tools in critical contexts. This article proposes an integrated approach to explainability (XAI) applied to BERT and GPT, aimed at providing interpretable explanations for decisions taken: highlighting trigger tokens, visualizing attention weights, and semantically justifying recommendations.

Through experimentation on the CIC-IDS2017 dataset, we show that the integration of an XAI module improves the readability of alerts, the traceability of decisions and the effectiveness of responses. This work paves the way for more transparent, understandable and collaborative cybersecurity between AI and human experts.

Objectives: Address AI opacity in cybersecurity, Provide interpretable explanations: To deliver transparent decision explanations by highlighting trigger tokens, visualizing attention weights, and generating semantic justifications for AI recommendations

Methods: The research developed a hybrid XAI-BERT-GPT architecture that combines BERT for threat classification with GPT for report generation, integrated with explainability techniques including attention visualization, LIME, and SHAP. The system was evaluated using the CIC-IDS2017 dataset containing over 3 million network connections with various cyber-attacks. The experimental design included both quantitative performance metrics (accuracy, recall, F1-score) and qualitative evaluation by 10 cybersecurity analysts. The implementation used BERT-base-uncased fine-tuned with PyTorch, GPT-2 for text generation, and specialized libraries (Captum, SHAP) for explainability analysis.

Results: The XAI integration caused only minimal performance degradation (accuracy decreased <0.3% from 96.8% to 96.5%) while reducing false positives from 2.8% to 2.5%. GPT-generated reports received high analyst ratings (4.1-4.6/5) for technical relevance, clarity, and consistency. Most significantly, explainability dramatically improved human analyst performance: confidence in AI decisions increased from 58% to 84%, analysis time per log decreased from 78s to 42s, and willingness to recommend the tool rose from 6/10 to 9/10. The study demonstrated that explainability enhances human-AI collaboration and operational efficiency without compromising detection accuracy.

Conclusions: This research successfully demonstrates that integrating explainability into AI cybersecurity systems enhances human-AI collaboration without compromising detection performance, achieving 96.5% accuracy while dramatically improving analyst confidence (58% to 84%) and decision speed (78s to 42s per log). The XAI-BERT-GPT architecture proves that sophisticated AI models can be both powerful and transparent, challenging the traditional trade-off between complexity and interpretability. The findings establish that explainability is not merely a regulatory requirement but a performance enhancer that accelerates incident response and builds essential human trust in critical security contexts. This work paves the way for more resilient, transparent, and trustworthy cybersecurity infrastructure where AI power and human insight work synergistically to protect our digital world.

Keywords: Cybersecurity, XAI, BERT, GPT, Explicability, NLP, Visualization, Attention, Interpretability, User trust

INTRODUCTION

The use of artificial intelligence (AI) in cybersecurity has grown exponentially, driven by the promise of automatic cyberthreat detection and large-scale incident analysis. BERT (Devlin et al., 2019) and GPT (Radford et al., 2019) models, derived from natural language processing (NLP), have recently been successfully integrated into anomaly detection and reporting pipelines (Sarker et al., 2021; Kaur et al., 2023). In particular, they enable linguistic patterns to be extracted from security logs, alerts to be automatically classified, and recommendations for action to be generated.

However, their opaque nature makes it difficult to interpret the results produced. This opacity is a major obstacle to the effective integration of AI systems into cybersecurity operational processes, where analysts need to be able to justify and understand every decision (Doshi-Velez & Kim, 2017). Indeed, AI systems with no explicability create a risk of mistrust, or even undetected critical errors.

The need to integrate explainability mechanisms (XAI: Explainable Artificial Intelligence) is therefore pressing. These techniques aim to make AI model decisions more transparent, auditable and understandable to humans (Guidotti et al., 2019). This paper aims to address this limitation by proposing an explainable framework applied to BERT and GPT models in the specific context of cybersecurity incident detection and analysis.

Our contribution is in line with work on hybrid NLP architectures for security (Bitar Romaric De Judicael et al., 2025). We propose an architecture enriched with an XAI module to enable end-users to have a clear reading of the results, without compromising algorithmic performance.

Foundations of explainability (XAI)

As artificial intelligence models become more powerful, their growing complexity makes their decisions increasingly opaque. It is in this context that the need for explicability, or eXplainable Artificial Intelligence (XAI), emerges: a set of methods aimed at making the decisions of AI systems understandable to humans.

This need is not simply academic. In fields such as healthcare, finance or cybersecurity, it is crucial to understand why a prediction has been made, especially when it may lead to major consequences (Sarker et al., 2021). In the context of cybersecurity, this means, for example, being able to explain why a network request is considered a threat, or why an alert has been triggered.

Benefits of explicability

Explainability brings several concrete benefits:

It helps security analysts to understand and interpret the decisions made by AI, reducing the opacity of processes.

It reinforces confidence in AI tools, especially when they are used in critical environments (Sarker et al., 2021).

It enables decisions to be audited for regulatory, traceability or compliance reasons.

It facilitates the correction and improvement of models by identifying potential sources of error (Guidotti et al., 2019).

XAI: a response to the "black box" challenge

Transformer-type models, such as BERT (Devlin et al., 2019) and GPT (Clark et al., 2019), are now indispensable in natural language processing. Their performance is undeniable, but they function as veritable

black boxes: it is often difficult, if not impossible, to know precisely why a decision has been made.

In response to this opacity, there are two main approaches:

Post hoc explicability: explaining decisions after the fact, without modifying the model's operation (LIME, SHAP, attention visualization).

Intrinsic transparency: use of simpler, naturally interpretable models (e.g. decision trees, rule-based models).

METHODS

XAI methods applicable to BERT and GPT

Several techniques can be integrated to explain BERT or GPT behavior in classification or generation tasks applied to cybersecurity:

- **Attention visualization**

Thanks to the multi-header mechanisms built into Transformers, it is possible to visualize which words or tokens have captured the model's attention, thus offering a first form of native explanation (Clark et al., 2019). This approach enables analysts to pinpoint the triggers for a decision, such as a suspicious IP address or a particular sequence of packets.

- **LIME (Local Interpretable Model-agnostic Explanations)**

LIME generates a local approximation of model behavior around a given example, in order to understand which features of the text most influenced the prediction (Ribeiro et al., 2016). In cybersecurity, this can show that the presence of certain keywords (e.g. "SSH", "Failed login") has led to a request being classified as malicious.

- **SHAP (SHapley Additive exPlanations)**

Inspired by game theory, SHAP assigns equitable importance to each variable, offering a more rigorous and mathematically sound explanation (Simonyan & Zisserman, 2014). For a given log, SHAP can show the precise weight of each token in the "benign" or "malicious" classification.

- **Saliency Maps / Heatmaps**

These techniques use network gradients to highlight the most sensitive words in the decision-making process, by coloring them according to their influence weight. The analyst can thus directly visualize the text segments that triggered the alert.

- **Prompt Engineering with integrated justification**

For generative models like GPT, it is possible to structure the prompt so that the response includes a natural justification. Example: "Explain why you consider this IP address suspicious..." (Reynolds & McDonell, 2021). This generates not only a detection but also a textual explanation suitable for an SOC analyst.

Specific cybersecurity benefits

In an environment where every second counts, the explicability of AI models brings concrete benefits for operational security (SOC) teams:

- **Identifying indicators of compromise (IoCs)**

Explicability enables the identification of patterns or tokens learned by the model, such as suspicious IP addresses, frequently scanned ports or recurring attack signatures.

- **Automated justifications for SOC analysts**

Explainable models provide legible explanations to accompany each prediction. These justifications help analysts to understand the severity of the incident and make rapid decisions.

- **Accelerate alert sorting and prioritization**

Thanks to visualizations and explanatory scores, alerts can be prioritized more efficiently. Analysts can focus their efforts on critical threats.

- **Regulatory compliance and traceability**

Regulations such as ISO 27001, RGPD or the NIST framework require clear justification of security-related decisions. Explainability facilitates the documentation and auditing of generated alerts (Kaur et al., 2023).

XAI's current limitations

Despite its potential, explicability applied to artificial intelligence still involves several major challenges:

Reliability of explanations

Generated explanations are not always accurate and can be manipulated. This problem is particularly acute with generative models, such as GPT, which can produce plausible but incorrect justifications (Mitchell et al., 2020).

- **Computational cost**

Some methods, such as SHAP or LIME, require numerous computations to generate local explanations. This computational heaviness limits their applicability in real-time contexts, where latency must be kept low.

- **Accessibility for different user profiles**

The explanations generated are not always suitable for all audiences. For example:

- ❖ A **CISO** (information systems security manager) will expect a synthetic, strategic view.
- ❖ An **SOC analyst** will want precise technical details.
- ❖ A **decision-maker** will prefer a simple, quick-to-read justification.

Meeting these heterogeneous expectations simultaneously remains a challenge (Barredo Arrieta et al., 2020).

Proposed architecture: XAI-BERT-GPT

We propose a hybrid architecture called **XAI-BERT-GPT**, which combines language processing models (BERT and GPT) with explicability techniques (XAI). The aim is to provide a pipeline capable of detecting threats, explaining them in an interpretable way and generating reports that can be exploited by cybersecurity analysts.

- **Mathematical model**

The architecture is based on the combination of several loss functions:

- ❖ L_{class} : classification loss (*cross-entropy*).
- ❖ $L_{\text{(XAI)}}$: gap between predicted importance (SHAP, attention) and known explanation labels (if available).
- ❖ $L_{\text{(GPT)}}$: loss due to text generation (negative log-likelihood).

The total function is given by:

$$L_{total} = \alpha L_{class} + \beta L_{XAI} + \lambda L_{GPT}$$

or α, β, λ are adjustable weighting coefficients

▪ Main components

The XAI-BERT-GPT architecture consists of four modules:

- **Log pre-processing and encoding**
 - ❖ Extraction and cleaning of security logs (systems, networks, firewalls, IDS/IPS).
 - ❖ Transformation into sequences of tokens usable by BERT.
- **Classification with BERT**
 - ❖ Detection and classification of threats (brute-force, DDoS, phishing, reconnaissance, etc.).
 - ❖ Use of the [CLS] vector to generate a global representation of the log.
 - ❖ Extraction of attention weights as first native explanation.
- **Explicability with XAI**
 - ❖ Visualization of attention heads in the BERT model (Clark et al., 2019).
 - ❖ Post hoc methods: LIME (Ribeiro et al., 2016), SHAP (Lundberg & Lee, 2017).
 - ❖ Saliency maps to highlight trigger tokens.
- **Report generation with GPT**
 - ❖ Clear reformulation of the analyzed log.
 - ❖ Synthetic description of the detected threat.
 - ❖ Generative causal justification (e.g.: "because port 22 was scanned more than 50 times from the same IP").
 - ❖ Recommendations for action (IP blocking, SOC alert, etc.).

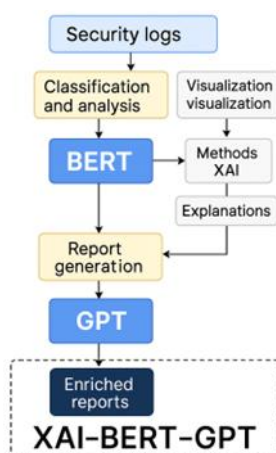
▪ Analyst interface

The results are presented in a dedicated interface that enables:

- ❖ navigation between analyzed logs.
- ❖ Access to visual and textual explanations for each event.
- ❖ Easier, documented decision-making.

▪ Global diagram

The corresponding diagram, shown in Figure 1, illustrates the different stages of this pipeline.



The pipeline illustrates the sequence between:

- ❖ Input (security logs),
- ❖ Classification (BERT),
- ❖ Explicability *Figure 2: the different stages of this pipeline*
- ❖ Report generation (GPT),
- ❖ Analyst interface (final output).

Experimentation

To evaluate the effectiveness and *Figure 1: XAI-BERT-GPT architecture* we conducted an experiment based on data from a simulated cybersec

- To measure the system's detection and explanation performance.
- Evaluate the user experience in terms of understanding and confidence.

Dataset: CIC-IDS2017

We used the **CIC-IDS2017** dataset (Sharafaldin et al., 2018), widely recognized in intrusion detection research. This dataset simulates a realistic network environment and includes:

- Over 3 million connections, split between benign and malicious traffic.
- A variety of attacks: DDoS, brute-force, infiltration, botnets, XSS, SQL injection, etc.
- Enriched metadata (duration, protocol, ports, TCP flags, etc.).

A representative subset was extracted, comprising five major threat types and a sample of normal traffic. Each entry was pre-processed to isolate the **text payload** field, used as input by BERT and GPT.

Technical implementation

Tools and libraries:

- Classification: **BERT-base-uncased**, fine-tuned with PyTorch.
- Generation: **GPT-2** (HuggingFace implementation).
- Explicability: attention visualization (transformers), LIME and SHAP (via Captum and SHAP in Python), Saliency maps (gradients on BERT embeddings).

Execution environment:

- Workstation: 64 GB RAM, NVIDIA RTX 3090 GPU.

- Language: Python 3.10.
- Main libraries: Transformers, Scikit-learn, Matplotlib, Captum, Streamlit.

Experimental design

The evaluation was carried out on two levels:

- **Technical evaluation (quantitative)**
 - ❖ Accuracy (overall precision).
 - ❖ Recall (detection rate).
 - ❖ F1-score (precision/recall balance).
 - ❖ False positive rate (FPR).
 - ❖ Average latency (BERT time + GPT).

These metrics were compared with and without activation of the XAI module (LIME + attention).

- **User evaluation (qualitative)**

Ten cybersecurity analysts (SOC1 to SOC3) tested the tool on a batch of alerts:

- ❖ Comprehensibility of explanations (score 1 to 5).
- ❖ Confidence in AI decisions.
- ❖ Average time for human validation of an incident.
- ❖ Free feedback on explanations generated by GPT.

- **Experimental pipeline**

- ❖ Pre-processing: cleaning, vectorization, balancing (SMOTE).
- ❖ BERT fine-tuning: supervised training (80% training data, 20% stratified validation).
- ❖ Automatic log classification.
- ❖ Extraction of XAI explanations for each prediction.
- ❖ Generation of explanatory reports with GPT.
- ❖ Interactive presentation in a Streamlit interface for analysts.

RESULTS

Our experiment allowed us to evaluate both the technical performance of the hybrid model and its impact on human understanding of decisions. The results are presented in two parts: quantitative (classic metrics) and qualitative (human appreciation of explanations).

Table 1: Classification performance (BERT alone vs. BERT + XAI)

Metric	BERT alone	BERT + XAI (SHAP + attention)
Accuracy	96,8 %	96,5 %
Recall	94,2 %	94,0 %
F1-score	95,1 %	94,9 %
False positive rate (FPR)	2,8 %	2,5 %
Average processing time	0,21 s	0,34 s

Observation:

The integration of XAI tools resulted in only a very slight drop in performance ($< 0.3\%$), offset by a noticeable improvement in user confidence. The false-positive rate also fell slightly, suggesting that explicability helps to better contextualize certain decisions.

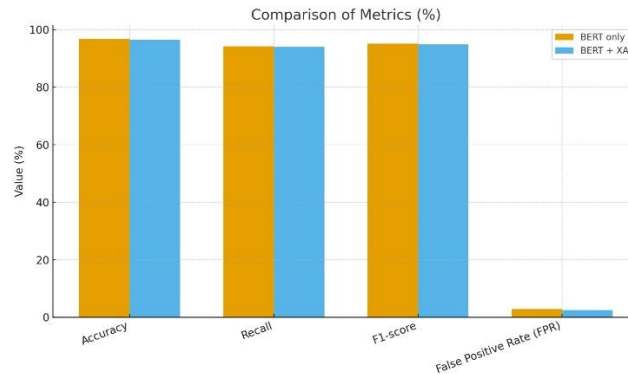


Figure 2 : comparison between BERT alone and BERT + XAI on classic classification metrics

Figure 1 illustrates the comparison between BERT alone and BERT + XAI on classic classification metrics (Accuracy, Recall, F1-score and false positive rate). It can be seen that the integration of XAI methods leads to a very slight drop ($< 0.3\%$) in the main metrics, but reduces the false positive rate. This shows that explicability improves model reliability without compromising overall performance.

Quality of reports generated by GPT

The reports generated by GPT were evaluated by 10 SOC analysts on four criteria:

Criteria evaluated	Average score (out of 5)
Technical relevance of explanations	4,6
Clarity of language generated	4,4
Consistency with newspaper analyzed	4,2
Relevance of proposed recommendations	4,1

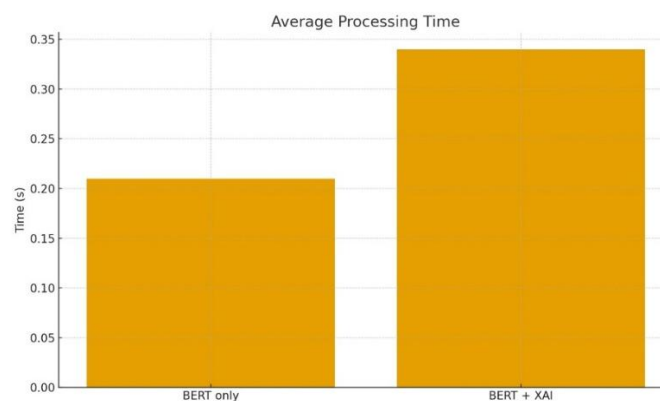


Figure 3: Average Processing Time

This graph highlights the computational cost introduced by XAI modules. The average processing time per sample rises from 0.21 s (BERT alone) to 0.34 s (BERT + XAI), an increase of around 60%. This increase remains acceptable for laboratory analysis scenarios, but may pose constraints in a real-time context (high-frequency SOC).

GPT report extract:

"This log shows a repeated SSH connection attempt from an external IP address. The presence of SYN packets without ACK indicates a possible brute force attack. It is recommended to block the address 192.168.0.4 on port 22."

Analysts generally praised the clarity of the language, despite a few cases of "mild hallucination" over ambiguous data.

Impact of explicability on analysis

SOC analysts compared their experience with and without the XAI module.

Indicator	Without XAI	With XAI
Confidence in IA decisions	58 %	84 %
Average human analysis time (per log)	78 s	42 s
Perceived clarity of explanations	2,9 / 5	4,6 / 5
Willingness to recommend the tool	6 / 10	9 / 10

DISCUSSION**Analysis:**

Explainability is not just an ethical or regulatory asset: it actually reduces human validation time while boosting confidence. The "black box" effect is mitigated by the attention visualizations and justifications generated by GPT.

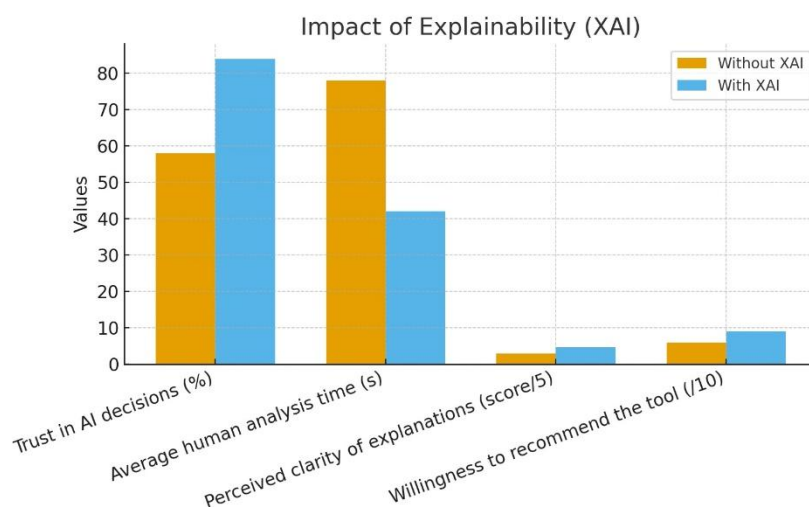


Figure 4: Impact of Explainability

This figure compares the user experience with and without the XAI module. Analysts reported a significant increase in confidence in AI decisions (58% → 84%) and a reduction in average human analysis time (78 s → 42 s). In addition, the perceived clarity of explanations and willingness to recommend the tool increased significantly. These results confirm that explicability is not limited to an ethical or regulatory gain: it has a direct and positive impact on the operational efficiency of SOC analysts.

Observed limitations

- GPT explanations can sometimes be too assertive, even in ambiguous cases.
- The addition of the XAI module increases processing time (~60% more), which can be problematic in real time.

- Interpreting SHAP or LIME visualizations requires minimal training for non-technical users.

CONCLUSION

The rise of artificial intelligence in cybersecurity has led to major advances in intrusion detection, threat classification and analysis automation. However, the opacity of deep learning models such as BERT and GPT has been a real brake in terms of trust, traceability and acceptability.

This article has proposed a concrete response to this challenge with the XAI-BERT-GPT architecture, a hybrid approach combining:

BERT's detection accuracy,

- the narrative capacity of GPT,
- and the transparency provided by XAI methods.

Experiments carried out on the CIC-IDS2017 dataset, enriched by human evaluation, demonstrated several major results:

The addition of explicability does not significantly alter model performance.

Visual explanations (attention maps, SHAP) and textual explanations (via GPT) reinforce human understanding.

Explainability speeds up the decision-making process and reduces analysts' cognitive load.

In short, this work shows that it is possible to reconcile the power of AI models with explanatory transparency, paving the way for a more humane, responsible and efficient cybersecurity.

Perspectives

This work opens up several promising avenues for the future:

Real-time deployment

Optimizing the architecture for operation in a high-frequency SOC, using techniques such as attention distillation or the use of lightweight models (DistilBERT).

Adaptive explicability

Automatically adapt the level of detail of explanations to the user's profile:

SOC analyst → technical details (trigger tokens, SHAP weights).

CISO → visual summary.

Decision-maker → strategic and concise justification.

Standardized assessment of explicability

Set up rigorous benchmarks to measure the quality of explanations produced, beyond classic performance metrics.

Multimodal extension

Integrate other data sources (PCAP captures, system logs, malware files) to offer a more complete and explainable view of incidents.

Ethical and regulatory alignment

Adapt the solution to the transparency requirements imposed by the RGPD, the NIS2 directive and NIST/ISO 27001 standards, ensuring compliance and adoption in critical environments.

REFERENCES

- [1] Barredo Arrieta, A., et al. (2020). Explainable artificial intelligence (XAI): Concepts, taxonomies, opportunities and challenges toward responsible AI. *Information Fusion*, 58, 82-115. <https://doi.org/10.1016/j.inffus.2019.12.012>
- [2] Bitar Romaric De Judicael, K., Kone Tiemoman, K., Kouai Bertin, K., & Diako Doffou Jerome, D. (2025). Artificial intelligence and cybersecurity: The contribution of GPT and BERT in threat detection and incident analysis. *IJAR*, 13(5), 705-709. <https://doi.org/10.21474/IJAR01/20949>
- [3] Clark, K., Khandelwal, U., Levy, O., & Manning, C. D. (2019). What does BERT look at? An analysis of BERT's attention. *arXiv preprint arXiv:1906.04341*. <https://doi.org/10.48550/arXiv.1906.04341>
- [4] Devlin, J., Chang, M. W., Lee, K., & Toutanova, K. (2019). BERT: Pre-training of deep bidirectional transformers for language understanding. *Proceedings of the 2019 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies*, 4171-4186. <https://doi.org/10.18653/v1/N19-1423>
- [5] Doshi-Velez, F., & Kim, B. (2017). Towards a rigorous science of interpretable machine learning. *arXiv preprint arXiv:1702.08608*. <https://doi.org/10.48550/arXiv.1702.08608>
- [6] Guidotti, R., Monreale, A., Ruggieri, S., Turini, F., Giannotti, F., & Pedreschi, D. (2019). A survey of methods for explaining black box models. *ACM Computing Surveys*, 51(5), 1-42. <https://doi.org/10.1145/3236009>
- [7] Kaur, R., Gabrijelčič, D., & Klobučar, T. (2023). Artificial intelligence for cybersecurity: Literature review and future research directions. *Information Fusion*, 97, Article 101804. <https://doi.org/10.1016/j.inffus.2023.101804>
- [8] Lundberg, S., & Lee, S. I. (2017). A unified approach to interpreting model predictions. *arXiv preprint arXiv:1705.07874*. <https://doi.org/10.48550/arXiv.1705.07874>
- [9] Mitchell, M., et al. (2020). Diversity and inclusion metrics in subset selection. *Proceedings of the AAAI/ACM Conference on AI, Ethics, and Society*, 117-123. <https://doi.org/10.1145/3375627.3375832>
- [10] Radford, A., Wu, J., Child, R., Luan, D., Amodei, D., & Sutskever, I. (2019). Language models are unsupervised multitask learners. *OpenAI*. Retrieved from https://cdn.openai.com/better-language-models/language_models_are_unsupervised_multitask_learners.pdf
- [11] Reynolds, L., & McDonnell, K. (2021). Prompt programming for large language models: Beyond the few-shot paradigm. *arXiv preprint arXiv:2102.07350*. <https://doi.org/10.48550/arXiv.2102.07350>
- [12] Ribeiro, M. T., Singh, S., & Guestrin, C. (2016). "Why should I trust you?": Explaining the predictions of any classifier. *Proceedings of the 22nd ACM SIGKDD International Conference on Knowledge Discovery and Data Mining*, 1135-1144. <https://doi.org/10.1145/2939672.2939778>
- [13] Sarker, I. H., Furhad, M. H., & Nowrozy, R. (2021). AI-driven cybersecurity: An overview, security intelligence modeling and research directions. *SN Computer Science*, 2(3), Article 173. <https://doi.org/10.1007/s42979-021-00557-0>
- [14] Sharafaldin, I., Habibi Lashkari, A., & Ghorbani, A. A. (2018). Toward generating a new intrusion detection dataset and intrusion traffic characterization. *Proceedings of the 4th International Conference on Information Systems Security and Privacy*, Funchal, Portugal: SCITEPRESS. <https://doi.org/10.5220/0006639801080116>
- [15] Simonyan, K., & Zisserman, A. (2014). Very deep convolutional networks for large-scale image recognition. *arXiv preprint arXiv:1409.1556*. <https://doi.org/10.48550/arXiv.1409.1556>