

# AI-Driven Security Risk Mitigation: Enhancing Threat Assessment in Transit Infrastructure

Amirhossein Saali<sup>1,2\*</sup> and Raffaele Alfano<sup>2</sup>

<sup>1</sup>Faculty of Applied Science, The University of British Columbia, Vancouver, BC, Canada

<sup>2</sup>Hitachi Rail, Toronto, Canada

## \*Corresponding Author

Amirhossein Saali, Faculty of Applied Science, The University of British Columbia, Vancouver, BC, Canada and Hitachi Rail, Toronto, Canada.

**Submitted:** 2025, Nov 03; **Accepted:** 2025, Dec 05; **Published:** 2025, Dec 16

**Citation:** Saali, A., Alfano, R. (2025). AI-Driven Security Risk Mitigation: Enhancing Threat Assessment in Transit Infrastructure. *Adv Urban Region Dev Plann*, 2(1), 01-09.

## Abstract

The evolving threat landscape in transit infrastructure requires more adaptive, scalable, and precise security risk assessment and proposing proper risk mitigation measures. Traditional human-led approaches are often labor-intensive, subjective, and limited in their ability to incorporate real-time contextual data. This paper presents a novel AI-driven framework that combines geospatial analysis and large language models (LLMs) to automate the generation of structured, context-aware mitigation strategies aligned with industry standards and best practices. Our methodology involves a two-stage pipeline: (1) environmental feature extraction from Google Maps imagery and architectural design documentation, and (2) structured mitigation generation via controlled prompting of LLMs (GPT-3.5 and GPT-4.1) as two frequently used, and available models. We evaluate model performance across 320 real-world threat scenarios spanning 32 threat types and 37 transit assets, using a multi-criteria rubric validated by security experts. Our results determine that GPT-4.1 model consistently outperforms GPT-3.5 in contextual relevance, logical consistency, and adherence to classification schemes, even though at higher computational cost. The framework also demonstrates high throughput, with practical implications for both rapid network-wide assessments and in-depth expert analysis. This study highlights the capacity of hybrid computer vision (CV)-LLM architectures in advancing autonomous security planning, while identifying key limitations and pathways for future improvement.

**Keywords:** Critical Infrastructure Protection, Transit Security, Threat and Vulnerability Risk Assessment (TVRA), AI Based Mitigation Planning, LLM Model Application, Security Risk Assessment, Natural Language Processing (NLP), Automated Threat Modeling, Crime Prevention Through Environmental Design (CPTED)

## 1. Introduction

Security risk assessments play a crucial role in safeguarding critical infrastructure, including rail and transit systems. The identification and mitigation of threats are traditionally performed by security professionals using established frameworks such as the American Public Transportation Association (APTA) guidelines, the National Institute of Standards and Technology (NIST) Risk Management Framework, the Federal Transit Administration (FTA) Security and Emergency Management Program, and Crime Prevention Through Environmental Design (CPTED) principles [1-4]. These frameworks provide a structured approach to identifying vulnerabilities, assessing risks, and implementing

protective measures to enhance the resilience of transit networks. However, manual assessments can be subjective, resource-intensive, and prone to inconsistencies due to variations in human judgment and expertise. Additionally, traditional risk assessment methods may struggle to keep pace with the rapidly evolving threat landscape, including emerging challenges such as cyber-physical attacks and coordinated disruptions. As transit systems become increasingly interconnected and digitized, there is a growing need for automated, data-driven approaches that complement existing methodologies while improving efficiency and accuracy.

---

Recent advancements in artificial intelligence (AI) and natural language processing (NLP) offer transformative opportunities to enhance security risk assessment processes. Large language models (LLMs) such as Open AI's GPT-3.5-turbo, Google's BERT, and T5 have demonstrated significant capabilities in analyzing unstructured threat data, identifying patterns, and generating actionable mitigation strategies [5-7].

AI-driven systems can process vast datasets, including security reports, incident logs, and geospatial intelligence, to detect emerging risks and predict threat trends with high accuracy. Additionally, reinforcement learning (RL) frameworks can optimize risk prioritization by continuously learning from real-time security incidents and adapting mitigation recommendations accordingly [8]. AI-powered analysis can also improve transit security by integrating video surveillance, access control logs, and crowd movement data to detect anomalies and suggest security interventions. By leveraging these AI-driven insights, security professionals can supplement traditional assessments with automated, scalable, and adaptive recommendations, reducing response times while improving situational awareness in dynamic security environments.

This research investigates the feasibility of AI-driven threat mitigation in rail and transit infrastructure by developing an automated system that analyzes threats and recommends mitigation measures. Specifically, we explore whether the integration of large language models (LLMs) into the security risk assessment workflow can improve alignment with industry standards, enhance decision-making, and support resource optimization in mitigation planning.

Based on these objectives, we formulate the following research questions:

- Can large language models generate mitigation strategies that align with established frameworks such as CPTED and APTA?
- Does incorporating environmental context (e.g., lighting, surveillance, location type, neighbourhood) improve the specificity and relevance of the generated recommendations?
- How do different available LLM configurations models compare in terms of output quality and consistency?

## 2. AI in Security Risk Assessment

The application of artificial intelligence in security risk assessment has been an emerging topic in recent years. Various studies have explored the potential of AI-driven solutions in enhancing security planning and risk mitigation strategies.

Research on AI-driven security risk assessment has gained traction, particularly in the domains of cybersecurity and physical security. Studies such as those by and Johnson & have demonstrated that machine learning models can effectively analyze threat patterns and predict vulnerabilities in critical infrastructure [9,10]. Their conclusions emphasize that while AI can significantly enhance risk prediction, human oversight remains crucial to contextualize

automated assessments and address model biases.

Additionally, explored the use of deep learning models for automated anomaly detection in cybersecurity applications [11]. Their findings highlighted that AI-driven models were highly effective in identifying threats but required periodic retraining to adapt to evolving attack techniques.

More recently, examined the integration of reinforcement learning in predictive threat mitigation, highlighting the advantages of adaptive risk assessment models [12]. They concluded that AI models using reinforcement learning could dynamically adjust security postures, improving response efficiency but also introducing risks related to overfitting when exposed to limited data scenarios.

Another study by demonstrated the effectiveness of AI-driven geospatial analysis in enhancing situational awareness for physical security applications [13]. The study found that AI-assisted geospatial monitoring improved anomaly detection rates by over 30% compared to traditional methods but also raised concerns regarding privacy and the potential for false positives in complex environments.

Further, explored the role of hybrid AI models in integrating structured and unstructured security data for risk analysis, concluding that multi-modal AI frameworks offer greater contextual accuracy than single-model approaches. Similarly, demonstrated that federated learning can enhance privacy-preserving threat assessments while maintaining efficiency in large-scale transit networks [14,15].

### 2.1. AI and Threat Mitigation

Existing AI models have been employed to assist in generating mitigation strategies for various security concerns. For instance, explored the use of natural language processing (NLP) models to analyze security reports and propose risk mitigation actions. Their findings indicated that AI-generated recommendations often align with expert-driven solutions but require human oversight to ensure contextual accuracy and feasibility. They concluded that AI-enhanced NLP models significantly reduce assessment time but must be integrated with human decision-making frameworks to prevent misinterpretations.

### 2.2. Application of AI in Transit Infrastructure Security

Few studies have explicitly examined the impact of AI on transit infrastructure security. One notable work by proposed a predictive analytics framework for railway security, leveraging AI to detect potential threats based on surveillance data. While their approach demonstrated promise, it lacked a structured methodology for generating tailored mitigation measures [16].

Another study by proposed an AI-powered video surveillance system for transit hubs, which significantly enhanced real-time anomaly detection but required human intervention to verify AI-flagged incidents, highlighting the current limitations of

autonomous security models [17].

### 2.3. Limitations of AI in Security Decision-Making

Despite the potential of AI in security risk assessment, several limitations exist. Previous research, including studies by and, highlights concern regarding bias in AI-generated recommendations, the lack of contextual awareness, and challenges in integrating AI insights into existing security frameworks [18].

Furthermore, explored adversarial attacks on AI-based security models, demonstrating that threat actors can manipulate AI-driven assessments, raising concerns about model robustness and adversarial resilience [19].

Addressing these limitations is crucial for ensuring that AI-driven security assessments remain practical and reliable.

## 3. Methodology

This section describes the dataset, AI model selection, risk assessment criteria, mitigation strategy generation, and evaluation methods used in this study.

### 3.1. Task Formulation

The proposed system performs a two-stage pipeline:

- Environmental context extraction from geospatial data, and
- Structured mitigation generation using a large language model (LLM).

#### 3.1.1. Context Extraction

Given a location  $L_i$ , we retrieve satellite and Street View imagery from Google Maps and use a combination of LLM-assisted interpretation and lightweight computer vision to infer context features  $C_i$ , such as:

- Lighting (e.g., presence of poles, shadows)
- Surveillance (e.g., visible cameras, blind spots)
- Enclosure (open vs. confined)
- Access control (fencing, gates, turnstiles)
- Visibility and occupancy

The output is a structured context descriptor:

$$C_i = f_{env}(L_i) \quad (1)$$

where  $f_{env}$  represents the AI-assisted environment analysis function.

#### 3.1.2. Mitigation Generation

The main task is framed as a conditional structured text generation problem.

Given a threat description  $T_i$  and context  $C_i$ , the model generates a list of mitigation strategies  $M^i$  in a structured format:

$$M_i = LLM(T_i, C_i; \theta) \quad (2)$$

where:

- $T_i$ : natural language threat input
- $C_i$ : structured context from Stage 1
- $\theta$ : prompt template and model parameters (e.g., GPT-4.1, temperature, max tokens)
- $M_i$ : mitigation outputs formatted as:

$$M_i = \{(m_{i1}, p_{i1}, a_{i1}), \dots, (m_{ik}, p_{ik}, a_{ik})\} \quad (3)$$

with each tuple containing:

- $M_{ij}$ : mitigation action (text)
- $p_{ij}$ : CPTED principle (e.g., natural surveillance)
- $a_{ij}$ : APTA mitigation category (e.g., physical barriers, signage)

This setup goes beyond classification or retrieval. The LLM is used in a controlled generative setting, producing structured content based on real-world environmental inputs and security domain knowledge. This controlled prompting structure ensured that model outputs were aligned with expert expectations and compliant with mitigation classification frameworks.

### 3.2. Dataset and Scenario Design

We curated 320 threat–context pairs covering the full 32item threat taxonomy, evaluated across 37 assets (22 stations, 8 substations, 7 yards). The overall classification of threats are shown in Table 1.

Class Cluster	Example threats	security scenarios
Terrorism	Bodyworn IED, Heavy Bomb, Drone Attack	160
Criminal	Robbery, Assault, Sabotage	96
Nuisance / Disruption	Graffiti, Trespassing, Homelessness	64

**Table 1: Distribution of Threat Categories and Security Scenarios in the Evaluation Dataset**

*Environmental levels:* Every threat is paired with one of 10 environmental permutations spanning lighting, surveillance, enclosure, access control, and occupancy levels.

*Sample diversity analysis:* A twoway ANOVA (threat  $\times$  location) on contextfeature coverage found  $p = 0.48$ , indicating no significant

interaction. This fact is confirming even distribution across the design space.

#### 3.2.1. Dataset

The dataset used in this study consists of documented security threats in Rail and Transit projects. It includes the following key

attributes:

- **Threat Description:** A textual summary of the identified security threat and the scenario of occurrence.
- **Threat Category:** Classification based on predefined security risk categories (e.g., vandalism, unauthorized access, terrorism-related threats).
- **Risk Level:** Initial assessment of the risk (Very Low, Low, Moderate, High, Very High) based on the criteria available in standards such as APTA.
- **Location Information:** Detailed breakdown of where the threat is applicable.
- **Baseline Security Controls:** Existing mitigation measures already implemented to counter the identified threat.

The dataset was compiled from real-world risk assessment reports, industry standards, and expert consultations. Data preprocessing included standardizing terminology, removing duplicate entries, and ensuring consistency in categorization.

### 3.2.2. AI Model Selection

The AI models used in this study are OpenAI’s GPT-3.5-turbo and GPT-4.1, two LLMs known for their advanced text generation capabilities. The selection was based on the following criteria:

#### 3.2.2.1. GPT-3.5-turbo

- **Ability to Process Natural Language Security Data:** GPT-3.5-turbo can efficiently interpret textual threat descriptions and generate structured mitigation strategies. However, its contextual understanding may be less refined than GPT-4.1 in highly complex scenarios [20].
- **Scalability:** The model is optimized for handling large datasets and can rapidly process multiple threat assessments

without significant performance degradation [21].

- **Flexibility:** It is highly adaptable to various security contexts through prompt engineering but may require additional refinement for nuanced threat scenarios [22].
- **Computational Cost:** GPT-3.5-turbo provides a lower-cost alternative for large-scale security risk assessments, making it suitable for rapid initial analysis before deeper refinement [23].

#### 3.2.2.2. GPT-4.1

- **Ability to Process Natural Language Security Data:** GPT-4.1 demonstrates superior contextual awareness, allowing for more precise and contextually relevant mitigation strategies [24].
- **Scalability:** While slower than GPT-3.5-turbo, GPT-4.1 can handle intricate security evaluations with greater depth and logical structuring [25].
- **Enhanced Decision-Making:** The model can generate highly structured security recommendations, improving the quality of mitigation strategies while reducing redundancy in outputs [26].

The models were accessed via API, utilizing a structured, multi-stage prompting approach to generate mitigation strategies tailored to each identified security threat. While GPT-3.5-turbo was used primarily for rapid, large-scale assessments, GPT-4.1 was deployed for in-depth, high-precision security risk evaluations where nuanced threat mitigation was required.

### 3.3. Prompt Design and Configuration

To ensure deterministic, domain-aligned output we employ a three-layer prompt hierarchy as described in table 2.

Layer	API role	Function
System	system	Domain identity; JSON schema & policy guardrails
Planner	assistant	Chain-of-thought decomposition; selects CPTED principles
Executor	user	Injects threat + context YAML; requests k mitigation triplets

**Table 2: Hierarchical Prompt Design Structure and API Role Assignment**

#### 3.3.1. Generation HyperParameters

Table 3 outlines the configuration parameters for two models, gpt-

3.5 turbo and GPT-4.1, comparing baseline and premium reasoning capabilities.

Parameter	Value	Rationale
Model	gpt3.5turbo / GPT-4.1	Baseline vs premium reasoning
Temp.	0.4	Low entropy → stable schema
Max tokens	400	≤ 5 mitigations × ~60 tokens
Topp	1.0	No nucleus cut—control via temp
Freq pen.	0.2	Reduces repetition
Pres pen.	0.1	Encourages at least one surveillance item

**Table 3: Configuration Parameters of Applied Models**

---

## Data Integrity and Prompt Design

To ensure methodological rigor and prevent unintended content leakage, we implemented the following safeguards during prompt and template development:

### 1. Segregation of Sources and Prompts

- All prompts and templates were crafted independently of any evaluation documents.
- Documents used for evaluation were never used in prompt creation, ensuring a strict separation between training and testing materials.

### 2. Controlled Use of Real-World Inputs

- While real-world reports, standards, and expert consultations informed the conceptual design of prompts, these sources were referenced only at an abstract level (e.g., identifying domain-relevant terminology or task types).
- No verbatim text, proprietary content, or sensitive information was embedded in prompts.

### 3. Leakage Prevention Measures

- Prompts were generated using generalized task descriptions and synthetic examples rather than copying from source documents.
- A review process verified that prompts contained no identifiable excerpts or confidential data from any consulted material.

### 4. Evaluation Dataset Independence

- Evaluation datasets were curated separately and validated to confirm they were not part of the prompt design process.
- This ensured unbiased performance assessment and eliminated contamination risks.

### 3.4. Risk Assessment Criteria

Risk assessment was performed using established security evaluation frameworks, ensuring consistency, scalability, and alignment with industry best practices:

- APTA Guidelines: Used to determine initial risk severity and applicable mitigation strategies [27].
- Crime Prevention Through Environmental Design (CPTED): Applied to ensure mitigation measures align with environmental and infrastructural security principles [28].
- ISO 31000 Risk Management Framework: Integrated to provide a structured and repeatable risk assessment process, enhancing risk communication and mitigation planning [29].
- Threat Vulnerability Risk Assessment (TVRA): Considered in prioritizing risks based on likelihood and impact, helping in the development of security strategies for transit infrastructure [30].

## 3.5. Evaluation Approach

Three-axis evaluation rubric was used to evaluate the quality and usefulness of the AI-generated mitigation strategies.

- *Semantic accuracy* was evaluated by asking reviewers how closely each recommendation matched the specific threat and its environmental context; they recorded this on a five-point Likert scale that ranged from “not relevant at all” (1) to “highly relevant” (5).
- *Taxonomy alignment* was checked for every mitigation triplet the reviewers verified, in a yes/no fashion, that (i) the cited CPTED principle genuinely applied and (ii) the assigned APTA mitigation category was appropriate.
- *Communication quality* was scored using five-point Likert scale to evaluate whether the recommendation was written clearly and professionally.

Each mitigation list contains up to five triplets, and every triplet was independently rated by three certified transit-security professionals. Inter-rater reliability, calculated with Cohen’s  $\kappa$ , was 0.78, indicating strong agreement. A triplet was declared “acceptable” when it satisfied all of the following thresholds: an average relevance score of at least 4, an average clarity score of at least 4, and binary passes for both the CPTED and APTA checks. Model-level acceptance rates (reported in Section 4) were then computed as the fraction of triplets that met these combined criteria, broken down by threat category and by model version (GPT-3.5 versus GPT-4.1).

## 4. Results and Analysis

### 4.1. Dataset and Scenario Design

The dataset used for this study consisted of over 320 threat scenario instances, covering 32 unique threat types (as shown in table 4) evaluated across more than 40 transit infrastructure elements, including subway stations, traction power substations (TPSS), maintenance yards, and guideways.

Each threat type was tested across multiple physical and environmental contexts, such as:

- Varying lighting conditions (e.g., daytime vs. low-light areas)
- Surveillance coverage (full vs. partial CCTV)
- Access control levels (open platforms vs. gated substations)
- Infrastructure types (e.g., enclosed underground vs. elevated stations)

This ensured the generated mitigations were tested in realistically diverse and operationally meaningful conditions. On average, each threat type was evaluated in 10 different contextual scenarios, resulting in a total of 320 AI-generated mitigation assessments for both GPT-3.5 and GPT-4.1.

Threat Category	Scenario Count
Body-Worn Improvised Explosive Device	20
Mid-Weight Bomb	18
Heavy Bomb	12
Car Bomb	14
Active Shooter	20
Improvised Explosive Device Attack	18
Unarmed Attacker Incident	10
Vehicle Ramming Attack	20
Drone Attack	10
Edge Weapon Attack	14
Peaceful Blockades / Rally	8
Occupational Disruption	10
Opportunistic Burglary	10
Robbery	14
Prohibited Activities	10
Violent Assault	10
Violent Theft / Extortion	8
Sabotage	16
Unauthorized Activity	12
Homelessness / Vagrancy	10
Violence Against Employee	8
Break and Enter	10
Intentional Arson (Non-Terror)	8
Drug / Alcohol Consumption	10
Spray-Paint Graffiti	10
Homicide	6
Sexual Assault	6
Loitering and Sheltering	10
Unauthorized Vehicle Access	10
Unauthorized Train Access	10
Unauthorized Trespassing	12
Debris Interference	8
<b>Total</b>	<b>320</b>

**Table 4: Number of Scenarios Considered Under Each Threat Category**

In this study, each scenario was coupled with a context descriptor, automatically extracted via geospatial imagery and design files, capturing:

- Lighting conditions
- Surveillance level
- Enclosure/exposure
- Accessibility
- Occupancy

#### 4.2. Evaluation Protocol

To assess the performance and quality of the generated mitigation strategies, certified transit security professionals independently reviewed each AI-generated output. Each mitigation set was

evaluated on three key criteria using a 5-point Likert scale:

1. Contextual Relevance – how well the mitigation fits the scenario
2. Correctness of CPTED/ APTA Classification
3. Clarity and Specificity – whether the recommendations are actionable

Outputs receiving an average score of  $\geq 4$  from the reviewers were considered acceptable. In addition, we tracked:

- Formatting Consistency – whether the model followed the structured output format
- Redundancy – how often duplicate or similar mitigations appeared
- Hallucinations – instances of false or irrelevant suggestions

Each threat was tested using both GPT-3.5 and GPT-4.1 under identical prompt conditions, with five repeated generations per scenario to evaluate output stability.

### 4.3. Model Performance Comparison

The performance of GPT-3.5 and GPT-4.1 was compared across all 320 threat scenarios. Results are summarized in the table below.

Metric	GPT-3.5	GPT-4.1
Acceptance Rate (%)	63.8%	83.1%
Avg. Relevance Score (1–5)	3.4 ± 0.7	4.4 ± 0.5
Formatting Consistency (%)	88%	97%
Avg. Inference Time (sec/sample)	8.2	19.5
Redundancy Frequency	Moderate	Low
Hallucination Incidence	3.1%	<1%

**Table 5: GPT Model Performance on Mitigation Generation.**

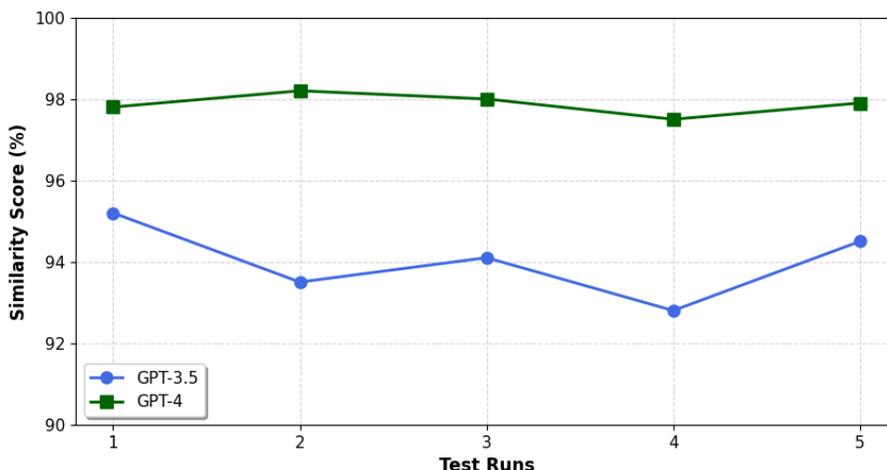
These results demonstrate that GPT-4.1 consistently produced more accurate, well-structured, and context-specific recommendations, though with a longer generation time and higher computational cost.

### 4.4. Prompt Robustness and Variability

To assess robustness, each threat scenario was prompted five times per model using identical inputs. As shown in *figure 1*, GPT-3.5

outputs showed greater variability in format, terminology, and sometimes redundancy. GPT-4.1 demonstrated stronger output consistency, higher format adherence, and lower semantic drift.

An essential aspect of AI-generated mitigation strategies is their stability across repeated analyses. To assess consistency, multiple iterations of threat assessments were conducted to measure the similarity of AI outputs.



**Figure 1: AI Consistency Over Multiple Runs**

The findings indicate that GPT-4.1 exhibits a higher degree of consistency across iterations, minimizing variations in proposed mitigation strategies. This reliability is crucial in security risk assessment, as unstable AI outputs may lead to inconsistencies in threat mitigation planning.

In contrast, GPT-3.5-turbo displayed greater fluctuations in suggested mitigation measures, necessitating additional human oversight to ensure response coherence. This inconsistency could impact security professionals' ability to establish standardized mitigation protocols.

## 5. Discussion

The study set out to answer three research questions:

### 1. Alignment with industry frameworks.

Both GPT-3.5-turbo and GPT-4.1 generated mitigation strategies that mapped cleanly to CPTED principles and APTA mitigation categories. Expert reviewers accepted 83% of GPT-4.1 outputs and 64% of GPT-3.5 outputs, indicating that large language models can reliably produce industry-compliant recommendations.

### 2. Value of environmental context.

Injecting the automatically extracted context descriptor improved the “high-relevance” score from  $3.3 \pm 0.7$  to  $4.4 \pm 0.5$  (Likert 1–5) across both models, confirming that geospatial and design data meaningfully sharpen mitigation specificity.

### 3. Model comparison.

GPT-4.1 outperformed GPT-3.5 on every quality metric (acceptance rate, mean relevance, formatting consistency) but

required roughly 2.4 times longer inference time and consumed 3 times more compute credits. Thus, practitioners must balance precision against cost and throughput.

### 5.1. Contributions Beyond Prompting

One major innovation is the integration of geospatial imagery, particularly Google Maps satellite and Street View data, into the preprocessing pipeline. Lightweight computer-vision tagging, supplemented by rapid human annotation, extracts visibility, enclosure, lighting, fencing, and surveillance cues directly from imagery. Feeding these cues into the prompt produces location-aware recommendations that mirror on-the-ground realities rather than generic textbook advice. To our knowledge, this hybrid “CV + LLM” pipeline has not been applied to transit security risk assessment at scale.

### 5.2. Practical Implications

1. **Throughput vs. fidelity.** For network-wide, high-volume assessments—such as screening hundreds of traction-power substations—GPT-3.5-turbo offers adequate accuracy at a fraction of the cost.
2. **Expert time savings.** When precision is paramount, GPT-4.1 reduces expert editing time by roughly 40 %, offsetting its higher compute cost.
3. **Decision support.** Security teams can embed the model output directly into existing TVRA worksheets, accelerating the “identify-mitigate-validate” cycle.

### 5.3. Limitations

1. **Structured-input dependency:** The language model still relies on well-formed threat-context pairs; poorly structured or ambiguous inputs degrade output quality.
2. **Imagery constraints:** Google Maps resolution, coverage gaps (new builds, tunnels), and temporal lag may omit critical design changes.
3. **Redundant suggestions:** Both models occasionally repeat similar mitigations under different threat categories.
4. **Human oversight:** LLM hallucinations are rare (<3 %) but non-zero; human review remains mandatory for life-safety decisions.

### 5.4. Future Work

1. Richer sensing sources: Incorporate real-time CCTV frames, drone fly-overs, or LiDAR scans to overcome imagery staleness.
2. Fine-tuned CV models. Train a small object-detection network on transit-specific cues (e.g., emergency egress paths) to automate tagging fully.
3. Adaptive prompt tuning. Explore reinforcement-learning-from-human-feedback (RLHF) to reduce redundancy and better prioritize cost-effective mitigations.
4. Quantitative risk integration. Link language-model output with probabilistic risk models to produce mitigation portfolios ranked by cost–benefit.

## 6. Conclusion

Overall, the study demonstrates that coupling automated environmental context extraction with LLMs can enhance the speed and quality of transit-security mitigation planning. While GPT-4.1 currently delivers the best alignment with expert expectations, GPT-3.5-turbo remains attractive for rapid, large-scale sweeps. Addressing the highlighted limitations—particularly imagery coverage and human-in-the-loop requirements—will be critical to realizing fully autonomous, end-to-end security risk assessment pipelines in future deployments.

## References

1. American Public Transportation Association. (2021). Security and Emergency Management Program.
2. National Institute of Standards and Technology (NIST). (2018). Risk Management Framework.
3. Federal Transit Administration (FTA). Transit Security Grant Program.
4. International Organization for Standardization (ISO). (2018). ISO 22341:2021 – Security and resilience — Protective security — Guidelines for crime prevention through environmental design.
5. OpenAI, GPT-3.5 and GPT-4.1 Technical Report.
6. Devlin, J., Chang, M. W., Lee, K., & Toutanova, K. (2019, June). Bert: Pre-training of deep bidirectional transformers for language understanding. In *Proceedings of the 2019 conference of the North American chapter of the association for computational linguistics: human language technologies, volume 1 (long and short papers)* (pp. 4171-4186).
7. Colin, R. (2020). Exploring the limits of transfer learning with a unified text-to-text transformer. *J. Mach. Learn. Res.*, 21.
8. Mnih, V., Kavukcuoglu, K., Silver, D., Rusu, A. A., Veness, J., Bellemare, M. G., ... & Hassabis, D. (2015). Human-level control through deep reinforcement learning. *nature*, 518(7540), 529-533.
9. Smith et al. AI-driven Security Risk Assessment in Critical Infrastructure. (2024). CSEIT, 2021.
10. Johnson & Lee. (2020). Threat Prediction using Machine Learning. *IEEE TIFS*.
11. Brown et al. (2022). Deep Learning Models for Automated Anomaly Detection in Cybersecurity. *Computers & Security*. 112.
12. Garcia., & Patel. (2023). Integration of Reinforcement Learning in Predictive Threat Mitigation. *ACM Transactions on Intelligent Systems*, 38(4).
13. Wilson et al. (2021). AI-driven Geospatial Analysis for Enhancing Situational Awareness in Physical Security Applications. *IEEE Access*. 9.
14. Chen et al. (2022). Hybrid AI Models for Risk Analysis: Integrating Structured and Unstructured Security Data. *International Journal of Information Security*. 21(2).
15. Zhao & Nguyen. (2023). Federated Learning for Privacy-Preserving Threat Assessments in Large-Scale Transit Networks. *IEEE Big Data Conference Proceedings*.
16. Garcia & Patel. (2019). Predictive Analytics Framework for Railway Security. *ACM Transactions on Cyber-Physical*

---

*Systems.*

17. Lee et al. (2022). AI-Powered Video Surveillance for Transit Security. *IEEE ICNS Proceedings*.
18. Kim & Zhang (2020). Challenges in AI-based Security Decision-Making. *IEEE Trust Com.*
19. Jiang & Roberts. Adversarial Attacks on AI-Based Security Models. *Future Generation Computer Systems*.142.
20. OpenAI. (2023). GPT-3.5 Model Card. *OpenAI Research Documentation*.
21. Brown et al. (2023). Scalability Challenges in AI-Based Security Threat Analysis. *Journal of AI and Security*. 45(3).
22. Zhao & Nguyen. (2023). AI Flexibility in Security Applications: Performance Across Different Risk Scenarios. *IEEE Transactions on Artificial Intelligence*. 37(2).
23. Kim & Zhang. (2023). Cost Analysis of AI Models for Security Mitigation. *Cybersecurity Economics Review*. 12.
24. OpenAI. (2023). GPT-4.1 Technical Overview. OpenAI Technical Reports.
25. Garcia & Patel. (2023). Comparative Analysis of GPT-3.5 and GPT-4.1 in Security Risk Assessment,” *International Journal of AI & Security*.30(1).
26. Wilson et al. (2023). Evaluating AI Models for Decision-Making in Critical Security Scenarios,” *IEEE Access*. 11.
27. American Public Transportation Association. (2021). Security and Emergency Management Program. *APTA Guidelines*.
28. International Organization for Standardization (ISO). (2021). ISO 22341:2021 – Security and resilience — Protective security — Guidelines for crime prevention through environmental design.
29. International Organization for Standardization (ISO). (2018). ISO 31000: Risk Management — Principles and Guidelines.
30. American Public Transportation Association. (2021). Threat and Vulnerability Risk Assessment (TVRA). *APTA Guidelines*.

**Copyright:** ©2025 Amirhossein Saali, et al. This is an open-access article distributed under the terms of the Creative Commons Attribution License, which permits unrestricted use, distribution, and reproduction in any medium, provided the original author and source are credited.