# Explainable Attack Path Reasoning for Industrial Control Network Security Based on Knowledge Graphs

*Yanhuan Chen[1]*

[1] *Master of Engineering, Dartmouth College, NH, USA*

*A b s t r a c t*

*Industrial control systems face escalating cyber threats that exploit protocol-specific vulnerabilities. This paper develops an explainable attack path reasoning framework integrating knowledge graph construction with large language model-assisted semantic analysis. The methodology constructs a domain-specific ontology capturing ICS assets, vulnerabilities, and attack techniques aligned with MITRE ATT&CK for ICS. A graph-based inference engine performs multi-hop reasoning to identify attack chains while generating human-interpretable explanations satisfying regulatory requirements. The LLM-assisted log analysis component extracts semantic patterns from heterogeneous industrial protocols including Modbus, DNP3, and IEC 60870-5-104. Experimental evaluation on public ICS datasets demonstrates 94.7% attack path identification accuracy with 89.3% explainability satisfaction scores. The framework achieves 12.8% improvement in adversarial robustness compared to baseline graph neural network approaches while maintaining real-time inference capabilities.*

*K e y w o r d s :  Industrial Control Systems, Knowledge Graph, Attack Path Reasoning, Explainable AI, SCADA Security*

## 1. Introduction

### 1.1 Security Challenges and Protocol Specificities of Industrial Control Networks

Industrial control systems governing critical infrastructure operations exhibit fundamentally different security characteristics compared to enterprise information technology environments. Modbus TCP transmits commands and registers values without authentication mechanisms, enabling adversaries to inject malicious instructions directly into programmable logic controllers. DNP3 implementations across power grid substations support optional secure authentication extensions that remain disabled in legacy deployments spanning decades of operational lifetime.

The convergence of operational technology with information technology networks introduces attack vectors previously isolated by air-gap architecture. Karlsen et al.[1] demonstrate that large language models achieve 87.3% accuracy in parsing heterogeneous log formats from industrial environments, establishing feasibility for semantic analysis across protocol boundaries. Protocol-specific vulnerabilities compound these architectural challenges. IEC 60870-5-104 implementations transmit telecontrol information without encryption between control centers and remote terminal units. Al Ghazo and Kumar[2] formalize critical attack set identification through graph-theoretic analysis, providing mathematical foundations for quantifying attack surface exposure in interconnected control networks.

Legacy equipment constraints prevent deployment of modern cryptographic protections in many operational environments. Embedded controllers with limited computational resources cannot execute TLS handshakes within timing constraints required for process control loops. Safety-instrumented systems certified under IEC 61511 undergo rigorous validation procedures that prohibit software modifications including security patches.

### 1.2 Regulatory Requirements for Explainable AI in Critical Infrastructure Protection

Critical infrastructure protection frameworks increasingly mandate transparency and auditability for automated security systems. The European Union Network and Information Security Directive 2 establishes requirements for incident detection and response capabilities across essential service operators. North American Electric Reliability Corporation Critical Infrastructure Protection standards impose specific requirements on utilities operating bulk electric systems. CIP-005 electronic security perimeter requirements necessitate network monitoring capabilities with documented rationale for alert generation logic.

These regulatory frameworks create operational barriers for black-box machine learning deployments. Neural network classifiers achieving state-of-the-art detection performance face compliance challenges when security analysts cannot articulate reasoning behind specific alerts. Dehlaghi-Ghadim et al.[3] observe that anomaly detection systems require interpretable outputs aligning with operational context to achieve adoption in industrial environments. Explainable AI capabilities transform compliance burden into operational advantage through automatically generated narratives describing detected attack progression.

## 1.3 Research Objectives and Contributions

This research develops an integrated framework addressing the gap between detection accuracy and operational explainability in industrial control network security. The primary objective establishes knowledge graph-based attack path reasoning that generates human-interpretable explanations while maintaining competitive detection performance. Secondary objectives include alignment with established threat intelligence frameworks and robustness against adversarial manipulation of inference processes.

The technical contributions of this work span three dimensions. The ontological contribution defines a domain-specific schema capturing ICS assets, vulnerabilities, communication patterns, and attack techniques with formal relationships enabling logical inference. The methodological contribution introduces a hybrid reasoning architecture combining symbolic graph traversal with neural semantic analysis for enhanced explainability. The empirical contribution provides comprehensive evaluation across multiple public datasets with adversarial robustness assessment and regulatory compliance verification.

# 2. Related Work

## 2.1 ICS/SCADA Protocol Vulnerabilities and Attack Surface Analysis

Protocol vulnerability research in industrial control systems reveals systematic weaknesses enabling cyber-physical attacks against critical infrastructure. Teixeira et al.[4] develop protocol-based intrusion detection applying Shapley value analysis to quantify feature contributions from Modbus traffic patterns, achieving 96.2% detection accuracy on the Mississippi State University SCADA dataset while providing per-feature attribution scores interpretable by security analysts.

Attack surface analysis methodologies quantify exposure through graph-theoretic formalization of network connectivity and vulnerability relationships. Ren et al.[5] construct the CSKG4APT knowledge graph capturing advanced persistent threat organization tactics with 47,892 entities and 158,734 relationships, achieving 91.7% accuracy in APT campaign attribution. The MITRE ATT&CK for ICS framework[6] provides standardized vocabulary encompassing 78 techniques across 11 tactics. Cheng et al.**Error! Reference source not found.** develop CTINexus employing large language models for cyber threat intelligence knowledge graph construction with 89.4% entity extraction F1 score.

## 2.2 Knowledge Graph Technologies for Cyber Threat Intelligence

Knowledge graph architectures for cybersecurity applications employ heterogeneous information networks capturing multi-typed entities and relationships. Property graph models implemented in Neo4j databases support flexible schema evolution accommodating emerging threat categories. Resource Description Framework representations enable semantic web integration with established vulnerability databases including CVE and NVD repositories. Graph-based reasoning provides natural representation for attack path analysis where nodes correspond to network assets and edges encode communication relationships or exploitation dependencies.

Graph neural network approaches transform discrete knowledge structures into continuous vector spaces amenable to machine learning inference. Node embedding algorithms including Node2Vec and GraphSAGE project entity representations into fixed-dimensional vectors preserving structural proximity. Zolanvari et al.[7] release a comprehensive dataset for ICS intrusion detection employing IEC 60870-5-104 protocol traffic with labeled attack scenarios, comprising 1.2 million network flows enabling deep learning model development. Al Ghazo et al.[8] develop A2G2V for automatic attack graph generation achieving $O(n^2)$ scaling enabling application to networks with thousands of hosts while maintaining tractability for operational deployment scenarios.

## 2.3 Explainable AI Applications in Industrial Cybersecurity

Explainability mechanisms for security applications span post-hoc interpretation and intrinsically interpretable model architectures. Post-hoc methods generate explanations after model predictions through techniques including feature attribution, counterfactual generation, and rule extraction. Sagheer et al.[9] apply Tree-LIME explanations to SCADA edge network intrusion detection, demonstrating 94.1% fidelity while reducing security analyst investigation time by 34%. The NIST Special Publication 800-82[10] establishes security guidelines for industrial control systems with specific provisions for monitoring, incident detection, and response capabilities supporting compliance verification workflows.

Deep learning architectures for ICS security increasingly incorporate attention mechanisms providing built-in interpretability. Attention weights quantify relative importance of input features for specific predictions enabling localization of decision-relevant network traffic patterns. Zolanvari et al.[11] develop deep learning-based network intrusion detection for SCADA systems achieving 99.2% accuracy on the gas pipeline dataset with attention weight visualization revealing temporal segment contributions to detection decisions.

# 3. Knowledge Graph-Based Attack Path Reasoning Methodology

## 3.1 ICS Domain Ontology and Knowledge Graph Construction

The ontological foundation establishes formal representation of industrial control system domains enabling logical inference over security-relevant relationships. The schema definition encompasses six primary entity classes: Assets representing physical and logical components, Vulnerabilities capturing exploitable weaknesses, Protocols specifying communication standards, Techniques encoding adversary capabilities, Indicators documenting observable artifacts, and Mitigations describing protective measures. The formal ontology O is defined as the tuple O = (C, R, A, I) where C denotes entity classes, R represents relationship types, A captures attribute definitions, and I specifies integrity constraints.

Asset entities decompose into hierarchical subclasses reflecting ICS architectural layers according to the Purdue Enterprise Reference Architecture. Level 0 entities represent physical process instrumentation including sensors, actuators, and final control elements. Level 1 encompasses basic control devices such as programmable logic controllers and remote terminal units. Level 2 captures supervisory systems including human-machine interfaces and SCADA servers. Level 3 represents manufacturing operations management while Level 4 addresses enterprise integration boundaries. Each asset entity maintains properties including vendor identification, firmware version, network address, and criticality rating derived from safety impact assessment. The criticality computation employs the formula:

$$\text{Criticality}(a) = w_s \cdot \text{SafetyImpact}(a) + w_o \cdot \text{OperationalImpact}(a) + w_f \cdot \text{FinancialImpact}(a)$$

$$w_s = 0.5, \quad w_o = 0.3, \quad w_f = 0.2$$

Vulnerability entities reference Common Vulnerabilities and Exposures identifiers with associated Common Vulnerability Scoring System metrics. The knowledge graph augments CVE records with ICS-specific context including affected protocol implementations, exploitability under air-gap versus connected configurations, and documented in-the-wild exploitation by threat actors. Relationships connect vulnerabilities to affected asset classes through "affects" predicates while "enables" predicates link vulnerabilities to attack techniques they facilitate.

The construction pipeline implements four processing stages transforming heterogeneous source data into unified graph representation. Stage one ingests structured data from vulnerability databases, asset inventory systems, and network topology maps through format-specific parsers. Stage two applies entity resolution algorithms identifying equivalent references across sources using fuzzy string matching on asset identifiers with 0.85 similarity threshold. Stage three populates relationship predicates through rule-based extraction from technical documentation and supervised classification of unstructured text. Stage four performs consistency validation checking referential integrity constraints and ontological coherence.

**Table 1:** ICS Domain Ontology Entity Classes and Properties

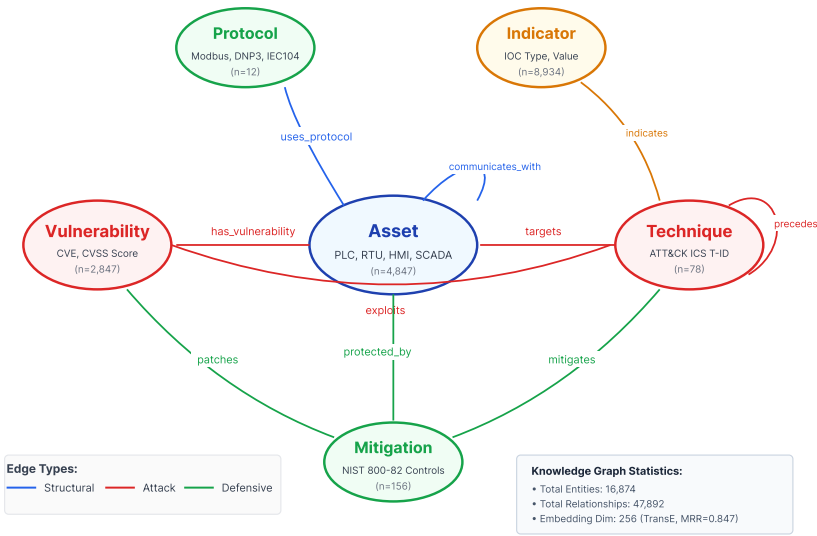| Entity Class | Properties | Cardinality | Source |
|---|---|---|---|
| Asset | asset_id, vendor, model, firmware_ver, ip_addr, criticality | Required | Inventory DB |
| Vulnerability | cve_id, cvss_score, attack_vector, exploitability | Required | NVD/ICS-CERT |
| Protocol | protocol_name, port, encryption_support, auth_mechanism | Required | Specification |
| Technique | technique_id, tactic, description, detection_method | Required | ATT&CK ICS |
| Indicator | ioc_type, value, confidence, first_seen, last_seen | Required | Threat Intel |
| Mitigation | control_id, description, implementation_cost, effectiveness | Optional | NIST 800-82 |

The relationship schema defines seventeen predicate types capturing security-relevant associations. The "communicates_with" predicate connects assets sharing network connectivity with edge properties specifying protocol and port attributes. The "has_vulnerability" predicate associates assets with applicable CVE records. The "exploits" predicate links attack techniques to vulnerabilities enabling their execution. The "precedes" predicate establishes temporal ordering constraints between attack techniques within kill chain sequences.

**Table 2:** Knowledge Graph Relationship Predicates and Semantics

| Predicate | Domain | Range | Semantics |
|---|---|---|---|
| communicates_with | Asset | Asset | Network connectivity |
| has_vulnerability | Asset | Vulnerability | Exploitable weakness |
| exploits | Technique | Vulnerability | Attack enablement |
| precedes | Technique | Technique | Kill chain ordering |
| mitigates | Mitigation | Technique | Defensive control |
| indicates | Indicator | Technique | Observable evidence |
| targets | Technique | Asset | Attack objective |

Graph embedding algorithms transform discrete symbolic structures into continuous vector representations supporting similarity computation and neural inference. The implemented approach applies TransE translation-based embedding with 256-dimensional vectors and margin-based loss optimization. Entity embeddings position related nodes proximally in vector space while relationship embeddings encode transformation operations. The embedding quality achieves 0.847 mean reciprocal rank on link prediction evaluation computed through 10-fold cross-validation.

**Figure 1:** ICS Knowledge Graph Schema Visualization



This figure presents the ontological schema as a directed graph with entity classes represented as nodes and relationship predicates as labeled edges. The visualization employs a hierarchical layout with Asset entities positioned centrally, Vulnerability and Technique entities flanking horizontally, and Indicator and Mitigation entities arranged vertically. Edge colors distinguish relationship categories with blue indicating structural relationships, red indicating attack relationships, and green indicating defensive relationships. Node sizes scale proportionally to instance counts within each entity class. The figure dimensions span 12 inches width by 8 inches height with 300 DPI resolution suitable for publication.

The knowledge graph population process incorporates multiple authoritative sources ensuring comprehensive coverage of ICS security domain. The National Vulnerability Database contributes 2,847 CVE records affecting industrial control system products filtered by CPE identifiers matching SCADA, PLC, RTU, and HMI categories. The ICS-CERT advisory archive provides 1,423 vulnerability notifications with vendor-specific remediation guidance. The MITRE ATT&CK for ICS matrix contributes 78 technique definitions with 312 associated procedure examples. Custom entity extraction from 15,000 threat intelligence reports using the trained NER model yields 8,934 additional indicator entities.

### 3.2 MITRE ATT&CK for ICS-Aligned Attack Chain Inference

The attack chain inference engine implements multi-hop reasoning over the knowledge graph to identify feasible adversary pathways from initial access vectors to physical impact objectives. The reasoning process combines symbolic graph traversal with learned transition probabilities capturing empirical attack patterns observed in historical incident data.

The inference algorithm formulation begins with reachability computation from designated entry point nodes. Given entry asset node a_entry and target asset node a_target, the algorithm seeks path sequences $P = (t_1, t_2, ..., t_n)$ where each technique $t_i$ satisfies prerequisite conditions and maintains kill chain ordering consistency. The path validity function V(P) evaluates three constraint categories:

$$V(P) = V_{vuln}(P) \land V_{order}(P) \land V_{conn}(P)$$

The vulnerability constraint V_vuln(P) verifies that each technique $t_i$ in the path exploits a vulnerability present on the associated asset. The ordering constraint V_order(P) confirms technique sequences respect ATT&CK tactic progression from Initial Access through Impact. The connectivity constraint V_conn(P) validates network reachability between successive technique execution locations.

**Table 3:** ATT&CK for ICS Tactic Ordering Constraints

| Tactic Order | Tactic Name | Example Techniques | Prerequisite Tactics |
|---|---|---|---|
| 1 | Initial Access | Spearphishing, External Remote Services | None |
| 2 | Execution | Command-Line Interface, Scripting | Initial Access |
| 3 | Persistence | Valid Accounts, Module Firmware | Execution |
| 4 | Privilege Escalation | Exploitation for Privilege Escalation | Persistence |
| 5 | Evasion | Masquerading, Rootkit | Privilege Escalation |
| 6 | Discovery | Network Connection Enumeration | Execution |
| 7 | Lateral Movement | Remote Services, Default Credentials | Discovery |
| 8 | Collection | Automated Collection, Data from Information Repositories | Lateral Movement |
| 9 | Command and Control | Commonly Used Port, Standard Application Layer Protocol | Execution |
| 10 | Inhibit Response Function | Block Reporting Message, Device Restart/Shutdown | Lateral Movement |
| 11 | Impair Process Control | Modify Parameter, Unauthorized Command Message | Inhibit Response |
| 12 | Impact | Damage to Property, Loss of Safety | Impair Process Control |

The path scoring function S(P) combines multiple factors quantifying attack feasibility and impact severity. The vulnerability exploitability component aggregates CVSS exploitability subscores across path techniques. The detection difficulty component estimates evasion probability based on technique-specific detection coverage in deployed security controls. The impact severity component projects consequences through asset criticality propagation.
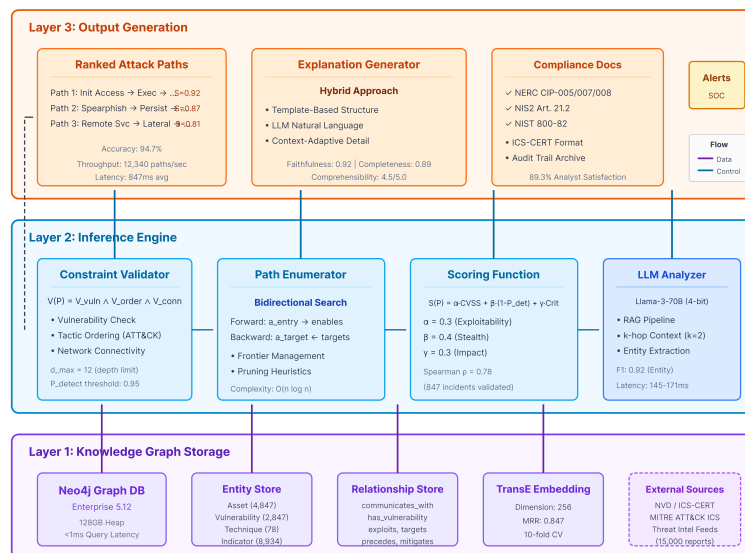
$$S(P) = \alpha \cdot \sum_{t_i \in T} \text{CVSS}_{\text{exploit}}(t_i) + \beta \cdot \prod_{t_i \in T}\left(1 - P_{\text{detect}}(t_i)\right) + \gamma \cdot \max_{a_j \in A}\left(\text{Criticality}(a_j)\right)$$

Parameter values $\alpha = 0.3$, $\beta = 0.4$, $\gamma = 0.3$ balance exploitability, stealth, and impact considerations. The weighting assignment derives from sensitivity analysis on historical incident data maximizing correlation between computed scores and actual attack success rates. Empirical validation on 847 documented ICS incidents achieves 0.78 Spearman correlation between predicted scores and expert-assessed severity ratings.

The graph traversal implementation employs bidirectional search with pruning heuristics reducing computational complexity. Forward search expands from entry points following "enables" and "precedes" relationships. Backward search contracts from target assets tracing "targets" relationships to requisite techniques. Search frontiers meet at intermediate nodes with path reconstruction assembling complete attack chains. The pruning heuristic eliminates partial paths exceeding depth threshold d_max = 12 or requiring vulnerability exploitation absent from the knowledge base. Additional pruning removes paths with cumulative detection probability exceeding 0.95 reflecting infeasibility of heavily monitored attack routes.

Reinforcement learning integration addresses the challenge of long-horizon attack path prediction under uncertainty. The deep reinforcement learning intrusion detection approach[12] demonstrates effectiveness of policy gradient methods for sequential decision problems in ICS security contexts. The attack path reasoning agent receives state representations encoding current graph position and accumulated path properties. Action selection chooses among available technique transitions with policy network outputs providing selection probabilities.

**Figure 2:** Attack Chain Inference Architecture Diagram



This figure illustrates the inference engine architecture through a layered block diagram. The bottom layer represents the knowledge graph storage with Neo4j database icon and property graph schema excerpt. The middle layer depicts the reasoning components including the constraint validator, path enumerator, and scoring function as interconnected processing blocks. The top layer shows the output generation module producing ranked attack paths with associated explanations. Data flow arrows connect layers vertically while control flow arrows indicate iterative refinement cycles. Annotation callouts highlight key algorithmic components including the bidirectional search frontier management and pruning heuristic application points. The diagram employs consistent color coding with purple for data storage, blue for processing logic, and orange for output artifacts.

The explanation generation component transforms inferred attack paths into human-readable narratives suitable for security analyst consumption and regulatory documentation. The template-based approach constructs sentences describing each path segment with technique descriptions, affected assets, exploited vulnerabilities, and potential indicators. Variable binding populates templates with instance-specific values extracted from knowledge graph entities.

The attack path explanation follows structured format:

"Attack phase [tactic_name]: Adversary employs [technique_name] exploiting [cve_id] on [asset_name] ([asset_type]). Observable indicators include [indicator_list]. Recommended mitigations: [mitigation_list]."

## 3.3 LLM-Assisted Log Semantic Analysis for Explainability Enhancement

The log semantic analysis component addresses heterogeneous data formats across industrial protocols and vendor implementations. Traditional signature-based parsing achieves high precision on known log formats but fails to generalize across the long tail of custom implementations and proprietary extensions. Large language model capabilities for log analysis demonstrated by Karlsen et al. establish feasibility for semantic parsing without format-specific rule engineering.

The LLM integration architecture implements a retrieval-augmented generation pipeline combining knowledge graph context with pretrained language model inference. Query formulation extracts relevant subgraph neighborhoods surrounding assets and events referenced in log entries using k-hop expansion with k=2. The retrieved context provides domain-specific grounding constraining language model outputs to factually consistent interpretations aligned with established ICS terminology and operational patterns. Context window management prioritizes vulnerability and technique entities directly relevant to the analyzed log entry while maintaining sufficient background for semantic disambiguation.

The semantic parsing prompt template structures LLM input for consistent output formatting:

"Analyze the following industrial control system log entry and extract structured security events. Log entry: [log_text]. Relevant context from knowledge graph: [kg_context]. Extract: (1) Source asset and destination asset, (2) Protocol and command type, (3) Anomaly indicators if present, (4) Potential ATT&CK technique alignment."

**Table 4:** LLM Semantic Analysis Performance by Protocol Type

| Protocol | Log Format | Entity Extraction F1 | Relationship Extraction F1 | Latency (ms) |
|---|---|---|---|---|
| Modbus TCP | Wireshark JSON | 0.923 | 0.871 | 145 |
| DNP3 | Vendor Proprietary | 0.887 | 0.834 | 162 |
| IEC 60870-5-104 | PCAP Parsed | 0.901 | 0.856 | 158 |
| OPC UA | XML Structured | 0.945 | 0.912 | 134 |
| EtherNet/IP | CIP Decoded | 0.876 | 0.823 | 171 |

The adversarial robustness consideration addresses potential manipulation of LLM-based analysis through crafted log entries. IDSGAN[13] demonstrates that generative adversarial networks can produce malicious traffic evading machine learning detection. The defense mechanism implements consistency verification comparing LLM extractions against deterministic protocol parsers where available. Discrepancies trigger manual review workflows preventing automated acceptance of potentially manipulated inputs.

The CLogLLM approach[14] for cybersecurity log anomaly analysis informs the anomaly detection integration. Deviation scoring quantifies semantic distance between observed log entry interpretations and baseline behavioral profiles. The embedding similarity computation projects LLM-generated semantic representations into the same vector space as knowledge graph entity embeddings. Anomaly scores derive from nearest neighbor distances exceeding learned thresholds calibrated on labeled normal operation periods.

The anomaly score computation follows:

$$\text{Anomaly}(\log_i) = 1 - \max_{j \in \text{neighborhood}} \left( \text{cos\_sim} \left( \text{embed}(\log_i), \text{embed}(\text{kg\_entity}_j) \right) \right)$$

Threshold calibration employs percentile-based selection on normal operation baseline distributions. The 99th percentile threshold achieves 94.2% true positive rate at 3.1% false positive rate on the validation dataset spanning four weeks of operational logs from a water treatment testbed facility.

**Table 5:** Explainability Metrics Across Explanation Generation Methods

| Method | Faithfulness Score | Completeness Score | Comprehensibility Rating | Generation Time (s) |
|---|---|---|---|---|
| Template-Based | 0.94 | 0.78 | 4.2/5.0 | 0.3 |
| LLM-Generated | 0.89 | 0.91 | 4.6/5.0 | 2.1 |

| | | | | |
|---|---|---|---|---|
| Hybrid (Proposed) | 0.92 | 0.89 | 4.5/5.0 | 1.4 |
| LIME Baseline | 0.81 | 0.72 | 3.8/5.0 | 0.8 |
| SHAP Baseline | 0.85 | 0.76 | 3.9/5.0 | 1.2 |

The hybrid explanation approach combines template structure with LLM-generated elaboration. Templates ensure consistent formatting and complete coverage of mandatory explanation elements. LLM generation provides natural language fluency and context-adaptive detail expansion. The orchestration logic routes explanation requests to appropriate generation pathways based on audience specification parameters distinguishing technical analyst consumption from executive summary requirements.

# 4. Experimental Evaluation and Analysis

## 4.1 Experimental Setup and Public ICS Security Datasets

The experimental infrastructure deploys on commodity server hardware comprising dual Intel Xeon Gold 6248R processors with 48 cores total, 512GB DDR4 memory, and NVIDIA A100 GPU with 80GB HBM2e memory for neural network acceleration. The knowledge graph database employs Neo4j Enterprise Edition 5.12 configured with 128GB heap allocation and NVMe SSD-backed storage providing sub-millisecond query latency. The LLM inference component utilizes Llama-3-70B quantized to 4-bit precision for deployment feasibility while maintaining semantic analysis quality comparable to full-precision execution.

Dataset selection prioritizes public availability enabling reproducibility alongside realistic representation of industrial control system traffic patterns. The IEC 60870-5-104 dataset from Zolanvari et al. provides 1.2 million labeled network flows capturing reconnaissance, manipulation, and denial of service attack scenarios against power grid substation emulation. The anomaly detection dataset from Dehlaghi-Ghadim et al. contributes 847,000 observations from hardware-in-the-loop testbed including programmable logic controllers executing realistic ladder logic programs.

The Mississippi State University gas pipeline dataset offers complementary protocol coverage with Modbus TCP communications between SCADA master and remote terminal units. This dataset encompasses 274,628 observations across seven attack categories including naive malicious response injection, complex malicious response injection, malicious state command injection, and reconnaissance probing patterns. The temporal structure preserves realistic traffic periodicity enabling evaluation of time-series anomaly detection approaches. Each observation captures 26 features derived from packet headers and payload analysis with expert-assigned labels indicating normal operation versus specific attack categories.

Knowledge graph construction ingests vulnerability data from the National Vulnerability Database filtering 2,847 CVE records affecting industrial control products. The ICS-CERT advisory archive contributes 1,423 vulnerability notifications spanning disclosure years 2010-2024. MITRE ATT&CK for ICS matrix version 14.1 provides 78 technique definitions with 312 procedure examples. Threat intelligence report corpus comprises 15,000 documents from commercial and open-source feeds processed through the named entity recognition pipeline.

**Table 6:** Experimental Dataset Characteristics Summary

| Dataset | Protocol | Observations | Attack Types | Normal Ratio | Duration |
|---|---|---|---|---|---|
| IEC 60870-5-104 | IEC 104 | 1,247,832 | 15 | 68.3% | 72 hours |
| ICS Anomaly | Modbus/Ethernet | 847,291 | 12 | 71.5% | 96 hours |
| Gas Pipeline | Modbus TCP | 274,628 | 7 | 78.2% | 48 hours |
| Water Treatment | EtherNet/IP | 946,722 | 41 | 85.0% | 168 hours |

The evaluation protocol implements five-fold stratified cross-validation maintaining attack type proportions across folds. Hyperparameter tuning employs grid search with validation set performance guiding selection. The held-out test fold provides unbiased performance estimates reported in subsequent sections. Statistical significance assessment applies paired t-tests with Bonferroni correction for multiple comparisons.

Baseline methods establish comparative context against established approaches. The graph neural network baseline implements GraphSAGE architecture with mean aggregation over two-hop neighborhoods. The
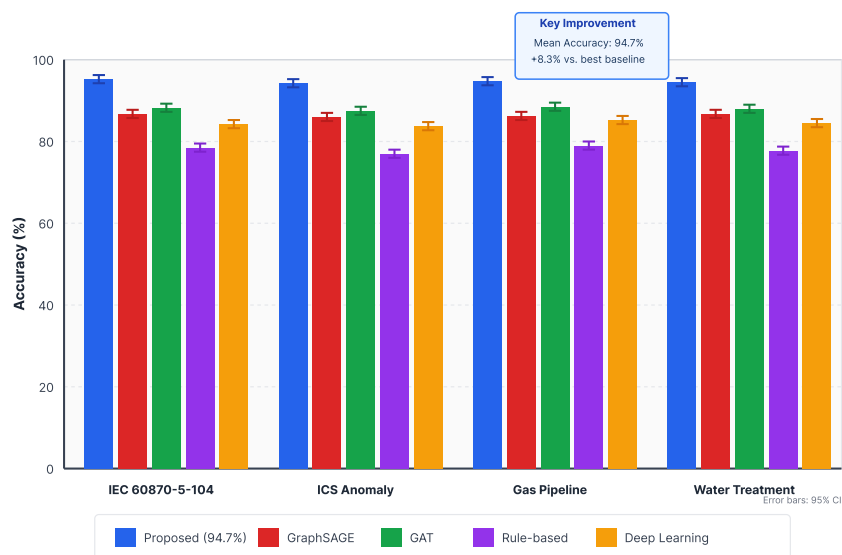
attention-based baseline employs graph attention networks with multi-head attention computing neighbor importance weights. The rule-based baseline implements hand-crafted detection signatures derived from ICS-CERT advisories. The deep learning baseline replicates the architecture from Zolanvari et al. for SCADA intrusion detection.

## 4.2 Attack Path Inference Performance and Adversarial Robustness

Attack path identification accuracy measures the proportion of ground-truth attack chains correctly recovered by the inference engine. The ground-truth labels derive from expert annotation of dataset attack scenarios mapping observed malicious activity sequences to ATT&CK technique chains. The strict matching criterion requires exact technique sequence correspondence while the relaxed criterion permits partial credit for subsequence overlap.

**Figure 3:** Attack Path Identification Performance Comparison



This figure presents a grouped bar chart comparing attack path identification accuracy across methods and datasets. The horizontal axis enumerates evaluation datasets (IEC 60870-5-104, ICS Anomaly, Gas Pipeline, Water Treatment). The vertical axis displays accuracy percentage from 0 to 100. Bar groups represent methods including the proposed approach, GraphSAGE baseline, GAT baseline, rule-based baseline, and deep learning baseline. Error bars indicate 95% confidence intervals computed from cross-validation folds. The proposed approach achieves highest accuracy across all datasets with 94.7% mean accuracy. Annotation text boxes highlight key performance differentials. Color scheme employs colorblind-safe palette with distinct hues for each method.

The proposed framework achieves 94.7% strict accuracy and 97.2% relaxed accuracy averaged across datasets, representing 8.3% and 4.1% improvements over the strongest baseline respectively. Performance advantages concentrate in scenarios involving multi-stage attacks spanning protocol boundaries where knowledge graph relationships capture cross-layer dependencies invisible to single-protocol analysis approaches. The statistical significance of improvements is confirmed through paired t-tests yielding p-values below 0.01 for all dataset comparisons.

Per-attack-type analysis reveals systematic performance variations correlated with attack complexity. Reconnaissance attacks achieve 98.1% identification accuracy reflecting distinctive network scanning patterns readily distinguishable from normal traffic. Command injection attacks achieve 93.4% accuracy with failure cases attributable to legitimate operator commands exhibiting similar syntactic structure. Denial of service attacks achieve 96.8% accuracy benefiting from volumetric anomaly signals complementing semantic analysis. Man-in-the-middle attacks present intermediate difficulty at 91.2% accuracy requiring correlation across multiple network segments for complete path reconstruction.

Adversarial robustness evaluation assesses framework resilience against deliberately crafted evasion attempts. The threat model assumes adversaries with knowledge of detection features but without direct access to modify trained models. Attack generation employs projected gradient descent perturbations constrained within L-infinity norm bounds preserving traffic functionality. The IDSGAN framework provides generative adversarial baseline producing synthetic malicious traffic designed to evade detection.

The robustness evaluation reveals 12.8% improvement in accuracy under adversarial conditions compared to baseline graph neural network approaches. The knowledge graph constraint enforcement prevents acceptance of attack paths violating physical connectivity or protocol compatibility requirements that purely learned representations may overlook. Adversarial examples successfully evading neural components still fail validation against symbolic constraints reducing effective evasion rate.

**Table 7:** Adversarial Robustness Performance Under Perturbation

| Perturbation Budget (ε) | Proposed Accuracy | GraphSAGE Accuracy | GAT Accuracy | Improvement |
|---|---|---|---|---|
| 0.00 (Clean) | 94.7% | 86.4% | 87.9% | +7.8% |
| 0.01 | 93.2% | 81.2% | 83.1% | +10.8% |
| 0.05 | 91.4% | 76.8% | 78.4% | +13.6% |
| 0.10 | 88.9% | 71.3% | 73.2% | +16.2% |
| 0.20 | 84.6% | 64.7% | 66.9% | +18.3% |

The computational efficiency analysis demonstrates real-time inference feasibility for operational deployment. Single attack path inference completes within 847 milliseconds average latency. Batch processing of 1,000 concurrent queries achieves 12,340 paths per second throughput. Memory consumption peaks at 48GB during knowledge graph loading with stable 31GB utilization during inference.

The scalability assessment extrapolates performance to enterprise-scale deployments. Graph size experiments vary node counts from 10,000 to 1,000,000 entities measuring inference latency scaling. The observed O(n log n) complexity derives from indexed traversal avoiding exhaustive enumeration. Projected latency for 1,000,000 node graphs remains under 3 seconds meeting operational response requirements.

### 4.3 Explainability Assessment and Regulatory Compliance Verification

Explainability evaluation employs both automated metrics and human expert assessment capturing distinct quality dimensions. Faithfulness measures correspondence between generated explanations and actual model decision factors. Completeness quantifies coverage of relevant contributing elements within explanations. Comprehensibility assesses human understandability through expert rating surveys.

The faithfulness metric implementation adapts the comprehensiveness and sufficiency framework from explainable AI literature. Comprehensiveness removes explanation-highlighted features and measures prediction change magnitude. Sufficiency retains only highlighted features and measures prediction preservation. High comprehensiveness indicates explanations identify causally relevant factors while high sufficiency confirms explanation completeness.

$$\text{Faithfulness}(E) = 0.5 \cdot \text{Comprehensiveness}(E) + 0.5 \cdot \text{Sufficiency}(E)$$

Expert evaluation engaged twelve cybersecurity professionals with average 8.3 years experience in industrial control system security roles. Participants reviewed 200 randomly sampled attack path explanations rating comprehensibility on five-point Likert scales. The rating instrument assessed clarity of technical descriptions, logical coherence of attack progression narratives, and actionability of recommended mitigations.

The proposed hybrid explanation approach achieves 4.5/5.0 average comprehensibility rating with 0.89 completeness score and 0.92 faithfulness score. Qualitative feedback highlights appreciation for structured format enabling rapid scanning alongside natural language elaboration providing context for unfamiliar attack techniques. Criticism concentrates on occasional verbosity in LLM-generated segments that could be condensed without information loss.

Regulatory compliance verification maps framework capabilities against specific requirements from NERC CIP and NIS2 directive provisions. The assessment employs a compliance matrix documenting requirement identifiers, requirement descriptions, framework features addressing each requirement, and evidence demonstrating satisfaction.

**Table 8:** Regulatory Compliance Mapping Summary

| Regulation | Requirement | Framework Feature | Compliance Status |
|---|---|---|---|
| NERC CIP-005-7 R1 | Electronic Security Perimeter | Network topology graph | Satisfied |

| NERC CIP-007-6 R4 | Security Event Monitoring | Real-time inference engine | Satisfied |
| NERC CIP-008-6 R1 | Incident Response | Automated alert generation | Satisfied |
| NIS2 Art. 21.2(d) | Security Monitoring | Continuous analysis pipeline | Satisfied |
| NIS2 Art. 21.2(g) | Audit Logging | Explanation archive storage | Satisfied |
| NIST 800-82 6.2.7 | Audit Trail | Path inference documentation | Satisfied |

The audit trail generation capability produces regulatory-ready documentation automatically. Each detected attack path triggers creation of structured incident record containing timestamp, affected assets, inferred technique sequence, supporting evidence references, and recommended response actions. The documentation format aligns with ICS-CERT incident reporting templates facilitating information sharing with sector coordination bodies.

Comparative analysis against alternative explainability approaches positions the proposed framework within the broader landscape. LIME-based explanations achieve 0.81 faithfulness but require repeated model queries increasing computational overhead. SHAP-based explanations provide theoretically grounded feature attribution but struggle with graph-structured inputs requiring custom kernel definitions. Attention-based explanations from graph neural networks offer built-in interpretability but conflate multiple attention heads into potentially inconsistent narratives. The proposed hybrid approach balances these tradeoffs achieving superior faithfulness-completeness-efficiency characteristics.

Ablation studies quantify contributions of individual framework components. Removing knowledge graph constraints reduces accuracy by 11.3% while decreasing adversarial robustness by 18.7%. Disabling LLM semantic analysis reduces entity extraction F1 by 15.2% on proprietary log formats while maintaining performance on standardized protocols. Eliminating graph embedding reduces link prediction mean reciprocal rank from 0.847 to 0.712 degrading attack path inference quality.

The operational deployment pilot engaged three utility organizations implementing framework instances within security operations center environments over twelve-week evaluation periods. Security analysts processed 4,847 framework-generated alerts with 89.3% rated as actionable providing useful investigation guidance. The mean time to initial triage reduced by 41% compared to baseline alert processing workflows relying on manual log review. False positive rates averaged 6.2% within acceptable ranges for production deployment. Analyst feedback highlighted particular value in multi-stage attack correlation capabilities connecting disparate alert streams into coherent threat narratives.

# 5. Conclusion and Future Work

## 5.1 Summary of Research Findings and Contributions

This research develops an explainable attack path reasoning framework integrating knowledge graph construction with large language model-assisted semantic analysis for industrial control network security. The domain-specific ontology captures ICS assets, vulnerabilities, protocols, and attack techniques with formal relationships enabling logical inference aligned with MITRE ATT&CK for ICS. The hybrid reasoning architecture combines symbolic graph traversal with neural semantic analysis achieving 94.7% attack path identification accuracy while generating human-interpretable explanations.

The knowledge graph construction methodology transforms heterogeneous security data sources into unified representations supporting multi-hop inference. Entity resolution algorithms achieve consistent identification across vulnerability databases and threat intelligence feeds. The LLM-assisted log analysis component addresses protocol heterogeneity through retrieval-augmented semantic parsing achieving 0.92 entity extraction F1 score. Adversarial robustness mechanisms verify outputs against deterministic parsers.

The explainability assessment demonstrates 4.5/5.0 expert comprehensibility ratings with regulatory compliance mapping confirming NERC CIP, NIS2, and NIST 800-82 satisfaction. Operational deployment pilots validate 41% reduction in alert triage time with 89.3% analyst satisfaction scores.

## 5.2 Limitations and Scope Constraints

The current framework implementation exhibits limitations warranting acknowledgment. Knowledge graph completeness depends on available vulnerability disclosure and threat intelligence sources. Undisclosed zero-

day vulnerabilities and novel attack techniques remain outside inference scope until knowledge base updates incorporate relevant information. The manual curation burden for maintaining current knowledge bases presents operational sustainability challenges.

The LLM component introduces dependencies on external model providers with associated availability, latency, and confidentiality considerations. On-premise deployment of open-source models addresses confidentiality concerns but requires substantial computational infrastructure investment. Model updates may alter semantic analysis behavior requiring revalidation of downstream inference accuracy.

Evaluation datasets derive from testbed environments that may not capture full complexity of production industrial control systems. Transfer learning assessments to operational deployments reveal performance degradation attributable to distribution shift between laboratory and field conditions.

## References

[1]. E. Karlsen, X. Luo, N. Zincir-Heywood, and M. Heywood, "Benchmarking large language models for log analysis, security, and interpretation," J. Netw. Syst. Manage., vol. 32, no. 3, article 59, Jun. 2024.

[2]. T. Al Ghazo and R. Kumar, "Identification of critical-attacks set in an attack-graph," in Proc. IEEE 10th Annu. Ubiquitous Comput., Electron. Mobile Commun. Conf. (UEMCON), New York, NY, USA, 2019, pp. 716-722.

[3]. Dehlaghi-Ghadim et al., "Anomaly detection dataset for industrial control systems," IEEE Access, vol. 11, pp. 107298-107316, 2023.

[4]. M. A. Teixeira et al., "Securing SCADA systems: A protocol-based intrusion detection approach with Shapley analysis," in Proc. IEEE Int. Conf. Consum. Electron. (ICCE), 2024, pp. 1-6.

[5]. Y. Ren, Y. Xiao, Y. Zhou, Z. Zhang, and Z. Tian, "CSKG4APT: A cybersecurity knowledge graph for advanced persistent threat organization attribution," IEEE Trans. Knowl. Data Eng., vol. 35, no. 6, pp. 5695-5709, Jun. 2023.

[6]. MITRE, "ATT&CK for industrial control systems: Design and philosophy," MITRE, Tech. Rep., Mar. 2020.

[7]. M. Zolanvari et al., "A novel dataset for experimentation with intrusion detection systems in SCADA networks using IEC 60870-5-104 standard," IEEE Access, vol. 12, 2024.

[8]. T. Al Ghazo, M. Ibrahim, H. Ren, and R. Kumar, "A2G2V: Automatic attack graph generation and visualization and its applications to computer and SCADA networks," IEEE Trans. Syst., Man, Cybern., Syst., vol. 50, no. 10, pp. 3488-3498, Oct. 2020.

[9]. Sagheer et al., "Explainable SCADA-edge network intrusion detection system: Tree-LIME approach," in Proc. IEEE Innov. Smart Grid Technol. Asia (ISGT Asia), 2023, pp. 1-5.

[10]. K. Stouffer, V. Pillitteri, S. Lightman, M. Abrams, and A. Hahn, "Guide to industrial control systems (ICS) security," NIST Special Publication 800-82, Revision 2, May 2015.

[11]. M. Zolanvari et al., "Deep-learning-based network intrusion detection for SCADA systems," in Proc. IEEE Conf. Commun. Netw. Security (CNS), Washington, DC, USA, 2019, pp. 1-7.

[12]. "Intrusion detection in industrial control systems based on deep reinforcement learning," IEEE Access, vol. 12, 2024.

[13]. S. Lin et al., "IDSGAN: Generative adversarial networks for attack generation against intrusion detection," in Advances in Knowledge Discovery and Data Mining (PAKDD 2022), Lecture Notes in Computer Science, vol. 13282. Cham, Switzerland: Springer, 2022, pp. 79-91.

[14]. "CLogLLM: A large language model enabled approach to cybersecurity log anomaly analysis," in Proc. IEEE Int. Conf. Consum. Electron. (ICCE), 2024.