# ASTRIDE: A Security Threat Modeling Platform for Agentic-AI Applications

Eranga Bandara[a], Amin Hass[b], Ross Gore[a], Sachin Shetty[a], Ravi Mukkamala[a], Safdar H. Bouk[a], Xueping Liang[c], Ng Wee Keong[d], Kasun De Zoysa[e], Aruna Withanage[f], Nilaan Loganathan[f]

[a]*Old Dominion University, Norfolk, VA, USA*
[b]*Accenture Technology Labs, Arlington, VA, USA*
[c]*Florida International University, USA*
[d]*Nanyang Technological University, Singapore*
[e]*University of Colombo, Sri Lanka*
[f]*Effectz.AI*

## Abstract

AI agent-based systems are becoming increasingly integral to modern software architectures, enabling autonomous decision-making, dynamic task execution, and multimodal interactions through large language models (LLMs). However, these systems introduce novel and evolving security challenges, including prompt injection attacks, context poisoning, model manipulation, and opaque agent-to-agent communication, that are not effectively captured by traditional threat modeling frameworks. In this paper, we introduce ASTRIDE, an automated threat modeling platform purpose-built for AI agent-based systems. ASTRIDE extends the classical STRIDE framework by introducing a new threat category, *A for AI Agent-Specific Attacks*, which encompasses emerging vulnerabilities such as prompt injection, unsafe tool invocation, and reasoning subversion, unique to agent-based applications. To automate threat modeling, ASTRIDE combines a consortium of fine-tuned vision-language models (VLMs) with the OpenAI-gpt-oss reasoning LLM to

*Email addresses:* cmedawer@odu.edu (Eranga Bandara),
amin.hassanzadeh@accenture.com (Amin Hass), rgore@odu.edu (Ross Gore),
sshetty@odu.edu (Sachin Shetty), mukka@odu.edu (Ravi Mukkamala), sbouk@odu.edu
(Safdar H. Bouk), xuliang@fiu.edu (Xueping Liang), awkng@ntu.edu.sg (Ng Wee
Keong), kasun@ucsc.cmb.ac.lk (Kasun De Zoysa), aruna@effectz.ai (Aruna
Withanage), nilaan@effectz.ai (Nilaan Loganathan)

perform end-to-end analysis directly from visual agent architecture diagrams, such as data flow diagrams(DFDs). LLM agents orchestrate the end-to-end threat modeling automation process by coordinating interactions between the VLM consortium and the reasoning LLM. Our evaluations demonstrate that ASTRIDE provides accurate, scalable, and explainable threat modeling for next-generation intelligent systems. To the best of our knowledge, ASTRIDE is the first framework to both extend STRIDE with AI-specific threats and integrate fine-tuned VLMs with a reasoning LLM to fully automate diagram-driven threat modeling in AI agent-based applications.

*Keywords:* LLM-Reasoning, Vision-Language-Model, Llama-Vision, OpenAI-gpt-oss, Threat-Modeling, STRIDE

## 1. Introduction

The rapid proliferation of AI agent-based systems—autonomous software entities empowered by large language models (LLMs) and capable of dynamic task execution, multimodal perception, and self-directed reasoning—has fundamentally transformed modern application architectures [1, 2]. These agents are increasingly embedded in workflows ranging from intelligent assistants to autonomous decision-making systems and multi-agent environments. However, their emerging behaviors and the reliance on complex language-driven interactions introduce novel security challenges that extend beyond the scope of traditional software threats. These include prompt injection attacks, context poisoning, reasoning subversion, and unintended tool invocation, all of which exploit the flexible and opaque nature of agent decision pipelines [3]. Existing threat modeling methodologies, including the widely adopted STRIDE framework, are not well-equipped to capture the emerging attack surfaces unique to agentic AI applications [4, 5]. These frameworks rely heavily on manual inspection of system architectures, demand significant domain expertise, and lack native support for AI-specific vulnerabilities such as prompt manipulation, context or memory poisoning, and inter-agent influence. Consequently, critical threats may go undetected until late in the development lifecycle, posing significant risks to the integrity, availability, and trustworthiness of AI-driven systems [6].

To address this gap, we introduce ASTRIDE, an automated threat modeling platform specifically designed for AI agent-based applications. ASTRIDE extends the classical STRIDE framework by introducing a new threat cate-

gory *A for AI Agent–Specific Attacks* which captures emerging security risks unique to agentic workflows. These include instruction manipulation (e.g., prompt injection), unsafe reasoning-driven tool use, and the misuse of agent memory or context windows. ASTRIDE combines a consortium of fine-tuned vision-language models (VLMs) [7] with the OpenAI-gpt-oss reasoning LLM [8] to perform end-to-end threat analysis directly from visual architecture diagrams, such as data flow diagrams. Each VLM is trained on a custom dataset of AI agent topologies annotated with both traditional STRIDE vectors and AI-specific threats, enabling robust component-level vulnerability detection. The VLM outputs are aggregated and synthesized by OpenAI-gpt-oss, which performs high-level reasoning to generate a coherent and explainable threat model. LLM agents [1] orchestrate the end-to-end threat modeling automation process by coordinating interactions between the VLM consortium and the reasoning LLM. To ensure practical deployment, all models are fine-tuned using the Unsloth library with QLoRA-based quantization, enabling low-latency inference on resource-constrained hardware [9].

Our evaluation shows that ASTRIDE provides an accurate, scalable, and interpretable solution for threat modeling in next-generation AI systems. To our knowledge, this is the first framework to both extend STRIDE with AI-specific threats and integrate fine-tuned VLMs with a reasoning LLM to fully automate visual, diagram-driven threat modeling for AI agent-based applications. The following are our main contributions of this research.

1. Propose a novel threat modeling framework for capturing the attack surfaces of agentic AI applications by extending the traditional STRIDE methodology.
2. Automate the end-to-end threat modeling process for agentic AI applications using a consortium of fine-tuned vision-language models and the OpenAI-gpt-oss reasoning model.
3. Utilize a fine-tuned vision-language model consortium to analyze threat model images (e.g., data flow diagrams) and generate ASTRIDE-based threat predictions for agentic AI applications.
4. Incorporate the OpenAI-gpt-oss reasoning model to synthesize and finalize threat assessments based on the outputs of the vision-language model consortium.

The rest of the paper is organized as follows: Section 2 presents related work and positions our approach in the context of existing AI-enabled threat

modeling frameworks. Section 3 describes the overall system architecture of the proposed platform. Section 4 outlines the core functionalities and operational workflow. Section 5 provides details on the implementation and performance evaluation. Finally, Section 6 concludes the paper and discusses directions for future research.

## 2. Related Work

Several researchers have explored AI-based approaches to threat modeling, each addressing different aspects of securing intelligent systems. Mauri and Damiani [10] extended STRIDE for AI/ML lifecycles, emphasizing tailored security assessments. The ADMIn framework [11] by Kumar et al. categorizes threats to AI software into attacks on datasets, models, and inputs. Yang et al.[12] proposed using LLMs to automate threat modeling in banking systems, while Mollaeefar et al.[13] introduced PILLAR, which integrates LLMs with the LINDDUN framework to automate privacy threat classification. In healthcare, STRIDE was adapted for LLM-based systems to model component-level threats [14]. Auspex [15] introduced 'tradecraft prompting' to embed threat modeling expertise into generative AI systems. Lastly, STRIDE-AI [16] by Zhou et al. offers an asset-centric model for identifying threats in ML components like datasets and training environments.

Table 1 provides a comparative analysis of existing frameworks across key dimensions, including support for fine-tuning on structured diagram datasets, runtime integration of LLMs and VLMs, vision-language threat modeling capabilities, reasoning LLM usage, and modular consortium-based orchestration of multiple specialized models. Most prior work lacks vision-language understanding of system diagrams or treats threat modeling as static text generation, offering limited ability to capture agent-specific attack surfaces such as prompt injection, context poisoning, or unsafe tool invocation addressed in ASTRIDE.

In contrast, our proposed ASTRIDE platform extends STRIDE with an additional category for AI Agent–Specific Attacks and uniquely combines fine-tuned vision-language models and the OpenAI-gpt-oss reasoning LLM to automate threat analysis directly from system diagrams. ASTRIDE supports fine-tuning on structured diagram datasets, ensures interoperability across multiple models, and introduces decision-level synthesis to resolve conflicting outputs. By integrating architectural understanding, multimodal reasoning , and safeguards against agent-specific risks such as prompt injec-

4

tion and unsafe tool invocation [17, 1], ASTRIDE advances the state of the art in automated threat modeling for agentic AI applications.

Table 1: AI-enabled Threat Modeling Framework Comparison

| Platform | Domain | Fine-tuning Support | Running LLM | Vision LM Support | Reasoning LLM Support | LLM Consortium Support |
|---|---|---|---|---|---|---|
| Astride | Threat Modeling | ✓ | Llama-3, Pixtrel, Qwen2 OpenAI-gpt-oss | ✓ | ✓ | ✓ |
| STRIDE for AI/ML [10] | AI/ML Systems | ✗ | ✗ | ✗ | ✗ | ✗ |
| ADMIn [11] | AI-based Software | ✗ | ✗ | ✗ | ✗ | ✗ |
| ThreatModeling-LLM [12] | Banking Systems | ✗ | GPT-based | ✗ | ✗ | ✗ |
| PILLAR [13] | Privacy (LINDDUN) | ✗ | GPT-4 | ✗ | ✗ | ✗ |
| LLM-based Healthcare STRIDE [14] | Healthcare Systems | ✗ | GPT-4 | ✗ | ✗ | ✗ |
| Auspex [15] | General Threat Modeling | ✗ | GPT-4 | ✗ | ✗ | ✗ |
| STRIDE-AI [16] | ML System Assets | ✗ | ✗ | ✗ | ✗ | ✗ |

## 3. System Architecture

Figure 1 illustrates the overall architecture of the proposed platform, which is composed of four core components: (1) Data Lake, (2) Vision-Language Models, (3) OpenAI-gpt-oss Reasoning LLM, and (4) LLM Agents. A brief description of each component is provided below.

### 3.1. Data Lake

Data Lake serves as the foundational infrastructure for the ASTRIDE platform, providing centralized storage, management, and access to the diverse and large-scale datasets required for training and fine-tuning VLM and reasoning LLMs for automated threat modeling [7]. Specifically, this layer aggregates a curated corpus of labeled threat modeling diagrams, including Data Flow Diagrams associated with threat vectors and trust boundaries. These datasets are critical for teaching VLMs to understand the structural and semantic patterns in system diagrams and to associate them with potential security threats.

### 3.2. Vision-Language Models

VLMs form the analytical layer of the ASTRIDE platform, enabling it to interpret complex threat modeling diagrams and autonomously extract threat insights. This layer consists of a consortium of fine-tuned VLMs, each trained on a large, labeled dataset of system architecture and data flow diagrams annotated with associated threat vectors, trust boundaries, and component-level metadata [7, 18]. Rather than relying on a single model,

the ASTRIDE platform employs a multi-model ensemble strategy, where each VLM independently analyzes input diagrams and outputs structured threat observations. These include AI agent-specific attacks, as well as potential STRIDE attacks spoofing, tampering, repudiation, information disclosure, denial of service, and privilege escalation threats - each mapped to specific components and data flows in the diagram. The use of diverse VLMs enhances resilience against individual model biases and ensures broader semantic coverage across a range of threat modeling styles and visual representations.
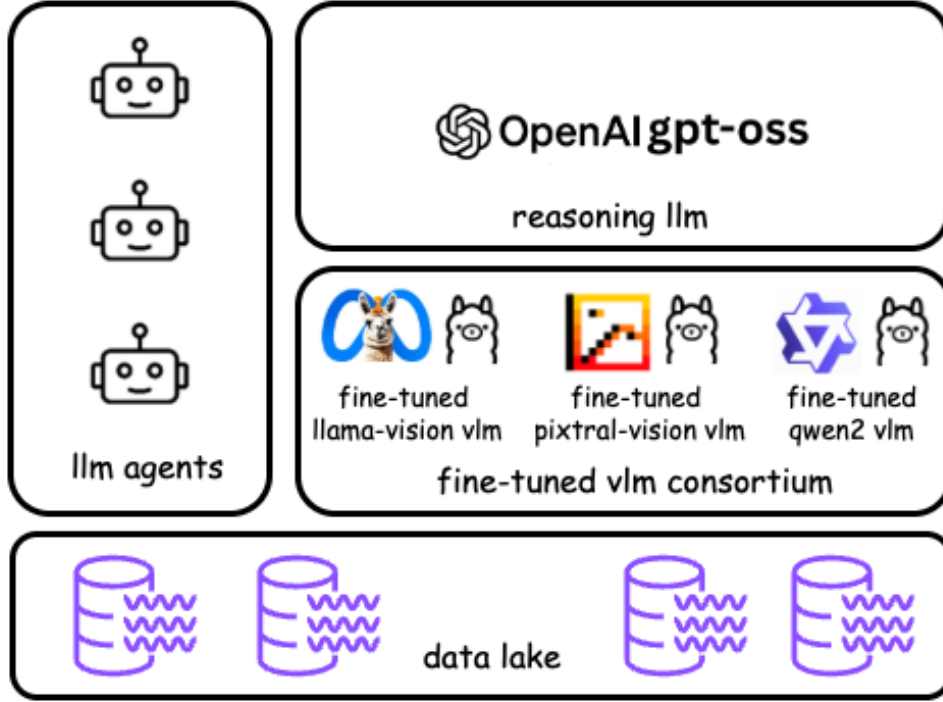


Figure 1: ASTRIDE system architecture.

### 3.3. OpenAI-gpt-oss Reasoning LLM

The OpenAI-gpt-oss Reasoning LLM [19, 8] represents the cognitive layer of the ASTRIDE platform, responsible for high-level reasoning, synthesis, and refinement of threat predictions derived from complex system diagrams. In the context of ASTRIDE, this layer performs the final aggregation and decision-making function in the Agentic-AI application threat modeling

pipeline. Predictions generated by the VLM consortium, each VLM offering a unique perspective on AI-Agent specific attacks and STRIDE-based attacks(spoofing, tampering, repudiation, information disclosure, denial of service, and elevation of privilege threats), are collected and structured by the LLM Agent Layer into a composite prompt [3, 20]. This prompt is then submitted to OpenAI-gpt-oss, which uses its reasoning abilities to synthesize a cohesive, prioritized, and contextually validated threat model.

By cross-analyzing the output of multiple VLMs, OpenAI-gpt-oss ensures that conflicting or incomplete assessments are reconciled into a consistent and trustworthy threat modeling representation. It also applies advanced reasoning to contextual cues (e.g., system roles, trust boundaries, and domain-specific risk patterns), enabling deeper insights than purely visual or rule-based systems can achieve [19, 21].
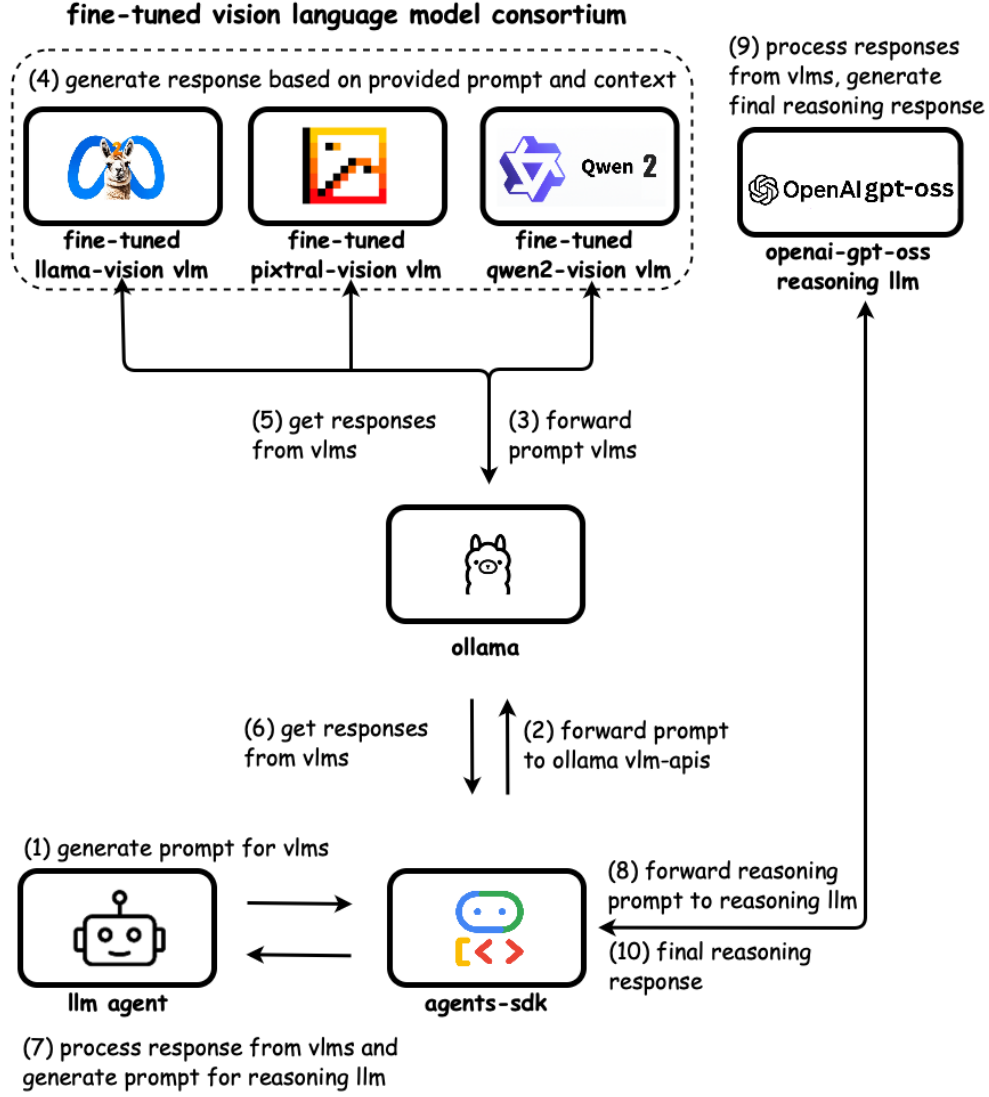
**fine-tuned vision language model consortium**

(4) generate response based on provided prompt and context

**fine-tuned llama-vision vlm**

**fine-tuned pixtral-vision vlm**

Qwen **2**

**fine-tuned qwen2-vision vlm**

(9) process responses from vlms, generate final reasoning response

OpenAI **gpt-oss**

**openai-gpt-oss reasoning llm**

(5) get responses from vlms

(3) forward prompt vlms

**ollama**

(6) get responses from vlms

(2) forward prompt to ollama vlm-apis

(1) generate prompt for vlms

**llm agent**

(8) forward reasoning prompt to reasoning llm

(10) final reasoning response

**agents-sdk**

(7) process response from vlms and generate prompt for reasoning llm

Figure 2: ASTRIDE threat prediction flow

*3.4. AI/LLM Agents*

AI/LLM Agents [1] serve as the orchestration layer of the ASTRIDE platform, managing coordination between the VLM consortium and the OpenAI-gpt-oss reasoning LLM [19, 22]. The agents generate customized prompts tailored to extract both STRIDE and Agentic-AI-specific threat vectors from the fine-tuned VLMs. Each VLM independently analyzes the input and re-

turns structured threat insights, as illustrated in Figure 2. Smart contracts then aggregate and reformat the VLM predictions into a unified reasoning prompt for OpenAI-gpt-oss, which synthesizes the output into a coherent and context-sensitive final threat model [8, 23].

## 4. Platform Functionality

There are four main functionalities of the platform: 1) Data Lake Setup, 2) Vision language model Fine-Tuning, 3) Threat prediction by Fine-tuned VLMs, and 4) Final Prediction by OpenAI-gpt-oss LLM. This section goes into the specifics of these functions.

### 4.1. Data Lake Setup

The Data Lake serves as the foundational layer of the ASTRIDE platform, designed to store, organize, and manage the comprehensive datasets required for fine-tuning VLMs on threat modeling tasks. Unlike traditional datasets focused on image classification or natural language, the ASTRIDE Data Lake contains visual system artifacts such as data flow diagrams, system architecture diagrams, and associated textual annotations, including labeled threat vectors, component descriptions, trust boundary definitions, and mitigation factors, which streamlines the fine-tuning process [24, 25].

### 4.2. Vision-Language Model Fine-Tuning

The next phase in the ASTRIDE workflow is the fine-tuning of VLMs using the curated, labeled dataset of threat modeling diagrams and threat annotations stored in the Data Lake. This step involves adapting general-purpose VLMs(e.g. *Llama-Vision, Pix2Struct, and Qwen2-VL* [26, 27]) to the specific domain of cybersecurity threat modeling. These models are fine-tuned to accurately associate system architecture and data flow diagram features with corresponding threat vectors, trust boundaries, system component roles, and mitigation. Through this fine-tuning pipeline, ASTRIDE transforms general VLMs into domain-specialized threat prediction agents capable of extracting actionable threats from real-world agentic-ai system diagrams with high precision and consistency, as illustrated in Figure 3.
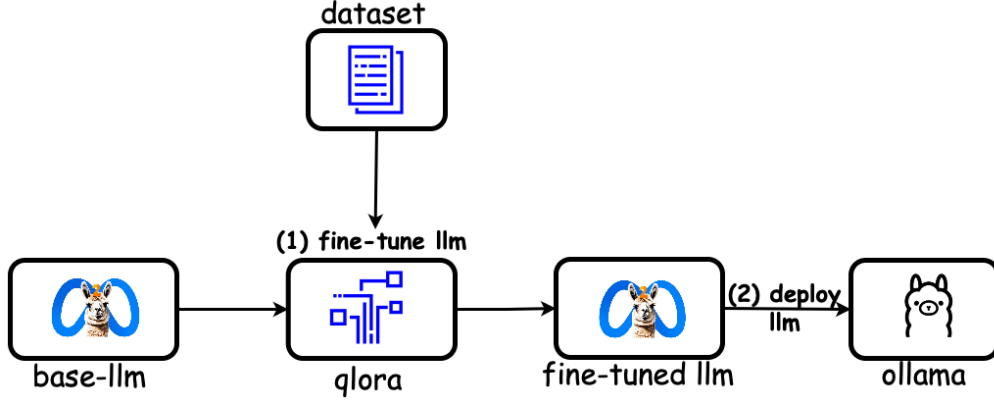
Figure 3: Fine-tune VLMs with Qlora and deploy with Ollama.

### 4.3. Threat Prediction by Fine-Tuned VLMs

Once the VLMs are fine-tuned using labeled threat modeling data, the next key functionality of the ASTRIDE platform is to extract threat predictions from new system diagrams using the VLM consortium. When a new data flow diagram or system architecture image is submitted, the LLM agent layer orchestrates the interaction with each fine-tuned VLM and initiates a distributed threat analysis process.

To ensure accurate and context-aware predictions, it employs custom prompt engineering, crafting model-specific instructions that inject system metadata, domain context, and task objectives alongside visual input []. These prompts instruct each VLM to analyze the image and extract threat vectors(e.g., AI-Agents specific threats and STRIDE threats) and to associate them with relevant system components and data flows [4].

Each model in the consortium independently processes the input diagram and produces a structured set of predictions, including identified threat types, associated components, severity levels, and recommended mitigations. These outputs are returned to the LLM agent layer, which collects, formats and prepares them for further synthesis using the OpenAI-gpt-oss Reasoning Layer [8].

### 4.4. Final Prediction and Reasoning by OpenAI-gpt-oss LLM

To ensure the precision, consistency, and reliability of the threat models generated, ASTRIDE employs a consensus-based reasoning mechanism powered by the OpenAI-gpt-oss LLM. Rather than relying on the output

10

of a single VLM, the platform aggregates threat predictions from multiple independently fine-tuned VLMs. These intermediate outputs, each offering different perspectives on the same visual input, are then synthesized by OpenAI-gpt-oss, which serves as the final reasoning and decision layer within the platform [8].

As a reasoning-centric LLM, OpenAI-gpt-oss is uniquely suited to interpret and compare the predictions from different VLMs. Contextual analysis is performed to validate, reconcile, and rank threats, ultimately generating a unified threat model that reflects the most accurate and complete understanding of the system diagram. To guide this reasoning process, LLM agents construct structured prompts that encapsulate each VLM's output, including identified threat vectors, affected components, and proposed mitigation strategies. An example of this prompt structure is shown in Figure 4.

```
1   prompt = ChatPromptTemplate(
2 ▾     messages=[
3           SystemMessagePromptTemplate.from_template(
4               "You are an expert in ASTRIDE-based threat modeling, which extends STRIDE to
                      include AI Agent-Specific Attacks. "
5               "You will be provided with threat predictions generated by multiple fine-tuned
                      vision-language models, "
6               "each analyzing a threat modeling diagram for an agentic AI application. These
                      diagrams include trust boundaries, memory stores, tool invocation modules,
                      and agent-to-agent communication channels. "
7               "Your task is to evaluate these predictions and reason about the most accurate,
                      comprehensive, and contextually appropriate ASTRIDE threat model. "
8               "In addition to traditional STRIDE categories (Spoofing, Tampering, Repudiation,
                      Information Disclosure, Denial of Service, Elevation of Privilege), "
9               "also assess for agent-specific threats such as Prompt Injection, Context
                      Poisoning, Reasoning Subversion, and Unsafe Tool Invocation. "
10              "Your goal is to synthesize a final structured threat model that reflects both
                      consistency and completeness across the predictions."
11          ),
12          HumanMessagePromptTemplate.from_template(
13              "Here are the ASTRIDE threat predictions for the given system diagram:\n"
14              "- **Pix2Struct-12B model**: {prediction_1}\n"
15              "- **Qwen2-VL-7B-Instruct model**: {prediction_2}\n"
16              "- **LLaVA-Vision-Instruct model**: {prediction_3}\n\n"
17              "Based on these predictions, please synthesize the most accurate final ASTRIDE
                      threat model. "
18              "For each identified threat, specify:\n"
19              "- The **ASTRIDE category** (e.g., Spoofing, Prompt Injection, etc.)\n"
20              "- The **affected component**\n"
21              "- A brief **threat description**\n"
22              "- Recommended **mitigation strategy**"
23          )
24      ]
25  )
```

Figure 4: Prompt for OpenAI-gpt-oss reasoning LLM.

## 5. Implementation and Evaluation

A functional prototype of the ASTRIDE platform has been implemented to validate the feasibility and performance of automated threat modeling from visual system representations. The core analytical layer is composed of three fine-tuned VLMs—*Llama-Vision, Pix2Struct, and Qwen2-VL* [26]—each trained to interpret architectural diagrams and extract ASTRIDE-based security insights. The fine-tuning process utilized a synthetically generated dataset comprising system architecture and data flow diagrams annotated with trust boundaries, threat vectors, and mitigation strategies. Each diagram in the dataset was paired with a structured JSON representation containing labeled components, identified threats (e.g.,

AI-Agents specific threats as well as STRIDE threats), and corresponding mitigation techniques.

Fine-tuning was performed using the `Unsloth` library on Google Colab, leveraging NVIDIA A100 and Tesla TPU resources to enable efficient training cycles. The dataset consisted of approximately 1,200 annotated records, automatically generated in the form of Mermaid diagrams [28] representing system architectures (e.g., data flow diagrams, component diagrams, and trust-boundary layouts). Each entry was structured in a conversational instruction-tuning format, as required by Unsloth [9], and included a `content` field (containing the Mermaid diagram prompt), a `type` field (describing the diagram category), and an `instruction` field guiding the model to identify STRIDE threats and recommend corresponding mitigations. The dataset was split into 2/3 for training, 1/6 for validation, and 1/6 for testing. Training was completed in approximately 1,627 seconds (27.12 minutes). The peak reserved memory during training was 14.605 GB, with actual training consumption reaching 5.853 GB, representing 39.69% of total memory usage and 99.03% of peak allocation. These metrics demonstrate that structured threat modeling using visual input and LLM fine-tuning can be accomplished efficiently, even on moderate-scale datasets. Post fine-tuning, all models were optimized using `QLoRA` [29], allowing quantized deployment on edge or resource-constrained hardware. These VLMs were then hosted on the `Ollama` framework, which provides an efficient runtime environment for inference [30]. The reasoning layer of the proposed ASTRIDE platform is implemented using the OpenAI-gpt-oss LLM. The LLM agent layer is implemented with the OpenAI-Agents-SDK [1]. Platform performance is evaluated in two key areas: Evaluation of fine-tuned VLMs and Evaluation of OpenAI-gpt-oss Reasoning LLM.

## 5.1. Evaluation of Fine-tuned VLMs

In this evaluation, we first assessed the training and validation loss during the fine-tuning process of the VLMs for diagram-based threat modeling. These metrics, visualized in Figure 5, demonstrate the models' progressive learning across training steps. Furthermore, Figure 6 captures multiple key training dynamics, including the loss difference, loss ratio, and loss derivatives over training steps, offering valuable insights into the model's convergence behavior and generalization performance. The consistently positive loss difference (validation loss exceeding training loss) suggests signs of overfitting, especially at steps with noticeable spikes. The loss ratio, ranging

13

from 1.0 to 3.0, highlights varying degrees of generalization, where a lower ratio reflects better alignment between training and validation performance. Additionally, the loss derivatives reveal rapid initial improvements followed by smaller, oscillating changes, indicating stabilization or saturation in the learning process.
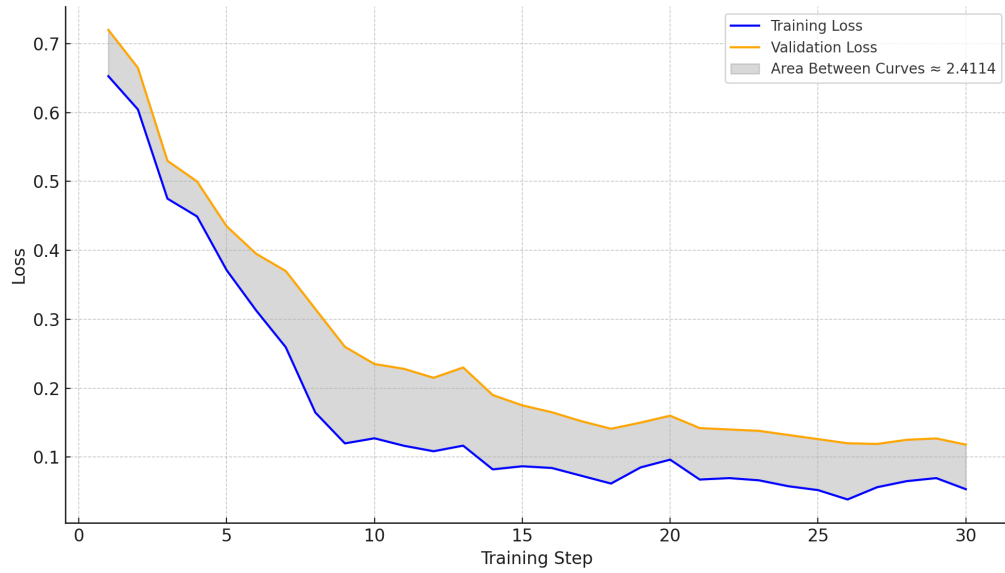


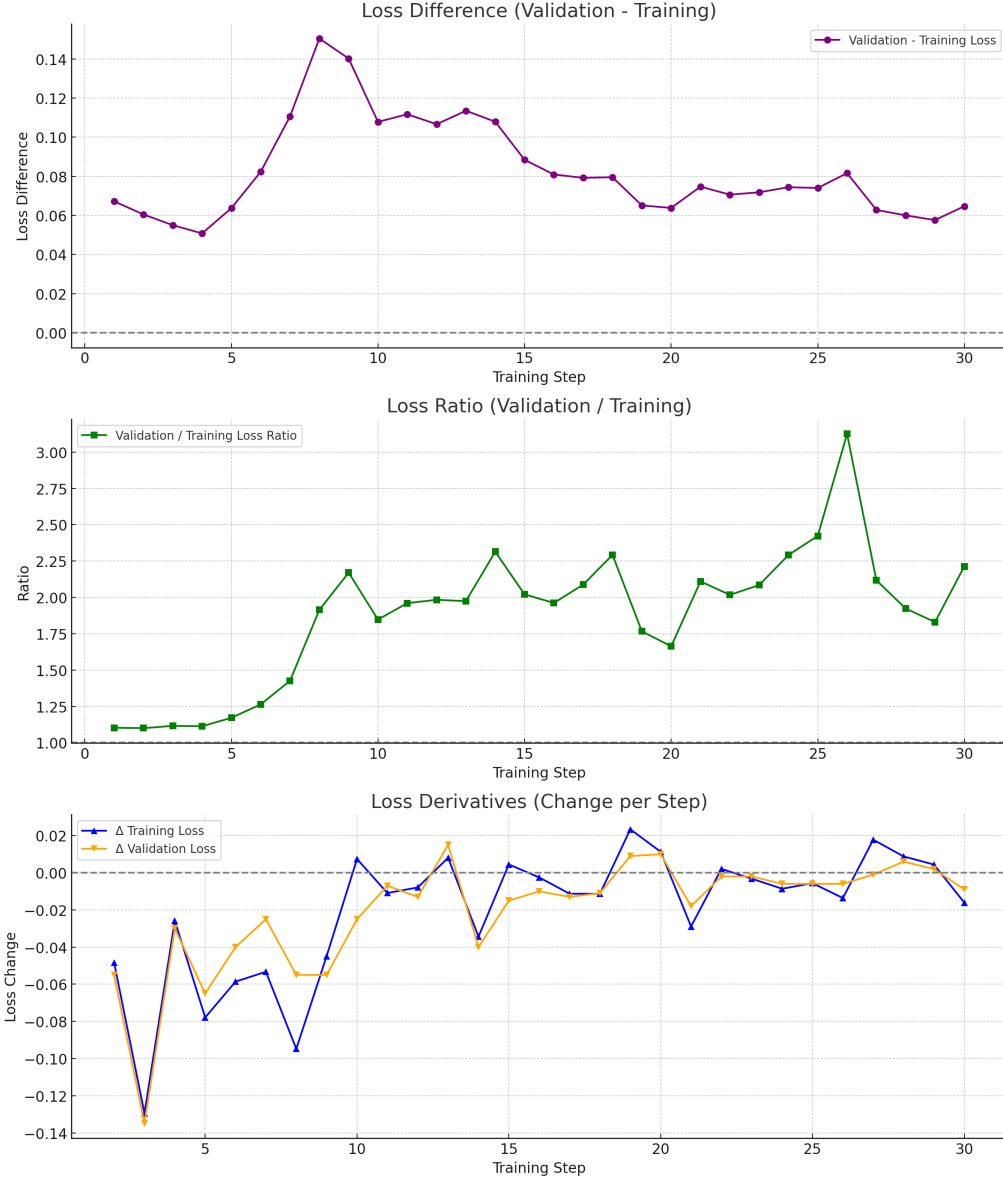Figure 5: Training loss and validation loss during fine-tuning of the Llama-3.2-11B-Vision-Instruct VLM.

Figure 6: Ratio of training to validation loss during the fine-tuning of the Llama-3.2-11B-Vision-Instruct VLM.

Figure 7 and Figure 8 present the predictions generated by two vision-language models, Llama-Vision and Pixtral-Vision [26], before and after fine-tuning for the given AI agent-based system architectures. For Llama-Vision,

the model initially detects only a single vulnerability—prompt injection at the prompt processor—without recognizing additional attack surfaces or providing detailed mitigations. After fine-tuning, it demonstrates substantial improvement by also identifying context poisoning in the reasoning core and unsafe tool invocation in the tool execution module. Furthermore, the model outputs context-aware mitigations, such as input sanitization, context integrity validation, and access control for API invocations. For Pixtral-Vision, the pre-fine-tuning prediction is similarly narrow, flagging only reasoning subversion at the intent and reasoning module. After fine-tuning, Pixtral-Vision expands its coverage to detect prompt injection in the NLU module, context poisoning in the contextual memory, and provides detailed, actionable mitigations—including prompt filtering with zero-trust validation, reasoning constraints with justification trails, and memory hashing with provenance tracking [17]. Taken together, these results demonstrate that fine-tuning significantly enhances the capabilities of VLMs to perform structured ASTRIDE-based threat modeling. Both models show improved coverage of AI-agent–specific threats, better interpretation of architectural components (e.g., reasoning core, memory store, external APIs), and stronger alignment with expert-annotated threat data. This highlights the effectiveness of fine-tuning in enabling multimodal models to generate context-aware, explainable, and actionable security insights for complex agentic AI systems.
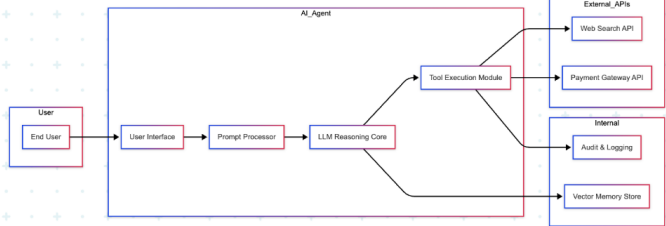
| Diagram |  |
| --- | --- |
| **Llama-Vision Prediction (Before fine-tuning)** | Prompt Injection<br>Component: Prompt Processor<br>Threat: Malicious input crafted by the user may inject unauthorized instructions that manipulate the agent's behavior. |
| **Llama-Vision Prediction (After fine-tuning)** | {"diagram":"AI Agent-Based Architecture","threats":[{"category":"A - AI Agent-Specific Attacks","component":"Prompt Processor","description":"Prompt injection attack could manipulate the agent's behavior by injecting unauthorized instructions through user inputs.","mitigation":"Implement prompt sanitization, input filtering, and context-aware validation before processing."},{"category":"A - AI Agent-Specific Attacks","component":"ReasoningCore","description":"Context poisoning can be introduced into the LLM's memory or prompt window, biasing decisions or enabling unauthorized tool use.","mitigation":"Enforce context integrity checks and isolate session memory."},{"category":"A - AI Agent-Specific Attacks","component":"ToolManager","description":"The LLM may invoke unauthorized tools or APIs based on manipulated reasoning, leading to misuse of external services.","mitigation":"Apply reasoning validation gates and enforce tool invocation access controls."}]} |

Figure 7: The prediction results of Llama-3.2-11B-Vision-Instruct VLM.

17

| | |
|---|---|
| **Diagram** |  |
| **Pixtral mlm prediction (Before fine-tuning)** | Reasoning Subversion<br>Component: Intent & Reasoning Module<br>Threat: The LLM may be manipulated through poisoned context or ambiguous input to make unsafe or biased decisions. |
| **Pixtral mlm prediction (After fine-tuning)** | {"diagram":"Voice-Based AI Agent System","threats":[{"category":"A - AI Agent-Specific Attacks","component":"NLU","description":"Prompt injection via voice commands can manipulate the LLM-based understanding to perform unintended or malicious actions.","mitigation":"Implement prompt filtering and intent validation, and apply a zero-trust policy for ambiguous commands."},{"category":"A - AI Agent-Specific Attacks","component":"DecisionModule","description":"Reasoning subversion can occur if context manipulation biases decision-making, potentially triggering unsafe tool use.","mitigation":"Add reasoning constraints, enforce contextual validation rules, and track agent decisions with justification trails."},{"category":"A - AI Agent-Specific Attacks","component":"ContextMem","description":"Context poisoning in long-term memory can corrupt future reasoning and lead to inappropriate behavior based on false context.","mitigation":"Validate all memory updates, apply memory hashing and provenance tracing."}]} |

Figure 8: The prediction results of Pixtral-Vision VLM.

## 5.2. Evaluation of OpenAI-gpt-oss Reasoning LLM

In this evaluation, we assessed the reasoning performance of the OpenAI-gpt-oss Reasoning LLM by comparing the individual predictions from the VLMs with the final prediction generated by OpenAI-gpt-oss. Figure 9 illustrates the predictions made by different VLMs for a threat modeling diagram, along with the final reasoning provided by OpenAI-gpt-oss. The results highlight OpenAI-gpt-oss's ability to synthesize and analyze the outputs from multiple VLMs, demonstrating its effectiveness and reliability in producing accurate final diagnoses. This evaluation underscores the reasoning LLM's role in refining and enhancing prediction precision, further validating its integration into the Astride framework.

| Threat diagram |  |
|---|---|
| Prediction of **Pixtral** | {"diagram":"Multi-Agent AI Collaboration Architecture","threats":[{"category":"A - AI Agent-Specific Attacks","component":"Planner","description":"Prompt injection from the human analyst may cause the planning agent to generate malicious or unintended task plans.","mitigation":"Apply prompt sanitization, intent validation, and policy-constrained planning logic."}]} |
| Prediction of **Qwen** | {"diagram":"Multi-Agent AI Collaboration Architecture","threats":[{"category":"A - AI Agent-Specific Attacks","component":"AgentMemory","description":"Context poisoning can occur if one agent writes misleading or malicious data to shared memory, impacting downstream reasoning.","mitigation":"Implement write access controls, data integrity validation, and traceable provenance for memory writes."}]} |
| Prediction of **Llama-vision** | {"diagram":"Multi-Agent AI Collaboration Architecture","threats":[{"category":"A - AI Agent-Specific Attacks","component":"Executor","description":"Unsafe tool invocation risk where the execution agent may call APIs without sufficient oversight due to reasoning subversion.","mitigation":"Enforce tool invocation policies, API usage whitelists, and reasoning validation filters."}]} |
| Reasoning of **OpenAI-o3** | Thought for 7 seconds<br>{"diagram":"Multi-Agent AI Collaboration Architecture","threats":[{"category":"Prompt Injection","component":"Planner","description":"A human analyst (or upstream agent) could craft prompts that trick the planning agent into producing malicious or unintended task plans.","mitigation":"Apply strict prompt sanitization and intent validation, constrain plan generation with guard-rail policies, and log all prompt inputs for audit."},{"category":"Context Poisoning","component":"AgentMemory","description":"An agent with write access could insert misleading or malicious data into shared memory, corrupting subsequent agents' reasoning.","mitigation":"Enforce fine-grained write permissions, validate data integrity and provenance on every write, and maintain immutable versioned snapshots of memory."},{"category":"Unsafe Tool Invocation","component":"Executor","description":"The execution agent may call external APIs or system commands without adequate oversight, leading to unintended side-effects or data leakage (reasoning subversion).","mitigation":"Require pre-execution policy checks, restrict calls to a vetted API whitelist, and run reasoning-validation filters before tool invocation."}]} |

Figure 9: Reasoning made by OpenAI-gpt-oss LLM.

# 6. Conclusions and Future Work

In this paper, we introduce ASTRIDE, a novel, fully automated threat modeling platform tailored for AI agent-based systems. ASTRIDE extends the traditional STRIDE framework by introducing a new threat category *A for AI Agent–Specific Attacks* to capture emerging security risks such as prompt injection, context poisoning, and unsafe tool invocation that are

unique to agentic workflows. We fine-tuned a consortium of VLMs, including Pix2Struct, Qwen2-VL, and Llama-Vision, using a custom dataset of annotated agent architecture diagrams containing trust boundaries, component interactions, data flows, and ASTRIDE-labeled threat vectors. These fine-tuned models are capable of interpreting complex visual representations and generating structured, component-level threat predictions. To enhance the reliability and transparency of the automation workflow, we employed LLM Agents to orchestrate secure and auditable interactions between the VLM consortium and the reasoning LLM. All VLMs were fine-tuned using the Unsloth library with QLoRA-based quantization, enabling resource-efficient deployment on consumer-grade hardware without sacrificing performance. Experimental evaluations using both synthetic and real-world architecture diagrams demonstrate that ASTRIDE significantly improves the accuracy of threat identification, reduces the dependence on human experts, and offers a scalable and explainable alternative to traditional threat modeling practices. For future work, we plan to expand the VLM ensemble with newer open-source multimodal models to further increase prediction performance and robustness.

## References

[1] D. B. Acharya, K. Kuppan, B. Divya, Agentic ai: Autonomous intelligence for complex goals–a comprehensive survey, IEEE Access (2025).

[2] E. Bandara, R. Gore, X. Liang, S. Rajapakse, I. Kularathne, P. Karunarathna, P. Foytik, S. Shetty, R. Mukkamala, A. Rahman, et al., Agentsway–software development methodology for ai agents-based teams, arXiv preprint arXiv:2510.23664 (2025).

[3] B. Cao, Y. Cao, L. Lin, J. Chen, Defending against alignment-breaking attacks via robustly aligned llm, arXiv preprint arXiv:2309.14348 (2023).

[4] R. Khan, K. McLaughlin, D. Laverty, S. Sezer, Stride-based threat modeling for cyber-physical systems, in: 2017 IEEE PES Innovative Smart Grid Technologies Conference Europe (ISGT-Europe), IEEE, 2017, pp. 1–6.

[5] E. Bandara, A. Hass, S. Shetty, R. Mukkamala, R. Gore, A. Rahman, S. H. Bouk, Deep-stride: Automated security threat modeling with

vision-language models, in: 2025 International Conference on Software, Telecommunications and Computer Networks (SoftCOM), 2025, pp. 1–7.

[6] D. S. Cruzes, M. G. Jaatun, K. Bernsmed, I. A. Tøndel, Challenges and experiences with applying microsoft threat modeling in agile development projects, in: 2018 25th Australasian Software Engineering Conference (ASWEC), IEEE, 2018, pp. 111–120.

[7] J. Zhang, J. Huang, S. Jin, S. Lu, Vision-language models for vision tasks: A survey, IEEE Transactions on Pattern Analysis and Machine Intelligence (2024).

[8] Y. Zhang, S. Mao, T. Ge, X. Wang, A. de Wynter, Y. Xia, W. Wu, T. Song, M. Lan, F. Wei, Llm as a mastermind: A survey of strategic reasoning with large language models, arXiv preprint arXiv:2404.01230 (2024).

[9] Y. Zheng, R. Zhang, J. Zhang, Y. Ye, Z. Luo, Z. Feng, Y. Ma, Llamafactory: Unified efficient fine-tuning of 100+ language models, arXiv preprint arXiv:2403.13372 (2024).

[10] L. Mauri, E. Damiani, Modeling threats to ai-ml systems using stride, Sensors 22 (17) (2022) 6662.

[11] V. Kumar, J. Mayo, K. Bahiss, Admin: Attacks on dataset, model and input. a threat model for ai based software, arXiv preprint arXiv:2401.07960 (2024).

[12] S. Yang, T. Wu, S. Liu, D. Nguyen, S. Jang, A. Abuadbba, Threatmodeling-llm: Automating threat modeling using large language models for banking system, arXiv preprint arXiv:2411.17058 (2024).

[13] M. Mollaeefar, A. Bissoli, S. Ranise, Pillar: an ai-powered privacy threat modeling tool, arXiv preprint arXiv:2410.08755 (2024).

[14] N. Nagaraja, H. Bahsi, Cyber threat modeling of an llm-based healthcare system.

[15] A. Crossman, A. R. Plummer, C. Sekharudu, D. Warrier, M. Yekrangian, Auspex: Building threat modeling tradecraft into an artificial intelligence-based copilot, arXiv preprint arXiv:2503.09586 (2025).

[16] L. Mauri, E. Damiani, Stride-ai: An approach to identifying vulnerabilities of machine learning assets, in: 2021 IEEE International conference on cyber security and resilience (CSR), IEEE, 2021, pp. 147–154.

[17] Z. Dong, Z. Zhou, C. Yang, J. Shao, Y. Qiao, Attacks, defenses and evaluations for llm conversation safety: A survey, arXiv preprint arXiv:2402.09283 (2024).

[18] E. Bandara, R. Gore, S. Shetty, R. Mukkamala, C. Rhea, A. Yarlagadda, S. Kaushik, L. De Silva, A. Maznychenko, I. Sokolowska, et al., Standardization of neuromuscular reflex analysis–role of fine-tuned vision-language model consortium and openai gpt-oss reasoning llm enabled decision support system, arXiv preprint arXiv:2508.12473 (2025).

[19] S. Agarwal, L. Ahmad, J. Ai, S. Altman, A. Applebaum, E. Arbus, R. K. Arora, Y. Bai, B. Baker, H. Bao, et al., gpt-oss-120b & gpt-oss-20b model card, arXiv preprint arXiv:2508.10925 (2025).

[20] E. Bandara, S. H. Bouk, S. Shetty, R. Gore, S. Kompella, R. Mukkamala, A. Rahman, P. Foytik, X. Liang, N. W. Keong, K. De Zoysa, Vindsec-llama — fine-tuned meta's llama-3 llm, federated learning, blockchain and pbom-enabled data security architecture for wind energy data platforms, in: 2025 International Wireless Communications and Mobile Computing (IWCMC), 2025, pp. 120–126. doi:10.1109/IWCMC65282.2025.11059510.

[21] G. Mondillo, M. Masino, S. Colosimo, A. Perrotta, V. Frattolillo, Evaluating ai reasoning models in pediatric medicine: A comparative analysis of o3-mini and o3-mini-high, medRxiv (2025) 2025–02.

[22] E. Bandara, R. Gore, A. Yarlagadda, A. H. Clayton, P. Samuel, C. K. Rhea, S. Shetty, Standardization of psychiatric diagnoses–role of fine-tuned llm consortium and openai-gpt-oss reasoning llm enabled decision support system, arXiv preprint arXiv:2510.25588 (2025).

[23] E. Bandara, S. Shetty, R. Mukkamala, R. Gore, P. Foytik, S. H. Bouk, A. Rahman, X. Liang, N. W. Keong, K. De Zoysa, et al., Model context contracts-mcp-enabled framework to integrate llms with blockchain smart contracts, arXiv preprint arXiv:2510.19856 (2025).

[24] H. Samo, K. Ali, M. Memon, F. A. Abbasi, M. Y. Koondhar, K. Dahri, Fine-tuning mistral 7b large language model for python query response and code generation: A parameter efficient approach, VAWKUM Transactions on Computer Sciences 12 (1) (2024) 205–217.

[25] E. Bandara, S. Shetty, R. Mukkamala, A. Rahman, P. Foytik, X. Liang, K. De Zoysa, N. W. Keong, Devsec-gpt — generative-ai (with custom-trained meta's llama2 llm), blockchain, nft and pbom enabled cloud native container vulnerability management and pipeline verification platform, in: 2024 IEEE Cloud Summit, 2024, pp. 28–35. doi:10.1109/Cloud-Summit61220.2024.00012.

[26] J. Zhang, J. Huang, S. Jin, S. Lu, Vision-language models for vision tasks: A survey, IEEE Transactions on Pattern Analysis and Machine Intelligence (2024).

[27] R. Gore, E. Bandara, S. Shetty, A. E. Musto, P. Rana, A. Valencia-Romero, C. Rhea, L. Tayebi, H. Richter, A. Yarlagadda, et al., Proof-of-tbi–fine-tuned vision language model consortium and openai-o3 reasoning llm-based medical diagnosis support system for mild traumatic brain injury (tbi) prediction, arXiv preprint arXiv:2504.18671 (2025).

[28] T. C. Davis, How do you know a mermaid when you see one? how do you see a mermaid when you know one?, Theatre Journal 71 (3) (2019) 257–288.

[29] T. Dettmers, A. Pagnoni, A. Holtzman, L. Zettlemoyer, Qlora: Efficient finetuning of quantized llms, Advances in Neural Information Processing Systems 36 (2024).

[30] T. Reason, E. Benbow, J. Langham, A. Gimblett, S. L. Klijn, B. Malcolm, Artificial intelligence to automate network meta-analyses: Four case studies to evaluate the potential application of large language models, PharmacoEconomics-Open (2024) 1–16.