

A Deep Learning based Cloud Intrusion Detection System

Video Demo & Explanation:  B22CS011_B22CS026_VCC_Project.mp4

Table of Contents

- Problem Statement
- Literature Reference & Existing Results
- Implementation
 - Step-by-step implementation
 - Creating the model
 - Deploying the model
 - Creating the deployment code
 - Creating the AI Agent VM
 - Running the startup script
 - Creating and deploying the exploitable web application
 - Creating the app hosting VM
 - Running the startup script
 - Using packet monitoring for feature extraction
 - Results
 - DL model
 - Cloud Web application deployment
 - Comparison and Analysis
 - Architecture Diagram
 - Future Scope

Problem Statement

Intrusion detection is one of the technologies that protects the cloud computing environment from malicious attacks. However, network traffic in the cloud computing environment is characterised by large scale, high

dimensionality, and high redundancy, and these characteristics pose serious challenges to the development of cloud intrusion detection systems (IDS). Deep learning (DL) techniques have shown considerable potential for intrusion detection. We want to devise and test the effectiveness of various DL models for intrusion detection in a real-world scenario through a self-created application hosted across cloud virtual machines (VMs).

Literature Reference & Existing Results

IDS for cloud environments have become increasingly complex due to the growing volume and complexity of cyber threats. Traditional signature-based IDS approaches often fail to detect such novel attacks. High-dimensional and large-scale network traffic associated with the cloud computing environment makes things more complicated. Recent research has turned towards DL-based methodologies for better anomaly detection.

Wang et al. (2022) proposed a hybrid intrusion detection framework that leverages a Stacked Contractive Autoencoder (SCAE) for unsupervised feature extraction and then utilise a Support Vector Machine (SVM) as the final layer for the classification task.

. The SCAE tries to reduce the high-dimensional network data into a more compact representation that is more informative for performing the classification. They benchmarked their model evaluations on NSL-KDD and KDD Cup 99 datasets, with a high accuracy of 0.92 and 0.99 in the respective datasets.

Building upon this, Dalal et al. (2023) introduced an Extremely Boosted Neural Network to predict multi-stage cyberattacks in a cloud computing context. They wanted to train their model on the Multi-Step Cyber-Attack Dataset (MSCAD) to create a model that is stage-aware, thus better suited for advanced persistent threats (APTs) and zero-day attacks. Their model consisted of an ensemble of neural networks, which was able to achieve a remarkable 99.72% detection accuracy on the MSCAD dataset, surpassing conventional approaches such as Bayesian Networks and QUEST decision trees.

Inspired by this, we decided to utilise a DL regularisation technique known as dropout, which we studied to have similar performance as ensemble learning but with a fraction of the network size.

Implementation

After observing these promising results in these studies, we developed the DL-based IDS using a deep neural network architecture. The model consisted of 3 hidden MLP layers, with dropout with $p=0.3$ applied on the first two hidden layers. A general logic in designing neural networks is to have the first hidden layer size to be greater than the input size ($=31$). Thus, the hidden layer size was chosen to be 64.

The model was trained and validated using the KDD Cup 99 dataset. The model exhibited a very high train accuracy of 0.9994 & validation accuracy of 0.9996 on the 5-class classification task.

Given the encouraging performance, we decided to move further with the model and test it in a real-world setting by testing it on network packets captured from a Google Cloud Platform (GCP) VM hosting an exploitable web application that we designed and deployed. This deployment allowed us to evaluate the model's reliability under realistic network conditions, with incoming traffic logged and processed continuously.

Step-by-step implementation

Creating the model

The Python Pytorch library was used to build the model, and the dataset was sourced from Kaggle. After sufficient data pre-processing, the model training was set for 20 epochs with early stopping in case validation accuracy decreases.

Deploying the model

Once the model was trained, it was deployed on a GCP VM so that it could be queried by the system running the application for prediction. The VM on which the model was deployed was separate from the VM running the app, since the DL model may be resource-intensive and thus may hinder application performance.

Creating the deployment code

After training, the trained model was saved as a .pt file. We wrote a Python script, 'agent.py', that listens (using socket programming) for queries to get the prediction for the data sent over the network.

It then uses the saved DL model to get the prediction(s) and sends them back to the client.

Creating the AI Agent VM

A GCP VM with the following configuration was created-

1. Region: us-central1-a
2. Machine Type: e2-medium (2 vCPU, 1 Core, 4GB memory)
3. OS: Ubuntu 24.04 LTS
4. Storage: 25 GB balanced persistent disk (more than the standard 10GB since the 'flask' python library is very large)
5. Firewall: allow HTTP & HTTPS traffic

The screenshot shows the 'Create an instance' wizard in the Google Cloud Platform. On the left, a sidebar lists various configuration sections: Machine configuration (selected), OS and storage, Data protection, Networking, Observability, Security, and Advanced. The main panel displays the 'Machine configuration' settings. It includes fields for 'Name' (set to 'agentvm'), 'Region' (set to 'us-central1 (Iowa)'), 'Zone' (set to 'us-central1-a'), and a 'General purpose' machine type (selected). Below this, a table lists various machine types with their descriptions, vCPUs, memory, and CPU platform. The 'E2' machine type is selected. To the right, a 'Monthly estimate' table shows the cost breakdown: \$24.46 for 2 vCPU + 4 GB memory, \$2.50 for 25 GB balanced persistent disk, and additional costs for Logging (\$0.04/hour), Monitoring (\$0.04/hour), and Total (\$26.96). Buttons at the bottom include 'Create' and 'Cancel'.

Creating the agent VM with the above configurations

The screenshot shows the Google Cloud Platform interface for creating a new VM instance. The left sidebar lists several configuration sections: Machine configuration, OS and storage, Data protection, Networking (selected), Observability, Security, and Advanced. The main area is titled "Networking" and contains a "Firewall" section with options to allow HTTP and HTTPS traffic. Below this are fields for "Network tags" and "Hostname". Under "Network performance configuration", there's a "Network bandwidth" section with an option to enable per-VM Tier_1 networking performance. On the right, a "Monthly estimate" table shows costs for 2 vCPU + 4 GB memory (\$24.46), 25 GB balanced persistent disk (\$2.50), Logging (Cost varies), Monitoring (Cost varies), and a total monthly estimate of \$26.96. Buttons at the bottom include "Create", "Cancel", and "Equivalent code".

Changing the network settings to allow internet communication between the agent VM and the app VM

Running the startup script

The following shell script was executed to install the required libraries and run the “AI agent” script. This script can be found in the project’s git repository.

The screenshot shows the Google Cloud Platform VM instances list. The top navigation bar includes "VM instances", "Create Instance", "Import VM", and "Refresh". The main area has tabs for "Instances", "Observability", and "Instance schedules". A message states: "Your project's VMs use global DNS names by default. To reduce the risk of cross-regional outages, we recommend you use zonal DNS instead. [Learn more](#)". The "VM instances" table lists one instance: "agent-vm" (Status: Running, Name: agent-vm, Zone: us-central1-a, Internal IP: 10.128.0.2 (nic0), External IP: 34.56.118.72 (nic0), Connect via SSH). The table includes columns for Status, Name, Zone, Recommendations, In use by, Internal IP, External IP, and Connect.

Start the agent VM and connect to the secure shell (SSH)

The screenshot shows an SSH-in-browser interface. The terminal window displays the following text:

```
Welcome to Ubuntu 24.04.2 LTS (GNU/Linux 6.8.0-1027-gcp x86_64)

 * Documentation: https://help.ubuntu.com
 * Management: https://landscape.canonical.com
 * Support: https://ubuntu.com/pro

System information as of Thu Apr 17 21:07:54 UTC 2025

System load: 1.09      Processes: 118
Usage of /: 12.5% of 23.17GB  Users logged in: 0
Memory usage: 10%          IPv4 address for ens4: 10.128.0.2
Swap usage: 0%

Expanded Security Maintenance for Applications is not enabled.

14 updates can be applied immediately.
4 of these updates are standard security updates.
To see these additional updates run: apt list --upgradable

Enable ESM Apps to receive additional future security updates.
See https://ubuntu.com/esm or run: sudo pro status

The programs included with the Ubuntu system are free software;
the exact distribution terms for each program are described in the
individual files in /usr/share/doc/*copyright.

Ubuntu comes with ABSOLUTELY NO WARRANTY, to the extent permitted by
applicable law.

b22cs011@agent-vm:~$ sudo apt update && sudo apt -y install git python3-pip virtualenv
git clone https://github.com/Anoop-CS011/VCC-Project.git
cd VCC-Project
virtualenv myProjectEnv
source myProjectEnv/bin/activate
pip install -r requirements.txt
python agent.py
```

Copy and execute the agent's shell script.

```
(myProjectEnv) b22cs011@agent-vm:~/VCC-Project$ python agent.py
Model loaded...
[Agent] Listening on 0.0.0.0:5000
[]
```

Once the agent script is loaded, it will be available for prediction

Creating and deploying the exploitable web application

Creating the app hosting VM

A GCP VM with the following configuration was created-

1. Region: us-central1-a
2. Machine Type: e2-medium (2 vCPU, 1 Core, 4GB memory)
3. OS: Ubuntu 24.04 LTS
4. Storage: 25 GB balanced persistent disk (more than the standard 10GB since the 'flask' python library is very large)
5. Firewall: allow HTTP & HTTPS traffic

The screenshot shows the 'Create an instance' wizard in the Google Cloud Platform. The left sidebar lists configuration steps: Machine configuration (e2-medium, us-central1-a), OS and storage (Ubuntu 24.04 LTS), Data protection (Snapshot schedules), Networking (2 firewall rules, 1 network interface), Observability (Install Ops Agent), Security, and Advanced. The main panel is titled 'Machine configuration' and shows the following details:

- Name: mainapp-vm
- Region: us-central1 (Iowa)
- Zone: us-central1-a
- General purpose selected (Compute optimized, Memory optimized, Storage optimized, GPUs are disabled)
- Machine types for common workloads, optimized for cost and flexibility:

Series	Description	vCPUs	Memory	CPU Platform
C4	Consistently high performance	2 - 192	4 - 1,488 GB	Intel Emerald
C4A	Arm-based consistently high performance	1 - 72	2 - 576 GB	Google Axion
N4	Flexible & cost-optimized	2 - 80	4 - 640 GB	Intel Emerald
C3	Consistently high performance	4 - 192	8 - 1,536 GB	Intel Sapphire
C3D	Consistently high performance	4 - 360	8 - 2,880 GB	AMD Genoa
E2	Low cost, day-to-day computing	0.25 - 32	1 - 128 GB	Intel Broadwell
N2	Balanced price & performance	2 - 128	2 - 864 GB	Intel Cascade
N2D	Balanced price & performance	2 - 224	2 - 896 GB	AMD Milan
T2A	Scale-out workloads	1 - 48	4 - 192 GB	Ampere Altra
T2D	Scale-out workloads	1 - 60	4 - 240 GB	AMD Milan
N1	Balanced price & performance	0.25 - 96	0.6 - 624 GB	Intel Haswell
- Buttons at the bottom: Create, Cancel, and Equivalent code.

On the right, a 'Monthly estimate' table shows costs for the selected configuration:

Item	Monthly estimate
2 vCPU + 4 GB memory	\$24.46
25 GB balanced persistent disk	\$2.50
Logging	Cost varies
Monitoring	Cost varies
Snapshot schedule	Cost varies
Total	\$26.96

Links for Compute Engine pricing and Cloud Operations pricing are also present.

Creating the app VM with the above configurations

Running the startup script

The following shell script was executed to install the required libraries and deploy the app. This script can be found in the project's git repository.

The screenshot shows the 'VM instances' list in the Google Cloud Platform. The table displays the following information:

Status	Name	Zone	Recommendations	In use by	Internal IP	External IP	Connect
Running	agent-vm	us-central1-a			10.128.0.4 (nic0)		SSH
Running	mainapp-vm	us-central1-a			10.128.0.3 (nic0)	34.58.115.20 (nic0)	SSH

Start the application VM and connect to the SSH.

```
Welcome to Ubuntu 24.04.2 LTS (GNU/Linux 6.8.0-1027-gcp x86_64)

 * Documentation: https://help.ubuntu.com
 * Management: https://landscape.canonical.com
 * Support: https://ubuntu.com/pro

System information as of Thu Apr 24 06:54:30 UTC 2025

System load: 1.57      Processes: 120
Usage of /: 52.4% of 23.17GB  Users logged in: 0
Memory usage: 7%          IPv4 address for ens4: 10.128.0.3
Swap usage: 0%

* Strictly confined Kubernetes makes edge and IoT secure. Learn how MicroK8s just raised the bar for easy, resilient and secure K8s cluster deployment.

https://ubuntu.com/engage/secure-kubernetes-at-the-edge

Expanded Security Maintenance for Applications is not enabled.

9 updates can be applied immediately.
To see these additional updates run: apt list --upgradable

2 additional security updates can be applied with ESM Apps.
Learn more about enabling ESM Apps service at https://ubuntu.com/esm

Last login: Sat Apr 19 18:34:43 2025 from 35.235.244.32
b22cs011@mainapp-vm:~$ sudo apt update && sudo apt -y install git python3-pip virtualenv
sudo apt -y install python3-scapy
git clone https://github.com/Anoop-cs011/VCC-Project.git
cd VCC-Project
virtualenv myProjectEnv
source myProjectEnv/bin/activate
pip install -r requirements.txt
unicorn --bind 0.0.0.0:8080 app:app --daemon
sudo ./myProjectEnv/bin/python monitor.py
```

Copy and execute the shell script to deploy the app.

Using packet monitoring for feature extraction

The monitoring script was created in Python. It uses the ‘scapy’ library for packet “sniffing”. Once a packet is sniffed, features are calculated and extracted from it. The features are stored until a ‘batch’ of packets is sniffed. Then, the ‘batch’ is sent to the agent VM over the network using socket programming. The prediction returned from the model is stored in a local file, which the app can read to display recent alerts.

The monitoring script also begins execution once the “app” shell script is executed.

```
(myProjectEnv) b22cs011@mainapp-vm:~/VCC-Project$ gunicorn --bind 0.0.0.0:8080 app:app --daemon
sudo ./myProjectEnv/bin/python monitor.py
[Monitor] Starting packet sniffer...
```

The monitoring script is also executed once the "app" shell script is executed.

```
Welcome to Ubuntu 24.04.2 LTS (GNU/Linux 6.8.0-1027-gcp x86_64)

* Documentation: https://help.ubuntu.com
* Management: https://landscape.canonical.com
* Support: https://ubuntu.com/pro

System information as of Thu Apr 24 10:38:45 UTC 2025

System load: 0.09 Processes: 117
Usage of /: 52.3% of 23.17GB Users logged in: 0
Memory usage: 7%
IPv4 address for ens4: 10.128.0.4
Swap usage: 0%

* Strictly confined Kubernetes makes edge and IoT secure. Learn how MicroK8s just raised the bar for easy, resilient and secure K8s cluster deployment.
https://ubuntu.com/engage/secure-kubernetes-at-the-edge

Expanded Security Maintenance for Applications is not enabled.

9 updates can be applied immediately.
To see these additional updates run: apt list --upgradable

2 additional security updates can be applied with ESM Apps.
Learn more about enabling ESM Apps service at https://ubuntu.com/esm

Last login: Sat Apr 19 18:34:41 2025 from 35.235.244.34
b22cs011@agent-vm:~$ cd VCC-Project/
b22cs011@agent-vm:~/VCC-Project$ source myProjectEnv/bin/activate
(myProjectEnv) b22cs011@agent-vm:~/VCC-Project$ python agent.py
Model loaded...
[Agent] Listening on 0.0.0.0:5000
[Agent] Connection from ('10.128.0.3', 48018)
[Agent] Connection from ('10.128.0.3', 59906)
```



```
Welcome to Ubuntu 24.04.2 LTS (GNU/Linux 6.8.0-1027-gcp x86_64)

* Documentation: https://help.ubuntu.com
* Management: https://landscape.canonical.com
* Support: https://ubuntu.com/pro

System information as of Thu Apr 24 10:38:24 UTC 2025

System load: 0.12 Processes: 111
Usage of /: 53.3% of 23.17GB Users logged in: 0
Memory usage: 10%
IPv4 address for ens4: 10.128.0.3
Swap usage: 0%

* Strictly confined Kubernetes makes edge and IoT secure. Learn how MicroK8s just raised the bar for easy, resilient and secure K8s cluster deployment.
https://ubuntu.com/engage/secure-kubernetes-at-the-edge

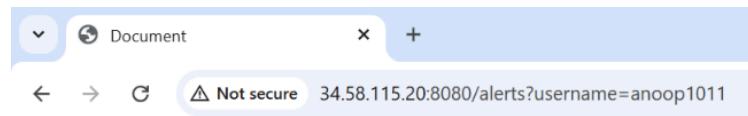
Expanded Security Maintenance for Applications is not enabled.

0 updates can be applied immediately.

2 additional security updates can be applied with ESM Apps.
Learn more about enabling ESM Apps service at https://ubuntu.com/esm

*** System restart required ***
Last login: Thu Apr 24 06:54:32 2025 from 35.235.244.32
b22cs011@mainapp-vm:~/VCC-Project/
b22cs011@mainapp-vm:~/VCC-Project$ source myProjectEnv/bin/activate
(myProjectEnv) b22cs011@mainapp-vm:~/VCC-Project$ gunicorn --bind 0.0.0.0:8080 app:app --daemon
sudo ./myProjectEnv/bin/python monitor.py
[Monitor] Starting packet sniffer...
```

The monitoring script is sniffing packets and sending batch requests to the agent VM



Recent Intrusion Alerts

- at 1745491298.8816972 → : from 10.128.0.3 → type : dos
- at 1745491299.076896 → : from 134.199.145.46 → type : dos
- at 1745491299.0821898 → : from 134.199.145.46 → type : dos
- at 1745491300.387726 → : from 35.235.244.34 → type : dos
- at 1745491300.3883903 → : from 10.128.0.3 → type : dos

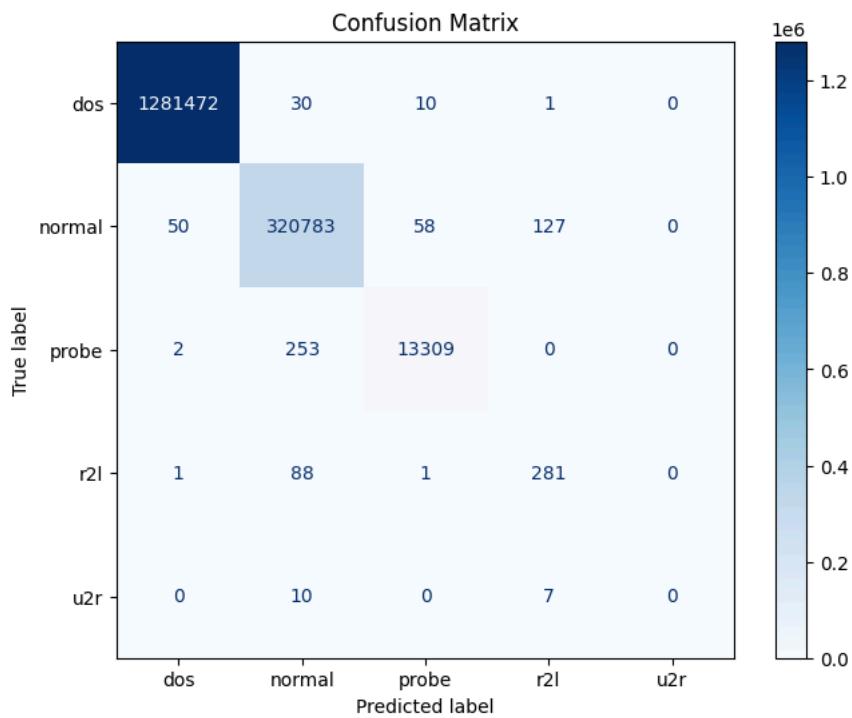
[Back to Home](#)

The recent intrusions can be seen from the application website after registering and logging in.

Results

DL model

As discussed earlier, the trained model was able to achieve a remarkable train accuracy of 0.9994 & validation accuracy of 0.9996 on the 5-class classification task, which is similar, if not better than the models studied from the referenced literature.



Cloud Web application deployment

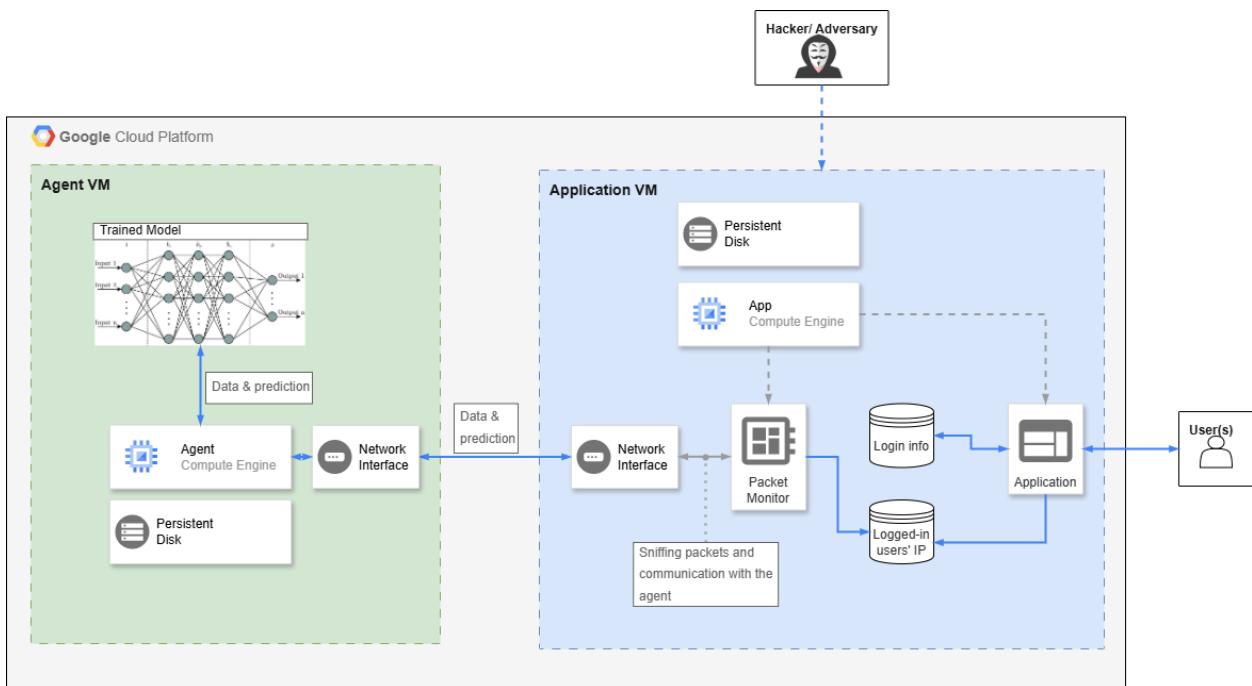
Given the high accuracy of the DL model, we were surprised to see that the model was not able to perform well on the live network data that was collected. The primary concern was that it was predicting normal network traffic as an intrusion, particularly as a DOS attack, which raised questions about the model's generalisation capability outside the dataset. This discrepancy could possibly be due to the skewed nature of the training dataset, which has an overrepresentation of the DOS samples compared to other classes. Thus, the model may have developed a bias towards DOS and begin misclassifying ambiguous or previously

unseen packets. The model probably overfits the benchmark data distribution and fails to adapt to the natural variability and noise present in real-world network traffic.

Comparison and Analysis

Such a good DL model not being able to perform well in live deployment is probably attributed to the way these models make predictions. They take one packet at a time as input and make their prediction. However, in real-world scenarios, the intrusive packets are meant to be similar to ordinary traffic. For example, in a Smurf DOS attack, ordinary ICMP echo packets are flooded into the network with a spoofed IP address of the target machine. Here, individual ICMP packets appear to be normal, but if we simply observe the abnormally large ICMP traffic, we can easily deduce that the system is under attack. Such a mechanism of detecting correlation across packets to predict intrusion can be trained using models like transformers, which utilise a cross-attention mechanism to help learn the relation between data points across the input data. These models can be trained to find patterns from a sequence of packets, rather than from one packet at a time.

Architecture Diagram



Future Scope

Future work may focus on building models that not only find patterns in individual packets but also across a sequence of packets. Using transformers might be an excellent method to realise this mechanism. However, they might be too resource-intensive, even for just inference, to keep the costs low for a feasible deployment. This could also be seen as a challenge to build a “tiny” transformer with reasonable performance in a resource-constrained environment.

The test production environment can be improved with a better design of the exploitable web application. This could be done by incorporating more features into the web application to simulate the real-world environment more accurately.

Different datasets like MSCAD & NSL-KDD can be explored and used to train and test more models. These models could possibly be used simultaneously for inference by extracting features used in multiple datasets.