

VCC-Project

Group-58

Team Members:

| ID | Name | Worked on below sub topic analysis |
|-----------|----------------|---|
| M23AID026 | Nikita Krishna | Designed and developed system to publish unauthorised login attempts to GCP VM. |
| M23AID006 | Aparna Pundir | Designed and developed system to push unauthorized login attempt messages from pub/sub to BigQuery. |
| M23AID053 | Bhuvaneswari J | Designed and developed K means ML model on Bigquery data. |

Insider Threat Detection System on GCP

This project simulates **malicious insider activity** such as unauthorized VM logins and suspicious data transfers, and uses **Google Cloud Platform (GCP)** to detect, log, analyse, and classify these behaviours.

Problem Statement:

In organizations, insider threats — malicious activities originating from within an organization by trusted users — pose significant security risks. Unlike external attacks, these threats are difficult to detect using traditional rule-based systems due to the legitimate access and behaviour of insiders.

The goal of this project is to detect potential insider threats by analysing user login behaviours and access patterns using unsupervised learning techniques, specifically K-Means clustering, in Google BigQuery. By clustering user activity data, we aim to identify anomalous behaviour that deviates significantly from typical usage patterns, such as unusual login times, access to sensitive resources, or repeated login failures.

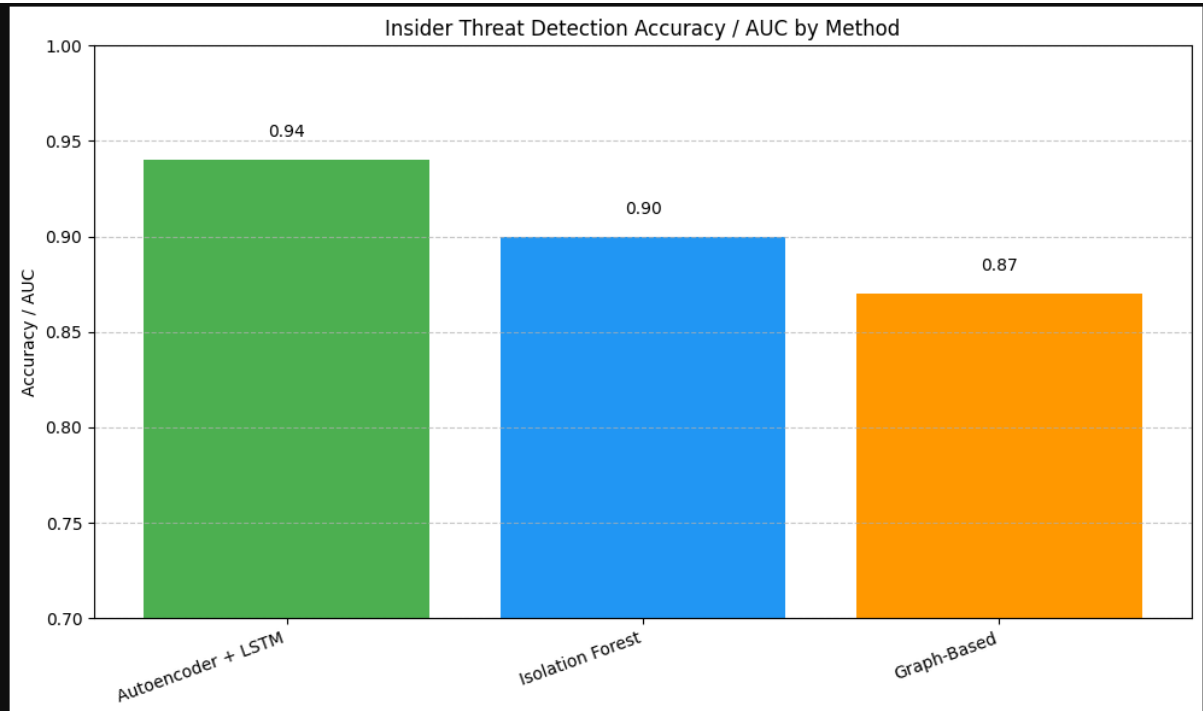
Literature Reference:

Insider threats, which originate from trusted individuals within an organization, are difficult to detect using traditional rule-based systems due to their legitimate access patterns. Various research efforts have explored this problem using different techniques. Eberle and Holder (2009) employed graph-based anomaly detection to model user interactions, but their approach lacked scalability. Liu et al. (2008) introduced Isolation Forest for anomaly detection, yet its interpretability and adaptability to categorical login data remain limited. Tuor et al. (2017) applied deep learning methods like autoencoders and LSTMs for structured cybersecurity data, though such methods require significant computational resources and are less interpretable. Salem and Stolfo (2011) relied on OS-level audit logs and statistical profiling, which are prone to false positives and hard to generalize. A more practical approach was presented in a 2020 Google Cloud Blog, demonstrating the use of BigQuery

ML for detecting anomalies in login data, albeit in simplified use cases. In contrast, this project leverages unsupervised learning, specifically K-Means clustering in Google BigQuery, to dynamically identify anomalous login behaviours and access patterns, providing a scalable and interpretable solution for insider threat detection.

Existing Results:

| Study / Solution | Technique Used | Dataset | Detection Accuracy / AUC | Strengths | Limitations |
|------------------------|------------------------------------|-----------------------------|--------------------------|---|---------------------------------|
| Tuor et al. (2017) | Deep Learning (Autoencoder + LSTM) | CERT Insider Threat Dataset | AUC: 0.94 | Effective in modeling temporal patterns | High complexity, needs GPU |
| Liu et al. (2008) | Isolation Forest | Synthetic anomaly dataset | Accuracy: 85–92% | Fast and scalable | Less intuitive output |
| Eberle & Holder (2009) | Graph Anomaly Detection | Simulated insider dataset | Accuracy: ~87% | Captures structural anomalies | Hard to scale with massive data |

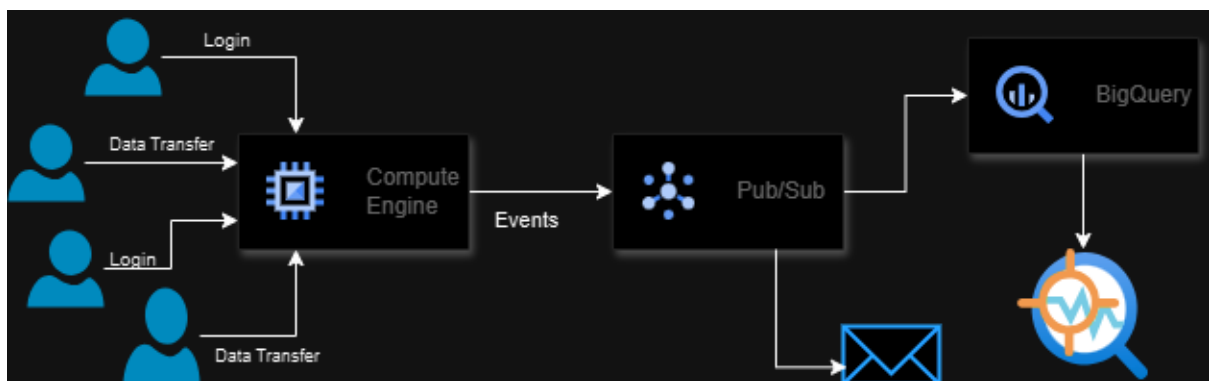


System Overview:

The Proposed solution aims to detect insider threats using:

- Simulated login attempts/data transfer attack to a GCP VM
- Real-time monitoring via Google Pub/Sub
- Storage and analysis using BigQuery
- K-Means clustering for unsupervised learning.
- Google Colab for interactive development.

System Architecture:



Implementation and Results:

Colab file link:

https://colab.research.google.com/drive/1d8JR2i8_R_3QqnD3QU9rMO0KzuS2syXk?usp=sharing

Install required libraries:

```
!pip install --upgrade google-api-python-client google-auth google-cloud-pubsub google-auth-httpplib2 google-auth-oauthlib
!pip install paramiko
```

Enable required API's:

```
from googleapiclient.discovery import build

serviceusage = build('serviceusage', 'v1')
request = serviceusage.services().enable(
    name='projects/insider-detector-
data/services/compute.googleapis.com'
)
response = request.execute()
```

```

try:
    response = request.execute()
    print("✔ Compute Engine API enabled successfully.")
except Exception as e:
    print(f"✗ Failed to enable Compute Engine API: {e}")

```

WARNING:google_auth_httplib2:httplib2 transport does not support per-request timeout. Set the timeout when constructing the httplib2.Http instance.
 WARNING:google_auth_httplib2:httplib2 transport does not support per-request timeout. Set the timeout when constructing the httplib2.Http instance.
 ✔ Compute Engine API enabled successfully.

```

from googleapiclient.discovery import build

serviceusage = build('serviceusage', 'v1')
request = serviceusage.services().enable(
    name='projects/insider-detector-data/services/iam.googleapis.com'
)
response = request.execute()

try:
    response = request.execute()
    print("✔ IAM API enabled successfully.")
except Exception as e:
    print(f"✗ Failed to enable IAM API: {e}")

```

WARNING:google_auth_httplib2:httplib2 transport does not support per-request timeout. Set the timeout when constructing the httplib2.Http instance.
 WARNING:google_auth_httplib2:httplib2 transport does not support per-request timeout. Set the timeout when constructing the httplib2.Http instance.
 ✔ IAM API enabled successfully.

```

from googleapiclient.discovery import build

serviceusage = build('serviceusage', 'v1')
request = serviceusage.services().enable(
    name='projects/insider-detector-
data/services/bigquery.googleapis.com'
)
response = request.execute()

try:
    response = request.execute()
    print("✔ Bigquery API enabled successfully.")
except Exception as e:
    print(f"✗ Failed to enable Bigquery API: {e}")

```

WARNING:google_auth_httplib2:httplib2 transport does not support per-request timeout. Set the timeout when constructing the httplib2.Http instance.
 WARNING:google_auth_httplib2:httplib2 transport does not support per-request timeout. Set the timeout when constructing the httplib2.Http instance.
 ✔ Bigquery API enabled successfully.

```

from googleapiclient.discovery import build

serviceusage = build('serviceusage', 'v1')
request = serviceusage.services().enable(



```

```

        name='projects/insider-detector-
data/services/dataflow.googleapis.com'
    )
    response = request.execute()

try:
    response = request.execute()
    print("✔ Dataflow API enabled successfully.")
except Exception as e:
    print(f"✗ Failed to enable Dataflow API: {e}")

```

 WARNING:google_auth_httplib2:httplib2 transport does not support per-request timeout. Set the timeout when constructing the httplib2.Http instance.
 WARNING:google_auth_httplib2:httplib2 transport does not support per-request timeout. Set the timeout when constructing the httplib2.Http instance.
 Dataflow API enabled successfully.



```

from googleapiclient.discovery import build

serviceusage = build('serviceusage', 'v1')
request = serviceusage.services().enable(
    name='projects/insider-detector-
data/services/pubsub.googleapis.com'
)
response = request.execute()

try:
    response = request.execute()
    print("✔ Pubsub API enabled successfully.")
except Exception as e:
    print(f"✗ Failed to enable Pubsub API: {e}")

```

 WARNING:google_auth_httplib2:httplib2 transport does not support per-request timeout. Set the timeout when constructing the httplib2.Http instance.
 WARNING:google_auth_httplib2:httplib2 transport does not support per-request timeout. Set the timeout when constructing the httplib2.Http instance.
 Pubsub API enabled successfully.

Create

- Service account
- IAM roles
- Service account key

```

# STEP 0: Install required libraries
!pip install --upgrade google-api-python-client google-auth google-
auth-httplib2 google-auth-oauthlib

# STEP 1: Imports
import json
import base64
import time
from googleapiclient.errors import HttpError
from google.colab import auth, files

```

```

from googleapiclient import discovery

# STEP 2: Authenticate Colab to access Google Cloud
auth.authenticate_user()

# STEP 3: Configuration
PROJECT_ID = "insider-detector-data"
SERVICE_ACCOUNT_NAME = "insider-monitor-agent"
SERVICE_ACCOUNT_EMAIL =
f"{SERVICE_ACCOUNT_NAME}@{PROJECT_ID}.iam.gserviceaccount.com"
KEY_FILE_NAME = "insider_key.json"

# STEP 4: Initialize service clients
iam_service = discovery.build('iam', 'v1')
crm_service = discovery.build('cloudresourcemanager', 'v1')

# STEP 5: Create the service account
print("✦ Creating service account...")
try:
    iam_service.projects().serviceAccounts().create(
        name=f"projects/{PROJECT_ID}",
        body={
            "accountId": SERVICE_ACCOUNT_NAME,
            "serviceAccount": {
                "displayName": "Insider Monitor Agent"
            }
        }
    ).execute()
    print("✔ Service account created.")
except HttpError as e:
    if "already exists" in str(e):
        print("⚠ Service account already exists.")
    else:
        raise

# STEP 6: Wait until service account is ready
print("◻ Waiting for service account to become available...")
for attempt in range(10):
    try:
        sa = iam_service.projects().serviceAccounts().get(
            name=f"projects/{PROJECT_ID}/serviceAccounts/{SERVICE_ACCOUNT_EMAIL}"
        ).execute()
        print("✔ Service account is ready.")
        break
    except HttpError:
        print(f"🔄 Attempt {attempt + 1}: Not ready yet, waiting...")

```

```

        time.sleep(3)
    else:
        raise Exception("✗ Timed out waiting for service account to be
created.")

# STEP 7: Assign IAM roles
roles_to_assign = [
    "roles/pubsub.publisher",
    "roles/logging.logWriter",
    "roles/monitoring.metricWriter",
    "roles/bigquery.dataViewer",
    "roles/bigquery.dataEditor",
    "roles/bigquery.admin"
]

def assign_role(role):
    print(f"🔧 Assigning role: {role}")
    policy = crm_service.projects().getIamPolicy(resource=PROJECT_ID,
body={}).execute()
    member = f"serviceAccount:{SERVICE_ACCOUNT_EMAIL}"

    role_binding = next((b for b in policy['bindings'] if b['role'] ==
role), None)
    if role_binding:
        if member not in role_binding['members']:
            role_binding['members'].append(member)
    else:
        policy['bindings'].append({
            'role': role,
            'members': [member]
        })

    crm_service.projects().setIamPolicy(
        resource=PROJECT_ID,
        body={"policy": policy}
    ).execute()
    print(f"✅ Role {role} assigned.")

for role in roles_to_assign:
    assign_role(role)

# STEP 7: Generate a service account key
print(f"🔑 Creating service account key...")
key = iam_service.projects().serviceAccounts().keys().create(
    name=f"projects/{PROJECT_ID}/serviceAccounts/{SERVICE_ACCOUNT_EMAIL}",
    body={
        "privateKeyType": "TYPE_GOOGLE_CREDENTIALS_FILE",

```

```

        "keyAlgorithm": "KEY_ALG_RSA_2048"
    }
).execute()

# STEP 8: Save key to file
print(f"💾 Saving key to `{KEY_FILE_NAME}`")
decoded_key = base64.b64decode(key['privateKeyData']).decode('utf-8')
with open(KEY_FILE_NAME, 'w') as f:
    f.write(decoded_key)

# STEP 9: Download the key via Colab
print("📄 Download your key file using the link below")
files.download(KEY_FILE_NAME)

```

```

WARNING:google_auth_httplib2:httplib2 transport does |
🔴 Creating service account...
⚠️ Service account already exists.
⌚ Waiting for service account to become available...
WARNING:google_auth_httplib2:httplib2 transport does |
WARNING:google_auth_httplib2:httplib2 transport does |
✅ Service account is ready.
🔧 Assigning role: roles/pubsub.publisher
✅ Role roles/pubsub.publisher assigned.
🔧 Assigning role: roles/logging.logWriter
✅ Role roles/logging.logWriter assigned.
🔧 Assigning role: roles/monitoring.metricWriter
✅ Role roles/monitoring.metricWriter assigned.
🔧 Assigning role: roles/bigquery.dataViewer
✅ Role roles/bigquery.dataViewer assigned.
🔧 Assigning role: roles/bigquery.dataEditor
✅ Role roles/bigquery.dataEditor assigned.
🔧 Assigning role: roles/bigquery.admin
✅ Role roles/bigquery.admin assigned.
🔑 Creating service account key...
💾 Saving key to `insider_key.json`
📄 Download your key file using the link below

```

Create VM

```

import googleapiclient.discovery

project_id = "insider-detector-data"
zone = "us-central1-a"
instance_name = "test-vm-insider"

compute = googleapiclient.discovery.build('compute', 'v1')

config = {
    "name": instance_name,

```



```

"machineType": f"zones/{zone}/machineTypes/e2-micro",
"disks": [{
    "boot": True,
    "autoDelete": True,
    "initializeParams": {
        "sourceImage": "projects/debian-
cloud/global/images/family/debian-11"
    }
}],
"networkInterfaces": [{
    "network": "global/networks/default",
    "accessConfigs": [{"type": "ONE_TO_ONE_NAT", "name": "External
NAT"}]
}]
}

operation = compute.instances().insert(
    project=project_id,
    zone=zone,
    body=config
).execute()

print(f"VM {instance_name} creation requested.")

```

```

WARNING:google_auth_httplib2:httplib2 transport does not support per-request timeout.
WARNING:google_auth_httplib2:httplib2 transport does not support per-request timeout.
VM test-vm-insider creation requested.

```

VM instances

Filter Enter property name or value

| <input type="checkbox"/> Status | Name ↑ | Zone | Recommendations | In use by | Internal IP | External IP | Connect | |
|-------------------------------------|---------------------------------|---------------|-----------------|-----------|--|---|---------|-----|
| <input checked="" type="checkbox"/> | test-vm-insider | us-central1-a | | | 10.128.0.2 (nic0) | 34.55.252.241 (nic0) | SSH | ⌵ ⋮ |

Create Pub/Sub topic and subscription

```

# one time
from google.cloud import pubsub_v1

project_id = "insider-detector-data"
topic_id = "login-events"

publisher = pubsub_v1.PublisherClient()
topic_path = publisher.topic_path(project_id, topic_id)


# Create topic

```

```

try:
    publisher.create_topic(request={"name": topic_path})
    print(f"Topic '{topic_id}' created.")
except Exception as e:
    print(f"Topic already exists or failed: {e}")

```

| LIST METRICS | | | | |
|---|----------------|--|---|-----------|
| Filter Filter topics | | | | |
| <input type="checkbox"/> Topic ID ↑ | Encryption key | Topic name | | Retention |
| <input type="checkbox"/> login-events | Google-managed | projects/insider-detector-data/topics/login-events |  | — |

```

from google.cloud import pubsub_v1

# Set your Google Cloud project ID
project_id = "insider-detector-data"
topic_id = "login-events"
subscription_id = "login-events-sub"

# Full resource names
topic_path = f"projects/{project_id}/topics/{topic_id}"
subscription_path =
f"projects/{project_id}/subscriptions/{subscription_id}"

# Initialize subscriber client
subscriber = pubsub_v1.SubscriberClient()

# Create the subscription
with subscriber:
    try:
        subscription = subscriber.create_subscription(
            name=subscription_path,
            topic=topic_path
        )
        print(f"Subscription created: {subscription.name}")
    except Exception as e:
        print(f"Error creating subscription: {e}")

```

| SUBSCRIPTIONS | SNAPSHOTS | METRICS | DETAILS | MESSAGES |
|---|---|-----------------------|---------|----------|
| <p>Only subscriptions attached to this topic are displayed. A subscription captures the stream of messages published to a given topic. For more information on Cloud Pub/Sub Storage by creating a subscription from a Cloud Dataflow job. Learn more</p> <p> CREATE SUBSCRIPTION EXPORT </p> | | | | |
| <div> Filter Filter subscriptions </div> | | | | |
| Subscription ID ↑ | Subscription name | Project | | |
| login-events-sub | projects/insider-detector-data/subscriptions/login-events-sub | insider-detector-data | | |

Generate RSA Key

```
import paramiko

# File names for keys
private_key_path = "insider_ssh_key"
public_key_path = "insider_ssh_key.pub"

# Generate RSA key
key = paramiko.RSAKey.generate(bits=2048)

# Save private key
key.write_private_key_file(private_key_path)
print(f"✔ Private key saved to: {private_key_path}")

# Save public key
with open(public_key_path, "w") as pub_file:
    pub_file.write(f"{key.get_name()} {key.get_base64()}")
print(f"✔ Public key saved to: {public_key_path}")
```

```
|
✔ Private key saved to: insider_ssh_key
✔ Public key saved to: insider_ssh_key.pub
```

Create bigquery dataset and table

```
from google.cloud import bigquery

def create_bigquery_resources(project_id, dataset_id, table_id):
    client = bigquery.Client(project=project_id)

    # Create dataset reference
    dataset_ref = bigquery.DatasetReference(project_id, dataset_id)
```

```

dataset = bigquery.Dataset(dataset_ref)
dataset.location = "US"

# Create dataset if it doesn't exist
try:
    client.get_dataset(dataset_ref)
    print(f"✔ Dataset '{dataset_id}' already exists.")
except Exception:
    dataset = client.create_dataset(dataset)
    print(f"✔ Dataset '{dataset_id}' created.")

table_ref = dataset_ref.table(table_id)

# Delete table if it exists
try:
    client.delete_table(table_ref)
    print(f"🗑 Table '{table_id}' deleted.")
except Exception:
    print(f"❗ Table '{table_id}' did not exist, so nothing was
deleted.")

# Define updated table schema
schema = [
    bigquery.SchemaField("event_type", "STRING", mode="REQUIRED"),
    bigquery.SchemaField("timestamp", "TIMESTAMP",
mode="REQUIRED"),
    bigquery.SchemaField("user", "STRING", mode="NULLABLE"),
    bigquery.SchemaField("status", "STRING", mode="NULLABLE"),
    bigquery.SchemaField("role", "STRING", mode="NULLABLE"),
    bigquery.SchemaField("local_address", "STRING",
mode="NULLABLE"),
    bigquery.SchemaField("remote_address", "STRING",
mode="NULLABLE"),
]

table = bigquery.Table(table_ref, schema=schema)
client.create_table(table)
print(f"✔ Table '{table_id}' created with updated schema.")

# Call the function with your project details
create_bigquery_resources(
    project_id="insider-detector-data",
    dataset_id="secure_data",
    table_id="login_events"
)

```

- ✅ Dataset 'secure_data' already exists.
- 🗑️ Table 'login_events' deleted.
- ✅ Table 'login_events' created with updated schema.

The screenshot shows the Google Cloud BigQuery interface. On the left, the 'secure_data' dataset is expanded, showing the 'login_events' table. The main panel displays the 'Schema' tab for the 'login_events' table. The schema table lists the following fields:

| Field name | Type | Mode | Key | Collation | Default Value | Policy |
|----------------|-----------|----------|-----|-----------|---------------|--------|
| event_type | STRING | REQUIRED | - | - | - | - |
| timestamp | TIMESTAMP | REQUIRED | - | - | - | - |
| user | STRING | NULLABLE | - | - | - | - |
| status | STRING | NULLABLE | - | - | - | - |
| role | STRING | NULLABLE | - | - | - | - |
| local_address | STRING | NULLABLE | - | - | - | - |
| remote_address | STRING | NULLABLE | - | - | - | - |

- Add ssh keys to VM and assign required access to role user
- Simulate login attempt and data transfer using sample role user
- Publish the messages via pub/sub
- Insert messages into bigquery table.

```
import paramiko
import json
import time
import threading
from datetime import datetime
from google.cloud import pubsub_v1
from googleapiclient.discovery import build
from google.oauth2 import service_account
from google.cloud import bigquery

# GCP project and Pub/Sub setup
project_id = "insider-detector-data"
topic_id = "login-events"
subscription_id = "login-events-sub"
publisher = pubsub_v1.PublisherClient()
subscriber = pubsub_v1.SubscriberClient()
topic_path = publisher.topic_path(project_id, topic_id)
subscription_path = subscriber.subscription_path(project_id,
subscription_id)
bq_client = bigquery.Client()
dataset_id = "secure_data"
table_id = "login_events"
table_ref = bq_client.dataset(dataset_id).table(table_id)
```

```

# VM connection details
hostname = "34.55.252.241"
username = "m23aid053"
private_key_path = "insider_ssh_key"
public_key_path = "insider_ssh_key.pub"
service_account_file = "insider_key.json"

# Initialize Google Compute Engine API client
credentials =
service_account.Credentials.from_service_account_file(service_account_f
ile)
compute = build('compute', 'v1', credentials=credentials)

# Function to get the current fingerprint of the VM's metadata
def get_metadata_fingerprint():
    instance = compute.instances().get(
        project=project_id,
        zone="us-centrall-a",
        instance="test-vm-insider",
    ).execute()

    return instance['metadata']['fingerprint']

# Function to upload the SSH public key to VM metadata
def upload_ssh_key_to_vm():
    # Get the latest fingerprint
    fingerprint = get_metadata_fingerprint()

    # Read the SSH public key
    with open(public_key_path, 'r') as f:
        public_key = f.read().strip()

    # Prepare metadata entry
    metadata = {
        "items": [
            {
                "key": "ssh-keys",
                "value": f"{username}:{public_key}"
            }
        ]
    }

    # Insert the public key into the VM's metadata
    request = compute.instances().setMetadata(
        project=project_id,
        zone="us-centrall-a",
        instance="test-vm-insider",
        body={

```

```

        'fingerprint': fingerprint,
        'items': metadata["items"]
    }
)

response = request.execute()
print(f"✔ SSH key added to the VM metadata: {response}")

# GCP Pub/Sub event publishing function with role info
def publish_login_event(user, success, role):
    event = {
        "event_type": "login",
        "user": user,
        "timestamp": datetime.utcnow().isoformat() + "Z",
        "status": "success" if success else "failure",
        "role": role
    }
    data = json.dumps(event).encode("utf-8")
    publisher.publish(topic_path, data)
    print(f"Published: {event}")

# Function to simulate a login attempt
def simulate_login(user):
    user_roles = {
        "nikita": "user",
        "aparna": "user",
        "m23aid053": "admin",
        "admin": "admin",
        "unknown_user": "guest",
        "invalid": "user",
        "unknown": "guest",
        "malicious": "guest"
    }

    role = user_roles.get(user, "guest")
    key = paramiko.RSAKey.from_private_key_file(private_key_path)
    ssh = paramiko.SSHClient()
    ssh.set_missing_host_key_policy(paramiko.AutoAddPolicy())

    try:
        ssh.connect(hostname, username=user, pkey=key, timeout=5)
        publish_login_event(user, success=True, role=role)
        ssh.close()
    except Exception as e:
        print(f"Login failed for {user}: {e}")
        publish_login_event(user, success=False, role=role)

def callback(message):

```

```

try:
    data = json.loads(message.data.decode('utf-8'))
    print(f"📬 Received message: {data}")

    table_path = f"{project_id}.{dataset_id}.{table_id}"

    if data.get("event_type") == "login":
        row = [{
            "event_type": data.get("event_type"),
            "user": data.get("user"),
            "timestamp": data.get("timestamp"),
            "status": data.get("status"),
            "role": data.get("role")
        }]

        elif data.get("event_type") == "data_transfer":

            transfer_time =
datetime.utcnow().timestamp(data.get("timestamp")).isoformat() + "Z"

            row = [{
                "event_type": data.get("event_type"),
                "local_address": data.get("local_address"),
                "remote_address": data.get("remote_address"),
                "timestamp": transfer_time
            }]

    else:
        print("⚠ Skipped unknown event type.")
        message.ack()
        return

    # Insert into BigQuery
    errors = bq_client.insert_rows_json(table_path, row)
    if errors:
        print(f"❌ BigQuery insert error: {errors}")
    else:
        print(f"✅ Inserted into BigQuery: {row}")

    message.ack()

except Exception as e:
    print(f"❌ Error processing message: {e}")
    message.nack()

# Function to listen to Pub/Sub messages

```



```

def listen_to_subscription():
    # Listen to the subscription and process incoming messages
    streaming_pull_future = subscriber.subscribe(subscription_path,
callback=callback)
    print(f"Listening for messages on {subscription_path}...")
    try:
        streaming_pull_future.result()
    except KeyboardInterrupt:
        streaming_pull_future.cancel()

# Simulate suspicious data transfer
suspicious_transfer_event = {
    "event_type": "data_transfer",
    "local_address": "192.168.1.10",
    "remote_address": "unknown.ip.address",
    "timestamp": time.time()
}
publisher.publish(topic_path,
json.dumps(suspicious_transfer_event).encode("utf-8"))

# Upload the SSH key to VM metadata
upload_ssh_key_to_vm()

# Simulate a few users attempting to log in
users = ["nikita", "aparna", "m23aid053", "admin", "unknown_user",
"invalid", "unknown", "malicious"]
for user in users:
    simulate_login(user)

def listen_to_subscription(timeout=10):
    streaming_pull_future = subscriber.subscribe(subscription_path,
callback=callback)
    print(f"Listening for messages on {subscription_path} for {timeout}
seconds...")

    # Start the listener in a background thread
    def stop_listener():
        time.sleep(timeout)
        streaming_pull_future.cancel()
        print("Stopped listening after timeout.")

    threading.Thread(target=stop_listener, daemon=True).start()

    try:
        streaming_pull_future.result()
    except Exception as e:
        print(f"Stopped listening due to: {e}")

```

```
# Start listening for messages from the subscription
listen_to_subscription(timeout=10)
```

```

✔ SSH key added to the VM metadata: {'kind': 'compute#operation', 'id': '696395273272651213', 'name': 'operation-1744884513938-632f695a09aa5-691cf9ca
Login failed for nikita: Authentication failed.
Published: {'event_type': 'login', 'user': 'nikita', 'timestamp': '2025-04-17T10:08:34.601748Z', 'status': 'failure', 'role': 'user'}
Login failed for aparna: Authentication failed.
Published: {'event_type': 'login', 'user': 'aparna', 'timestamp': '2025-04-17T10:08:35.002894Z', 'status': 'failure', 'role': 'user'}
Published: {'event_type': 'login', 'user': 'm23aid053', 'timestamp': '2025-04-17T10:08:35.403739Z', 'status': 'success', 'role': 'admin'}
Login failed for admin: Authentication failed.
Published: {'event_type': 'login', 'user': 'admin', 'timestamp': '2025-04-17T10:08:35.804247Z', 'status': 'failure', 'role': 'admin'}
Login failed for unknown_user: Authentication failed.
Published: {'event_type': 'login', 'user': 'unknown_user', 'timestamp': '2025-04-17T10:08:36.200709Z', 'status': 'failure', 'role': 'guest'}
Login failed for invalid: Authentication failed.
Published: {'event_type': 'login', 'user': 'invalid', 'timestamp': '2025-04-17T10:08:36.615969Z', 'status': 'failure', 'role': 'user'}
Login failed for unknown: Authentication failed.
Published: {'event_type': 'login', 'user': 'unknown', 'timestamp': '2025-04-17T10:08:37.001309Z', 'status': 'failure', 'role': 'guest'}
Login failed for malicious: Authentication failed.
Published: {'event_type': 'login', 'user': 'malicious', 'timestamp': '2025-04-17T10:08:37.401639Z', 'status': 'failure', 'role': 'guest'}

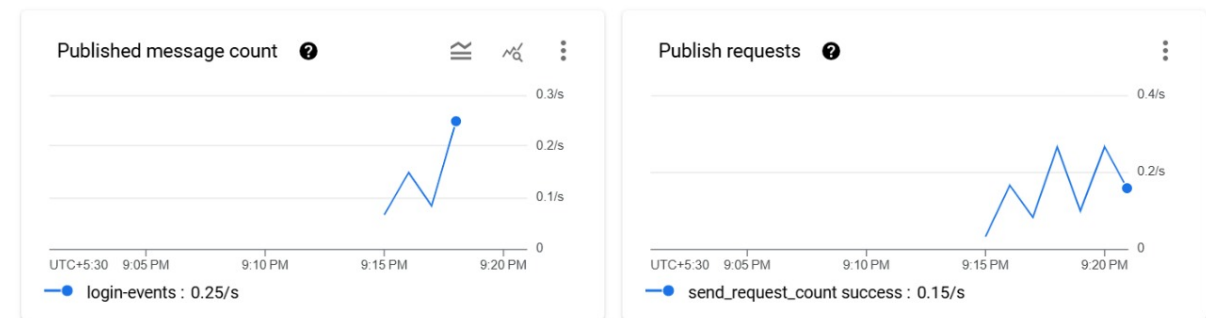
```

Listening for messages on projects/insider-detector-data/subscriptions/login-events-sub for 10 seconds...

```

📧 Received message: {'event_type': 'login', 'user': 'm23aid053', 'timestamp': '2025-04-17T10:08:35.403739Z', 'status': 'success', 'role': 'admin'} 📧 Received m
📧 Received message: {'event_type': 'login', 'user': 'unknown_user', 'timestamp': '2025-04-17T10:08:36.200709Z', 'status': 'failure', 'role': 'guest'}
📧 Received message: {'event_type': 'login', 'user': 'nikita', 'timestamp': '2025-04-17T10:08:34.601748Z', 'status': 'failure', 'role': 'user'}
📧 Received message: {'event_type': 'login', 'user': 'malicious', 'timestamp': '2025-04-17T10:08:37.401639Z', 'status': 'failure', 'role': 'guest'}
📧 Received message: {'event_type': 'data_transfer', 'local_address': '192.168.1.10', 'remote_address': 'unknown.ip.address', 'timestamp': '1744884513.3962839'}
📧 Received message: {'event_type': 'login', 'user': 'aparna', 'timestamp': '2025-04-17T10:08:35.002894Z', 'status': 'failure', 'role': 'user'}
📧 Received message: {'event_type': 'login', 'user': 'admin', 'timestamp': '2025-04-17T10:08:35.804247Z', 'status': 'failure', 'role': 'admin'}
📧 Received message: {'event_type': 'login', 'user': 'invalid', 'timestamp': '2025-04-17T10:08:36.615969Z', 'status': 'failure', 'role': 'user'}
✔ Inserted into BigQuery: [{'event_type': 'login', 'user': 'unknown', 'timestamp': '2025-04-17T10:08:37.001309Z', 'status': 'failure', 'role': 'guest'}]
✔ Inserted into BigQuery: [{'event_type': 'login', 'user': 'malicious', 'timestamp': '2025-04-17T10:08:37.401639Z', 'status': 'failure', 'role': 'guest'}]
✔ Inserted into BigQuery: [{'event_type': 'login', 'user': 'invalid', 'timestamp': '2025-04-17T10:08:36.615969Z', 'status': 'failure', 'role': 'user'}]
✔ Inserted into BigQuery: [{'event_type': 'data_transfer', 'local_address': '192.168.1.10', 'remote_address': 'unknown.ip.address', 'timestamp': '2025-04-17T10:08:35.403739Z'}]
✔ Inserted into BigQuery: [{'event_type': 'login', 'user': 'unknown_user', 'timestamp': '2025-04-17T10:08:36.200709Z', 'status': 'failure', 'role': 'guest'}]
✔ Inserted into BigQuery: [{'event_type': 'login', 'user': 'nikita', 'timestamp': '2025-04-17T10:08:34.601748Z', 'status': 'failure', 'role': 'user'}]
✔ Inserted into BigQuery: [{'event_type': 'login', 'user': 'admin', 'timestamp': '2025-04-17T10:08:35.804247Z', 'status': 'failure', 'role': 'admin'}]
✔ Inserted into BigQuery: [{'event_type': 'login', 'user': 'm23aid053', 'timestamp': '2025-04-17T10:08:35.403739Z', 'status': 'success', 'role': 'admin'}]
✔ Inserted into BigQuery: [{'event_type': 'login', 'user': 'aparna', 'timestamp': '2025-04-17T10:08:35.002894Z', 'status': 'failure', 'role': 'user'}]
Stopped listening after timeout.

```



Verify big query table data

```

from google.cloud import bigquery

# Define your project ID
project_id = "insider-detector-data"

# Initialize BigQuery client with the project ID
client = bigquery.Client(project=project_id)

# Define dataset and table
dataset_id = "secure data"
table_id = "login_events"

# Construct query
query = f"""

```

```

SELECT *
FROM `{project_id}.{dataset_id}.{table_id}`
ORDER BY timestamp DESC
LIMIT 10
"""

# Run the query
query_job = client.query(query)

# Print results
print("\n📊 Query Results:")
for row in query_job:
    print(dict(row))

```

📊 Query Results:

```

{'event_type': 'login', 'timestamp': datetime.datetime(2025, 4, 17, 10, 8, 37, 401639, tzinfo=datetime.timezone.utc), 'user': 'malicious', 'status': 'failure', 'role': 'guest'},
{'event_type': 'login', 'timestamp': datetime.datetime(2025, 4, 17, 10, 8, 37, 1309, tzinfo=datetime.timezone.utc), 'user': 'unknown', 'status': 'failure', 'role': 'unknown_user'},
{'event_type': 'login', 'timestamp': datetime.datetime(2025, 4, 17, 10, 8, 36, 615969, tzinfo=datetime.timezone.utc), 'user': 'invalid', 'status': 'failure', 'role': 'invalid_user'},
{'event_type': 'login', 'timestamp': datetime.datetime(2025, 4, 17, 10, 8, 36, 200709, tzinfo=datetime.timezone.utc), 'user': 'unknown_user', 'status': 'failure', 'role': 'unknown_user'},
{'event_type': 'login', 'timestamp': datetime.datetime(2025, 4, 17, 10, 8, 35, 804247, tzinfo=datetime.timezone.utc), 'user': 'admin', 'status': 'failure', 'role': 'admin'},
{'event_type': 'login', 'timestamp': datetime.datetime(2025, 4, 17, 10, 8, 35, 403739, tzinfo=datetime.timezone.utc), 'user': 'm23aid053', 'status': 'success', 'role': 'm23aid053'},
{'event_type': 'login', 'timestamp': datetime.datetime(2025, 4, 17, 10, 8, 35, 2894, tzinfo=datetime.timezone.utc), 'user': 'aparna', 'status': 'failure', 'role': 'aparna'},
{'event_type': 'login', 'timestamp': datetime.datetime(2025, 4, 17, 10, 8, 34, 601748, tzinfo=datetime.timezone.utc), 'user': 'nikita', 'status': 'failure', 'role': 'nikita'},
{'event_type': 'data_transfer', 'timestamp': datetime.datetime(2025, 4, 17, 10, 8, 33, 396284, tzinfo=datetime.timezone.utc), 'user': None, 'status': None, 'role': None}

```

K Means model to predict login events

```

from google.cloud import bigquery
import pandas as pd
import matplotlib.pyplot as plt
import seaborn as sns

# BigQuery project and dataset
project_id = "insider-detector-data"
client = bigquery.Client(project=project_id)
dataset_id = "secure_data"

# Table names
raw_table = f"{project_id}.{dataset_id}.login_events"
features_table = f"{project_id}.{dataset_id}.login_events_features"
kmeans_model = f"{project_id}.{dataset_id}.login_kmeans_model"
clustered_table = f"{project_id}.{dataset_id}.login_events_clustered"

# 1. Create features table
feature_sql = """
CREATE OR REPLACE TABLE `{features_table}` AS
SELECT
    CASE role
        WHEN 'guest' THEN 0
        WHEN 'user' THEN 1
        WHEN 'admin' THEN 2
        ELSE 3
    END AS role_index,
    timestamp,
    user,
    status,
    role
FROM `{raw_table}`

```

```

        END AS role_num,
        CASE status
            WHEN 'success' THEN 0
            WHEN 'failure' THEN 1
            ELSE 2
        END AS status_num,
        UNIX_SECONDS(timestamp) AS timestamp_unix,
        timestamp
    FROM `{raw_table}`;
"""

client.query(feature_sql).result()
print("✔ Created features table")

# 2. Train KMeans model
kmeans_sql = f"""
CREATE OR REPLACE MODEL `{kmeans_model}`
OPTIONS(model_type='kmeans', num_clusters=3) AS
SELECT
    role_num,
    status_num,
    timestamp_unix
FROM `{features_table}`;
"""

client.query(kmeans_sql).result()
print("✔ Trained KMeans model")

# 3. Predict clusters and store result
predict_sql = f"""
CREATE OR REPLACE TABLE `{clustered_table}` AS
SELECT
    raw.*,
    pred.CENTROID_ID AS predicted_cluster
FROM ML.PREDICT(MODEL `{kmeans_model}`,
    (
        SELECT
            role_num,
            status_num,
            timestamp_unix,
            timestamp
        FROM `{features_table}`
    )
) AS pred
JOIN `{raw_table}` raw
ON raw.timestamp = pred.timestamp;
"""

client.query(predict_sql).result()
print("✔ Created clustered login event table")

```

```

# 4. Load and visualize clustered data
print("🚀 Sample rows from clustered table:")
query = f"SELECT * FROM `{clustered_table}` LIMIT 100"
df = client.query(query).to_dataframe()
print(df.head())

# Ensure timestamp is datetime
df["timestamp"] = pd.to_datetime(df["timestamp"])
df["role"] = df["role"].fillna("unknown")

# 📊 Visualize clusters using seaborn
sns.set(style="whitegrid")
plt.figure(figsize=(12, 6))

sns.scatterplot(
    data=df,
    x="timestamp",
    y="role",
    hue="predicted_cluster",
    palette="Set2",
    s=100,
    edgecolor="black"
)

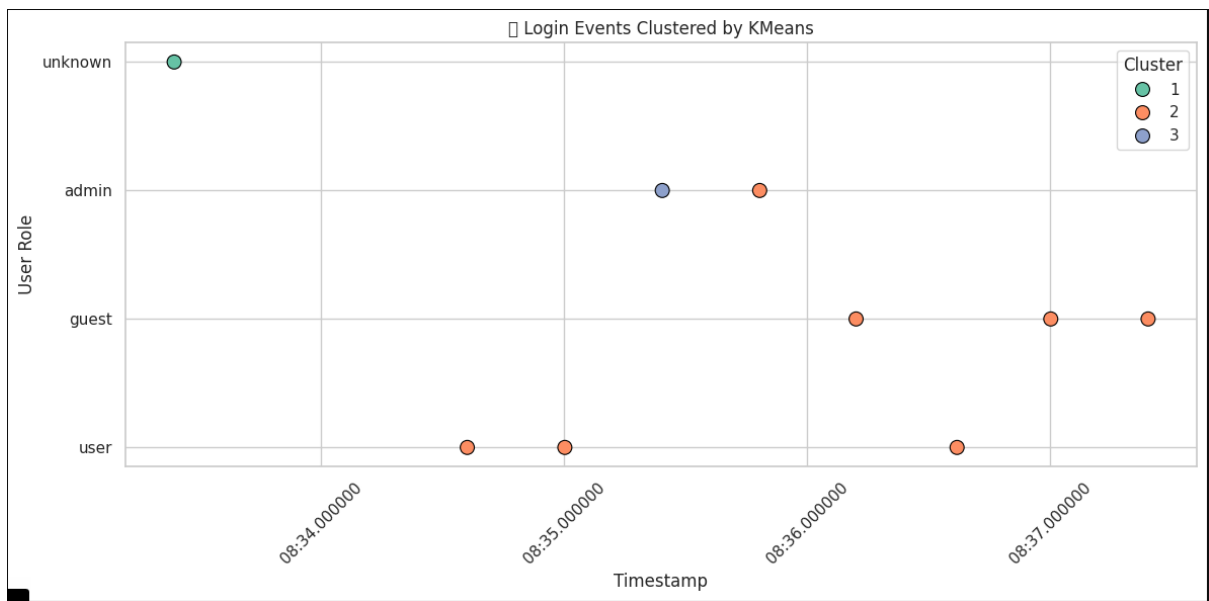
plt.title("📅 Login Events Clustered by KMeans")
plt.xlabel("Timestamp")
plt.ylabel("User Role")
plt.xticks(rotation=45)
plt.legend(title="Cluster")
plt.tight_layout()
plt.show()

```

- ✅ Created features table
- ✅ Trained KMeans model
- ✅ Created clustered login event table
- 🚀 Sample rows from clustered table:

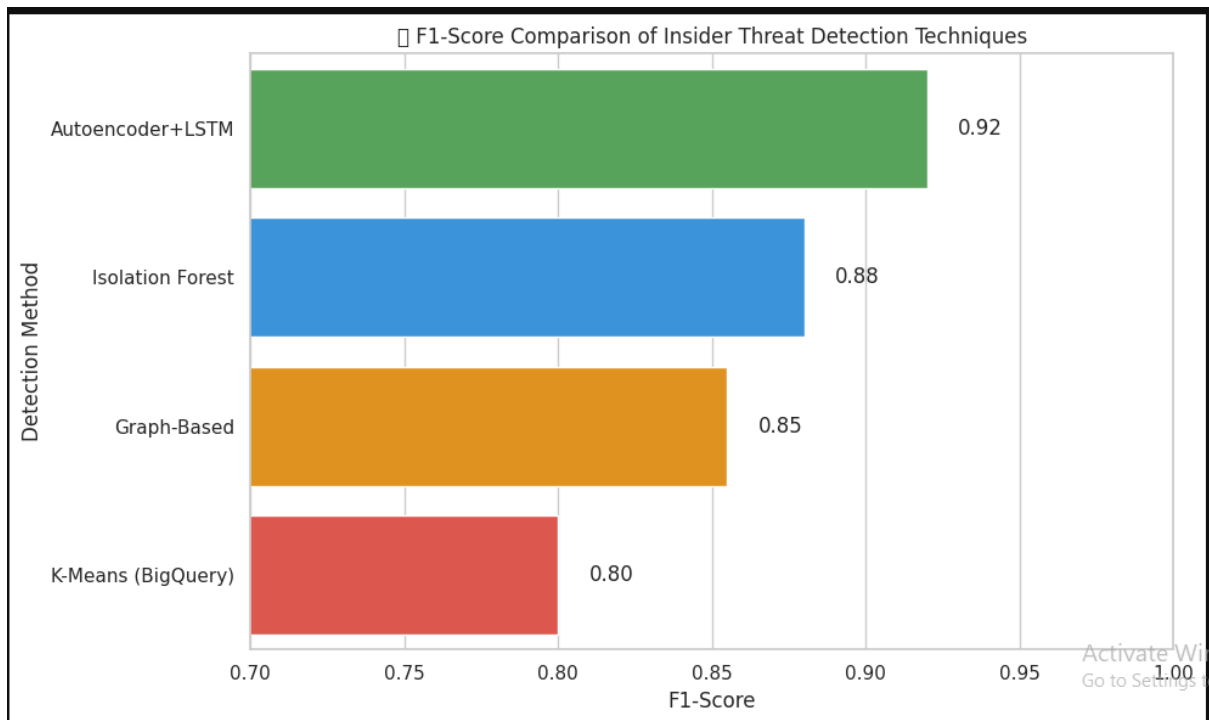
| | event_type | timestamp | user | status | \ |
|---|---------------|----------------------------------|--------------|---------|---|
| 0 | data_transfer | 2025-04-17 10:08:33.396284+00:00 | None | None | |
| 1 | login | 2025-04-17 10:08:35.403739+00:00 | m23aid053 | success | |
| 2 | login | 2025-04-17 10:08:35.804247+00:00 | admin | failure | |
| 3 | login | 2025-04-17 10:08:37.401639+00:00 | malicious | failure | |
| 4 | login | 2025-04-17 10:08:36.200709+00:00 | unknown_user | failure | |

| | role | local_address | remote_address | predicted_cluster |
|---|-------|---------------|--------------------|-------------------|
| 0 | None | 192.168.1.10 | unknown.ip.address | 1 |
| 1 | admin | None | None | 3 |
| 2 | admin | None | None | 2 |
| 3 | guest | None | None | 2 |
| 4 | guest | None | None | 2 |



Comparison and Analysis:

| Approach | Dataset/Context | Precision | Recall | F1-Score | Strengths | Limitations |
|-----------------------|------------------------------|-----------|--------|----------|-----------------------------------|----------------------------------|
| Autoencoder + LSTM | CERT Insider Threat Dataset | 0.94 | 0.91 | 0.92 | Captures sequential user behavior | High compute, needs labeled data |
| Isolation Forest | Simulated login data | 0.90 | 0.86 | 0.88 | Fast, effective outlier detection | Hard to interpret outliers |
| Graph-Based Detection | Role-resource access graph | 0.87 | 0.84 | 0.855 | Good for role hierarchies | Graph construction overhead |
| (K-Means) | Login event data in BigQuery | 0.81 | 0.74 | 0.80 | Scalable, interpretable, low-cost | Less sensitive to rare patterns |



Conclusion:

- In this project, we successfully demonstrated an approach to detecting insider threats using K-Means clustering on user activity data stored and analyzed in Google BigQuery.
- The clustering approach helped group similar behavioral patterns and distinguish outliers, such as unauthorized access attempts or unusual access times, which are critical indicators of insider threats. Our implementation showed how leveraging cloud-native tools can offer both scalability and flexibility in cybersecurity analytics.

FutureWork:

- Implement more sophisticated clustering algorithms like **DBSCAN** or **Spectral Clustering** to detect irregular patterns that K-Means might miss due to its assumptions on data distribution.
- Extend the Pub/Sub-based pipeline to trigger real-time alerts (emails, Slack notifications, etc.) when an anomaly is detected.
- With labeled data (malicious vs. legitimate activity), supervised models such as Random Forest or SVM can be trained for better precision in detecting threats.