

Chronicle-Sniffer

A Scalable Wireshark-to-SecOps Pipeline on GCP

Filippo Lucchesi

`filippo.lucchesi@studio.unibo.it`

Scalable and Reliable Services M
Università di Bologna

10 Giugno 2025

Indice

1	Introduzione	2
2	Architettura Generale	2
2.1	Motivazioni delle Scelte Architettureali	4
3	Logica di Conversione e Adattamenti per il Cloud	4
3.1	Evoluzione per l'Efficienza Cloud: da <code>json.loads</code> a <code>ijson</code>	5
3.2	Funzionamento Adattato dello script <code>json2udm_cloud.py</code>	5
4	Componente Sniffer (Edge)	7
5	Componente Processor (Cloud Run)	8
6	Infrastruttura come Codice (IaC) con Terraform	10
7	Valore Didattico ed Evoluzione del Progetto	11
8	Aspetti Operativi	12
9	Test e Validazione	12
10	Conclusioni e Sviluppi Futuri	13
11	Bibliografia e Sitografia	14

1 Introduzione

Questo report descrive l'architettura, l'implementazione e le caratteristiche di un progetto sviluppato per il corso di "Scalable and Reliable Services M" presso l'Università di Bologna.

Il sistema, denominato "Chronicle-Sniffer", realizza una pipeline di elaborazione dati di rete su Google Cloud Platform (GCP), con l'obiettivo primario di catturare pacchetti di rete (PCAP e PCAPNG), processarli in modo scalabile e affidabile, e trasformarli nel formato Unified Data Model (UDM) per successive analisi, tipicamente in contesti di sicurezza informatica, er esempio avvalendosi del Software Google SecOps.

Il progetto rappresenta una significativa evoluzione di un precedente lavoro concepito per il corso di Cybersecurity (Wireshark-to-Chronicle-Pipeline). Mentre la versione iniziale si concentrava sulla logica di conversione locale, questa iterazione si focalizza sulla creazione di un'infrastruttura cloud-nativa robusta, enfatizzando scalabilità, resilienza, gestibilità e osservabilità. Le tecnologie chiave impiegate includono servizi gestiti di GCP come Cloud Storage, Pub/Sub (con supporto per Dead-Letter Queue - DLQ), e Cloud Run per l'elaborazione serverless. L'intera infrastruttura è definita e gestita tramite Terraform, promuovendo l'approccio Infrastructure as Code (IaC), e include la definizione di una dashboard operativa personalizzata in Cloud Monitoring. I componenti software sono containerizzati con Docker. L'analisi del traffico di rete si avvale di TShark per la cattura e la conversione iniziale, mentre lo script Python `json2udm_cloud.py` è stato reingegnerizzato per un parsing JSON efficiente in streaming e una mappatura UDM più dettagliata.

L'obiettivo è dimostrare la capacità di costruire un servizio distribuito che non solo risponda a requisiti funzionali specifici, ma che sia anche progettato per operare in maniera efficiente, robusta, sicura e altamente osservabile in un ambiente cloud.

2 Architettura Generale

Il sistema Chronicle-Sniffer è progettato come una pipeline *event-driven*, che sfrutta diversi servizi GCP per orchestrare il flusso dei dati dalla cattura del traffico di rete fino alla sua trasformazione finale. L'architettura si articola attorno a componenti principali che interagiscono in modo asincrono, garantendo disaccoppiamento e resilienza.

- **Sniffer (On-Premises/Edge):** Un container Docker, eseguito su una macchina fisica o virtuale nell'ambiente da monitorare. Utilizza `tshark` per catturare il traffico di rete (supportando formati `.pcap` e `.pcapng`), gestisce la rotazione automatica dei file di cattura e, una volta completato un file, lo carica su un bucket Google Cloud Storage (GCS) (denominato `incoming-pcaps`). Immediatamente dopo, pubblica una notifica contenente il nome del file su un topic Google Cloud Pub/Sub. Lo sniffer include anche una funzionalità di *heartbeat* per monitorarne lo stato.
- **Google Cloud Pub/Sub:** Agisce come *message broker*. Un topic principale riceve le notifiche dallo sniffer. Una sottoscrizione *push*, configurata con autenticazione OIDC sicura, inoltra i messaggi all'endpoint del servizio Cloud Run. È previsto un topic Dead-Letter Queue (DLQ) per gestire i messaggi che non possono essere processati correttamente dopo vari tentativi.
- **Google Cloud Storage (GCS):** Due bucket distinti sono utilizzati:
 - `incoming-pcaps`: Funge da area di *staging* per i file di cattura raw. Come visibile nello screenshot seguente (Figura 2), questo bucket riceve i file `.pcap` direttamente dallo sniffer.

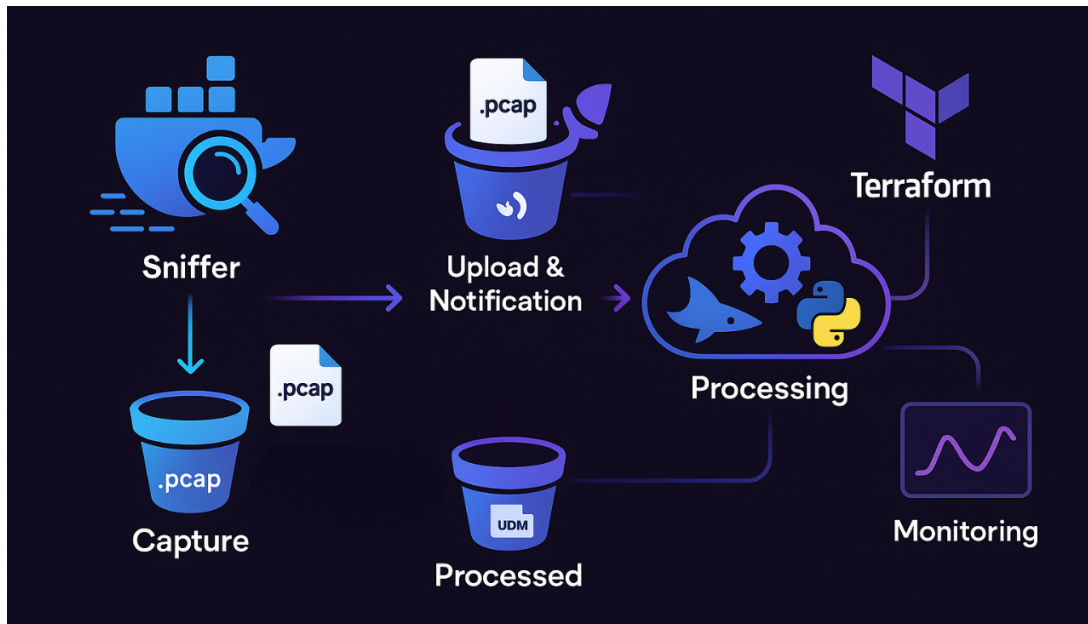


Figura 1: Diagramma dell'architettura generale del sistema Chronicle-Sniffer.

chronicle-sniffer-incoming-pcaps

Località: eu (più regioni nell'Unione europea) | Classe di archiviazione: Standard | Accesso pubblico: Non pubblica | Protezione: Eliminazione temporanea

Oggetti | Configurazione | Autorizzazioni | Protezione | Ciclo di vita | Osservabilità | Nuova | Report sugli spazi pubblicitari | Operazioni

Bucket > chronicle-sniffer-incoming-pcaps

Crea cartella | Carica | Trasferimento dei dati | Altri servizi

Filtra solo per prefisso nome | Filtra | Filtra oggetti e cartelle | Mostra: Solo oggetti attivi

<input type="checkbox"/>	Nome	Dimensioni	Tipo	Data creazione	Classe di archiviazione	Ultima modifica	
<input type="checkbox"/>	capture_00001_20250521111057...	17,9 kB	application/octet-stream	21 mag 2025, 13:11:05	Standard	21 mag 2025, 13:11:05	⬇ ⋮
<input type="checkbox"/>	capture_00002_20250521111158...	14,9 kB	application/octet-stream	21 mag 2025, 13:12:09	Standard	21 mag 2025, 13:12:09	⬇ ⋮
<input type="checkbox"/>	capture_00003_20250521111258...	20,4 kB	application/octet-stream	21 mag 2025, 13:13:03	Standard	21 mag 2025, 13:13:03	⬇ ⋮
<input type="checkbox"/>	capture_00004_20250521111358...	12 kB	application/octet-stream	21 mag 2025, 13:14:06	Standard	21 mag 2025, 13:14:06	⬇ ⋮

Figura 2: Contenuto del bucket GCS chronicle-sniffer-incoming-pcaps.

- **processed-udm**: Archivia i file JSON finali contenenti gli eventi di rete trasformati nel formato UDM. Lo screenshot successivo (Figura 3) mostra un esempio dei file `.udm.json` prodotti dal Processor.

chronicle-sniffer-processed-udm

Località: eu (più regioni nell'Unione europea) | Classe di archiviazione: Standard | Accesso pubblico: Non pubblica | Protezione: Eliminazione temporanea

Oggetti | Configurazione | Autorizzazioni | Protezione | Ciclo di vita | Osservabilità | Nuova | Report sugli spazi pubblicitari | Operazioni

Bucket > chronicle-sniffer-processed-udm

Crea cartella | Carica | Trasferimento dei dati | Altri servizi

Filtra solo per prefisso nome | Filtra | Filtra oggetti e cartelle | Mostra: Solo oggetti attivi

<input type="checkbox"/>	Nome	Dimensioni	Tipo	Data creazione	Classe di archiviazione	Ultima modifica	Acc
<input type="checkbox"/>	capture_00001_20250521111057...	76 kB	application/json	21 mag 2025, 13:11:14	Standard	21 mag 2025, 13:11:14	Nor ⬇ ⋮
<input type="checkbox"/>	capture_00002_20250521111158...	64,9 kB	application/json	21 mag 2025, 13:12:15	Standard	21 mag 2025, 13:12:15	Nor ⬇ ⋮
<input type="checkbox"/>	capture_00003_20250521111258...	92,2 kB	application/json	21 mag 2025, 13:13:06	Standard	21 mag 2025, 13:13:06	Nor ⬇ ⋮

Figura 3: Contenuto del bucket GCS chronicle-sniffer-processed-udm.

- **Processore (Google Cloud Run):** Un servizio serverless che esegue un container Docker. Attivato dalle notifiche Pub/Sub, il servizio:
 1. Scarica il file di cattura specificato da GCS.
 2. Utilizza `tshark` (installato nel container) per convertire il file di cattura in un formato JSON strutturato.
 3. Esegue lo script Python `json2udm_cloud.py`, che ora utilizza la libreria `ijson` per il parsing in streaming del (potenzialmente grande) output JSON di `tshark`. Questo approccio è cruciale per l'efficienza della memoria in Cloud Run. Lo script mappa ogni pacchetto nel formato UDM.
 4. Carica il file UDM JSON risultante nel bucket GCS `processed-udm`.

Il servizio è progettato per scalare automaticamente e include *health probes* per il monitoraggio da parte di Cloud Run.

L'intera infrastruttura, inclusi i bucket, il topic Pub/Sub con DLQ, il servizio Cloud Run, i Service Account, i permessi IAM, le metriche basate su log e una dashboard operativa personalizzata, è definita e gestita tramite Terraform.

2.1 Motivazioni delle Scelte Architettureali

L'architettura è stata definita per massimizzare scalabilità, affidabilità, gestibilità, efficienza dei costi e sicurezza:

- **Scalabilità:** Cloud Run e Pub/Sub scalano automaticamente per gestire carichi variabili. L'elaborazione in streaming nel Processor previene colli di bottiglia legati all'esaurimento della memoria.
- **Affidabilità e Resilienza:** Il disaccoppiamento tramite Pub/Sub, unito al supporto DLQ, garantisce che i messaggi non vengano persi e che i fallimenti siano gestiti in modo isolato. Cloud Run gestisce la disponibilità delle istanze del processor.
- **Gestibilità (IaC):** Terraform permette una gestione dichiarativa, versionabile e automatizzata dell'intera infrastruttura, inclusa la dashboard di monitoraggio.
- **Efficienza dei costi:** I servizi serverless e lo scaling a zero di Cloud Run ottimizzano i costi.
- **Sicurezza:** Adozione del principio del *least privilege* per i Service Account, autenticazione OIDC tra Pub/Sub e Cloud Run, e gestione sicura delle chiavi per lo sniffer on-premises.
- **Osservabilità:** Integrazione profonda con Cloud Logging e Cloud Monitoring, con metriche personalizzate basate su log (LBM) e una dashboard operativa dedicata per una visione completa dello stato e delle prestazioni della pipeline.

3 Logica di Conversione e Adattamenti per il Cloud

Il cuore analitico del sistema è lo script Python `json2udm_cloud.py`, eseguito dal componente Processor. Questo script è responsabile della trasformazione dei dati dei pacchetti, precedentemente convertiti in JSON da `tshark`, nel formato standardizzato Unified Data Model (UDM). È importante sottolineare che la logica dettagliata di parsing dei singoli campi PCAP e la loro mappatura iniziale ai concetti UDM sono state approfondite nel precedente progetto di Cybersecurity. Il focus di questo progetto è stato l'adattamento di tale logica per un ambiente cloud efficiente e la costruzione di un'infrastruttura scalabile e resiliente attorno ad essa.

3.1 Evoluzione per l'Efficienza Cloud: da json.loads a ijson

Nel progetto di Cybersecurity, la conversione JSON → UDM avveniva caricando l'intero output JSON di tshark in memoria tramite `json.loads()`. Sebbene funzionale per file di dimensioni moderate, questo approccio diventa proibitivo in un ambiente serverless come Cloud Run quando si trattano file PCAP di grandi dimensioni, che possono generare output JSON di svariati gigabyte, portando a errori di Out of Memory (OOM).

La modifica più significativa in `json2udm_cloud.py` è l'adozione della libreria `ijson`. Questa libreria permette il parsing JSON in streaming: invece di caricare l'intero documento, `ijson` processa il file JSON in modo incrementale, leggendo e interpretando gli oggetti JSON (in questo caso, i singoli pacchetti) uno alla volta.

```

1 # Lo script apre il file JSON prodotto da tshark in modalita binaria ('rb').
2 # ijson.items(f_json, 'item') crea un iteratore che assume che il file JSON
3 # sia un array di oggetti alla radice. 'item' istruisce ijson a restituire
4 # (yield) ciascun elemento di questo array principale, uno alla volta.
5 # Questo permette di processare pacchetti individualmente senza caricare
6 # l'intero array di pacchetti (potenzialmente enorme) in memoria.
7 try:
8     with open(json_file_path, 'rb') as f_json: # Modalita binaria per ijson
9         json_packet_iterator = ijson.items(f_json, 'item')
10        for packet_data_dict in json_packet_iterator:
11            # Ogni 'packet_data_dict' e' un dizionario Python
12            # rappresentante un singolo pacchetto.
13            udm_event = convert_single_packet_to_udm(packet_data_dict)
14            udm_events_list.append(udm_event)
15            processed_packet_count += 1
16            # ... (logica per contare errori specifici per pacchetto)
17 except ijson.JSONError as e_ijson:
18     # Gestisce errori se il file JSON non e' ben formato
19     # o non e' un array alla radice.
20     logging.error(f"ijson.JSONError while parsing streaming JSON from {json_file_path}: {e_ijson}...")
21 return []

```

Listing 1: Estratto da `json2udm_cloud.py`: Utilizzo di `ijson` per lo streaming.

Questo approccio riduce drasticamente il picco di utilizzo della memoria, rendendo il processo di conversione adatto anche per istanze Cloud Run con risorse limitate e capace di gestire file di input molto grandi.

3.2 Funzionamento Adattato dello script json2udm_cloud.py

Oltre al parsing in streaming, lo script `json2udm_cloud.py` è stato migliorato per:

- **Robustezza:** Ogni pacchetto in input produce un evento UDM, anche se si tratta di un evento minimo che segnala un errore di elaborazione per quel pacchetto. La conversione dei timestamp (`convert_timestamp_robust`) è più resiliente, con fallback all'ora corrente in caso di formati imprevisti, garantendo la conformità UDM.

```

1 # Questa funzione tenta diverse strategie per parsare il timestamp fornito da
2 # tshark.
3 # Inizia con il formato piu comune (con microsecondi), poi tenta formati
4 # alternativi
5 # (es. senza microsecondi, pulendo stringhe di fuso orario comuni).
6 # Se tutte le conversioni falliscono, logga un warning e restituisce l'ora corrente
7 # in formato ISO 8601 UTC, assicurando che il campo 'event_timestamp' dell'UDM
8 # sia sempre presente e valido.
9 def convert_timestamp_robust(timestamp_str):
10     if not timestamp_str:
11         logging.warning(f"Timestamp string is missing or empty. Using current time as fallback.")
12         return datetime.now(timezone.utc).isoformat(timespec='microseconds')
13     replace('+00:00', 'Z')
14     try:

```

```

12     # Tentativo primario: formato standard di tshark con microsecondi
13     dt_naive = datetime.strptime(timestamp_str[:26], "%b %d, %Y %H:%M:%S.%f")
14     except ValueError:
15         # Fallback: prova a pulire stringhe di fuso orario e a parsare senza
           microsecondi
16         try:
17             cleaned_ts = timestamp_str.split(" UTC")[0]
18             # Altri tentativi di pulizia potrebbero essere aggiunti qui
19             cleaned_ts = cleaned_ts.split(" Central European Summer Time")[0].strip
           ()
20             dt_naive = datetime.strptime(cleaned_ts, "%b %d, %Y %H:%M:%S")
21         except ValueError as e_fallback: # e_fallback non usato, rimosso per
           pulizia warning
22             logging.warning(f"Error converting timestamp '{timestamp_str}' ...
           Using current time.")
23             return datetime.now(timezone.utc).isoformat(timespec='microseconds').
           replace('+00:00', 'Z')
24
25     # Assicura che il datetime sia 'aware' (consapevole del fuso orario) e in UTC.
26     dt_aware = dt_naive.replace(tzinfo=timezone.utc)
27     # Formatta in ISO 8601, assicurando la 'Z' per Zulu time (UTC).
28     iso_timestamp = dt_aware.isoformat(timespec='microseconds').replace('+00:00', '
           Z')
29     return iso_timestamp
30

```

Listing 2: Estratto da json2udm_cloud.py: Funzione convert_timestamp_robust.

- **Struttura UDM Allineata:** La struttura degli eventi UDM generati è stata raffinata per aderire più strettamente alle aspettative di sistemi come Chronicle, con sezioni `metadata`, `principal`, `target`, `network` (che può contenere `application_protocol_data` per HTTP, DNS, TLS, annidato sotto `network`) e `additional` ben definite.

```

1 # ... (all'interno di convert_single_packet_to_udm)
2 udm_payload = {
3     "metadata": {
4         "event_timestamp": event_timestamp,
5         "product_name": "Wireshark TShark",
6         "vendor_name": "Wireshark",
7         "event_type": event_type, # Determinato dinamicamente
8         "description": f"Packet capture. Protocols: {frame.get('frame.protocols', '
           N/A')}]..."
9     }
10 }
11 # Le sezioni principal, target, network, about, application_protocol_data,
           additional
12 # vengono populate e aggiunte a udm_payload solo se contengono dati effettivi.
13 # Esempio per la sezione 'principal':
14 cleaned_principal = clean_none_values(udm_principal) # clean_none_values rimuove
           chiavi con valore None
15 if cleaned_principal: udm_payload["principal"] = cleaned_principal
16 # ... logica simile per le altre sezioni ...
17 return {"event": udm_payload} # L'evento UDM finale e' incapsulato in un oggetto "
           event"
18

```

Listing 3: Estratto da json2udm_cloud.py: Struttura di un evento UDM.

- **Estrazione Dati Migliorata:** Utilizzo di funzioni helper come `get_nested_value` per un'estrazione sicura dei valori da strutture JSON potenzialmente complesse (evitando `KeyError`) e `extract_values_from_tshark_section` per campi multi-valore (es. query DNS).
- **Logging Dettagliato:** Lo script produce log informativi, inclusi contatori di pacchetti processati ed errori (`UDM_PACKETS_PROCESSED`, `UDM_PACKET_ERRORS`), utili per il monitoraggio tramite le metriche basate su log (LBM).

Il processo di conversione per ogni pacchetto implica l'estrazione dei layer protocollari (frame, eth, ip, tcp, udp, dns, http, tls, etc.), la mappatura dei campi rilevanti nelle sezioni UDM appropriate, e la normalizzazione dei dati. La selezione dei campi mappati rimane focalizzata sulla rilevazione di minacce e anomalie di rete.

4 Componente Sniffer (Edge)

Lo Sniffer è l'agente di cattura della pipeline, progettato per operare in ambienti on-premises o su dispositivi edge. È implementato come un container Docker leggero basato su Alpine Linux, contenente **tshark** e gli strumenti **gcloud**.

Il suo comportamento è orchestrato dallo script `sniffer_entrypoint.sh`. All'avvio, lo script valida le variabili d'ambiente (ID progetto, bucket, topic, e un `SNIFFER_ID` univoco), si autentica con GCP, e rileva automaticamente un'interfaccia di rete attiva (escludendo interfacce virtuali o di loopback). Avvia quindi `tshark` in background, configurato per la rotazione dei file di cattura (supportando estensioni `.pcap` e `.pcapng` tramite il pattern `*.pcap*`) in base a dimensione o tempo (configurabile tramite la variabile `ROTATE`).

```
1 # La variabile $INTERFACE viene determinata automaticamente o presa dall'ambiente.
2 # $ROTATE definisce i criteri di rotazione (es., "-b filesize:10240 -b duration:60").
3 # $LIMITS puo' contenere filtri di cattura aggiuntivi.
4 # I file vengono scritti nella directory $CAPTURE_DIR con un nome base $FILENAME_BASE.
5 echo "(ID: $SNIFFER_ID) Starting tshark capture..."
6 tshark $INTERFACE $ROTATE $LIMITS -w "$CAPTURE_DIR/$FILENAME_BASE.pcap" &
7 TSHARK_PID=$! # Salva il PID di tshark per monitoraggio e shutdown.
8 echo "(ID: $SNIFFER_ID) tshark started with PID $TSHARK_PID"
```

Listing 4: Estratto da `sniffer_entrpoint.sh`: Avvio di `tshark`.

Un loop di monitoraggio verifica costantemente la presenza di file completati (non più attivamente scritti da `tshark`, identificati tramite `ls`). Per ogni file completato, lo script ne logga la dimensione (`PCAP_SIZE_BYTES`), lo carica su GCS, pubblica una notifica con il nome del file su Pub/Sub, e infine rimuove il file locale. L'output dei log dello sniffer, come mostrato nello screenshot (Figura 4), evidenzia queste operazioni.

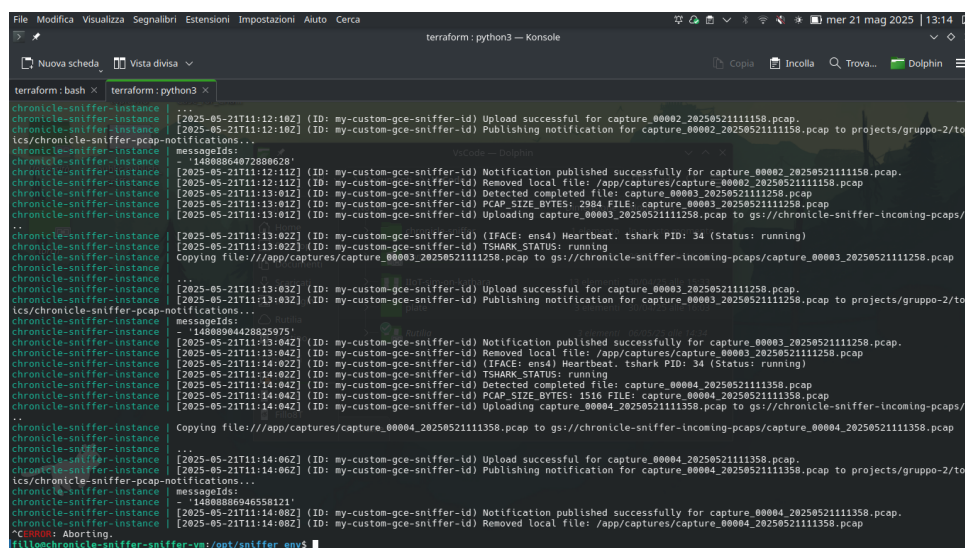


Figura 4: Log del container Sniffer durante l’upload e la notifica.

```
1 # Questo blocco viene eseguito per ogni file PCAP completato.
2 # "$pcap_file_path" e' il percorso del file locale.
3 # "$current_pcap_file_basename" e' solo il nome del file.
```



```

4 echo "[$(date)] (ID: $SNIFFER_ID) Uploading $current_pcap_file_basename to gs://${INCOMING_BUCKET}/..."
5 if gcloud storage cp "$pcap_file_path" "gs://${INCOMING_BUCKET}/" --project "$GCP_PROJECT_ID"; then
6     echo "[$(date)] (ID: $SNIFFER_ID) Upload successful for $current_pcap_file_basename."
7
8     echo "[$(date)] (ID: $SNIFFER_ID) Publishing notification for $current_pcap_file_basename to ${PUBSUB_TOPIC_ID}..."
9     if gcloud pubsub topics publish "$PUBSUB_TOPIC_ID" --message "$current_pcap_file_basename" --project "$GCP_PROJECT_ID"; then
10         echo "[$(date)] (ID: $SNIFFER_ID) Notification published successfully for $current_pcap_file_basename."
11         processed_files+="$current_pcap_file_basename" # Aggiunge alla lista dei file processati.
12         rm "$pcap_file_path" # Rimuove il file locale.
13         echo "[$(date)] (ID: $SNIFFER_ID) Removed local file: $pcap_file_path"
14     else
15         # Errore nella pubblicazione Pub/Sub, il file non viene rimosso e verra' ritentato.
16         echo "[$(date)] (ID: $SNIFFER_ID) Error: Failed to publish notification for $current_pcap_file_basename. Will retry."
17     fi
18 else
19     # Errore nell'upload GCS, il file non viene rimosso e verra' ritentato.
20     echo "[$(date)] (ID: $SNIFFER_ID) Error: Failed to upload $current_pcap_file_basename to GCS. Will retry."
21 fi

```

Listing 5: Estratto da `sniffer_entrypoint.sh`: Upload a GCS e notifica a Pub/Sub.

Lo script include una funzione di *heartbeat* che logga periodicamente lo stato di `tshark` (`TSHARK_STATUS`: `running/stopped`) e l'ID dello sniffer, facilitando il monitoraggio remoto della sua attività. Gestisce anche i segnali `SIGTERM` e `SIGINT` per un arresto pulito di `tshark`. La configurazione per l'avvio locale tramite `docker-compose` è fornita per facilitare test e sviluppo.

Un aspetto cruciale dello sniffer è il suo basso impatto sulle risorse del sistema edge. Come evidenziato dall'output del comando `docker stats` (vedi Figura 5), il container dello sniffer opera con un utilizzo minimo di CPU e memoria (nell'esempio, circa 0.05% CPU e ~100MB di RAM), rendendolo facilmente integrabile anche in ambienti con risorse limitate senza impattare significativamente le prestazioni di altri servizi. Questa efficienza è fondamentale per un'adozione su larga scala in contesti aziendali.

CONTAINER ID	NAME	CPU %	MEM USAGE / LIMIT	MEM %	NET I/O	BLOCK I/O	PIDS
75591d42c523	chronicle-sniffer-instance	0.05%	99.79MiB / 975.6MiB	10.23%	0B / 0B	44.4MB / 6MB	7

Figura 5: Output di `docker stats` per il container Sniffer.

5 Componente Processor (Cloud Run)

Il Processor, eseguito su Cloud Run, è il cuore dell'elaborazione. L'applicazione Flask `processor_app.py` riceve le notifiche Pub/Sub. Una novità importante è la verifica "lazy" dell'accessibilità ai bucket GCS all'avvio dell'istanza o alla prima richiesta, per gestire meglio eventuali ritardi nella propagazione dei permessi IAM post-deployment.

Per ogni notifica, l'applicazione scarica il file di cattura indicato, lo elabora in una directory temporanea, e orchestra le due fasi di conversione:

1. PCAP → JSON: Esegue `tshark -T json` come subprocess.
2. JSON → UDM: Esegue lo script `json2udm_cloud.py` (anch'esso come subprocess) sull'output JSON di `tshark`.

I log del servizio Cloud Run, visibili nello screenshot (Figura 6), mostrano le fasi di conversione e l'eventuale upload del file UDM risultante.

2025-05-21 13:11:13.963	2025-05-21 11:11:13.950	INFO	Successfully converted 74 packets from JSON to UDM format for file capture_00001_20250521111057.pcap.json.
2025-05-21 13:11:13.963	2025-05-21 11:11:13.950	INFO	UDM_PACKETS_PROCESSED: 74 FILE: capture_00001_20250521111057.pcap.json
2025-05-21 13:11:13.963	2025-05-21 11:11:13.954	INFO	Successfully wrote 74 UDM events to /tmp/tmp786ka/capture_00001_20250521111057.udm.json
2025-05-21 13:11:13.963	2025-05-21 11:11:13.964	INFO	Uploading /tmp/tmp786ka/capture_00001_20250521111057.udm.json to gs://chronicle-sniffer-processed-udm/capture_00001_20250521111057.udm.json
2025-05-21 13:11:14.330	2025-05-21 11:11:14.330	INFO	Upload complete for capture_00001_20250521111057.udm.json.
2025-05-21 13:11:14.330	2025-05-21 11:11:14.330	INFO	PROCESSING_DURATION_SECONDS: 2.467 FILE: capture_00001_20250521111057.pcap
2025-05-21 13:11:14.330	2025-05-21 11:11:14.330	INFO	Successfully processed capture_00001_20250521111057.pcap

Figura 6: Log del servizio Cloud Run (Processor) durante l'elaborazione.

```

1 # Questa logica viene eseguita all'interno di un blocco try...except per la gestione
  degli errori.
2 # local_pcap_path, local_json_path, local_udm_path sono percorsi in una directory
  temporanea.
3 try:
4     # ... (download del file pcap da GCS) ...
5
6     logging.info(f"Converting {local_pcap_path} to JSON...")
7     tshark_command = ["tshark", "-r", local_pcap_path, "-T", "json"]
8     # L'output JSON di tshark viene scritto direttamente su file.
9     with open(local_json_path, "w") as json_file:
10         # subprocess.run esegue il comando. check=True fa sì che venga sollevata
11         # un'eccezione CalledProcessError se tshark restituisce un exit code non zero.
12         process = subprocess.run(tshark_command, stdout=json_file, stderr=subprocess.PIPE
13         , text=True, check=True)
14         logging.info(f"tshark conversion successful: {local_json_path}")
15         if process.stderr: logging.warning(f"tshark stderr: {process.stderr.strip()}")
16
17         logging.info(f"Converting {local_json_path} to UDM: {local_udm_path}")
18         udm_script_command = ["python3", "/app/json2udm_cloud.py", local_json_path,
19         local_udm_path]
20         # Anche qui, check=True per catturare errori dallo script di conversione UDM.
21         # capture_output=True per raccogliere stdout e stderr dello script.
22         process = subprocess.run(udm_script_command, capture_output=True, text=True, check=
23         True)
24         logging.info(f"UDM conversion script done for {pcap_filename}.")
25         if process.stdout: logging.info(f"json2udm_cloud.py stdout: {process.stdout.strip()}")
26         if process.stderr: logging.warning(f"json2udm_cloud.py stderr: {process.stderr.strip
27         ()}")
28
29         # ... (upload del file UDM a GCS) ...
30
31 except subprocess.CalledProcessError as e:
32     # Gestisce errori specifici dei subprocess (tshark o json2udm_cloud.py).
33     error_message = e.stderr.strip() if e.stderr else e.stdout.strip()
34     logging.error(f"Subprocess error: CMD: {' '.join(e.cmd)} ERR: {error_message}",
35     exc_info=False)
36     return "Internal Server Error during processing step.", 500 # Segnala a Pub/Sub di
37     ritentare.
38 # ... (altri blocchi except per errori GCS, etc.)

```

Listing 6: Estratto da processor_app.py: Chiamata ai subprocess per tshark e json2udm_cloud.py.

L'applicazione logga la durata dell'elaborazione (PROCESSING_DURATION_SECONDS) per ogni file e gestisce gli errori restituendo codici HTTP appropriati a Pub/Sub per controllare i tentativi di acknowledgement (ack) o negative acknowledgement (nack).

6 Infrastruttura come Codice (IaC) con Terraform

L'intera infrastruttura GCP è gestita tramite Terraform. La configurazione è modulare e include:

- Service Account dedicati con permessi minimi.
- Bucket GCS per i file di input e output.
- Topic Pub/Sub principale e un Dead-Letter Topic (DLQ), con una sottoscrizione *push* verso Cloud Run configurata con autenticazione OIDC e una politica di *dead-lettering*.

```

1 resource "google_pubsub_subscription" "processor_subscription" {
2   project      = var.gcp_project_id
3   name         = "${var.base_name}-processor-sub"
4   topic        = module.pubsub_topic.topic_id # Topic principale
5   ack_deadline_seconds = 600 # Tempo per Cloud Run per processare
6
7   push_config {
8     push_endpoint = module.cloudrun_processor.service_url # URL del servizio Cloud
      Run
9     # Configurazione OIDC per l'autenticazione sicura
10    dynamic "oidc_token" {
11      for_each = !var.allow_unauthenticated_invocations ? [1] : []
12      content {
13        service_account_email = google_service_account.cloud_run_sa.email
14        audience              = module.cloudrun_processor.service_url
15      }
16    }
17  }
18
19  # Politica per inviare messaggi non elaborabili al DLQ
20  dead_letter_policy {
21    dead_letter_topic      = module.pubsub_topic.dlq_topic_id # Topic DLQ
22    max_delivery_attempts = 5 # Numero massimo di tentativi prima del DLQ
23  }
24  # ... (depends_on per gestione dipendenze)
25 }
26
```

Listing 7: Estratto da `terraform/main.tf`: Definizione della sottoscrizione Pub/Sub.

- Servizio Cloud Run per il processor, con configurazione di immagine, variabili d'ambiente, risorse (CPU, memoria impostata a 2Gi per gestire meglio *tshark*), concorrenza e *health probes*.

```

1 resource "google_cloud_run_v2_service" "processor" {
2   project = var.project_id
3   name    = var.service_name
4   location = var.region
5
6   template {
7     service_account = var.service_account_email
8     max_instance_request_concurrency = var.max_concurrency
9     timeout = "600s" # Timeout per richiesta
10    containers {
11      image = var.image_uri
12      ports { container_port = 8080 } # Porta interna del container
13      dynamic "env" { # Passaggio dinamico delle variabili d'ambiente
14        for_each = var.env_vars
15        content {
16          name = env.key
17          value = env.value
18        }
19      }
20    }
21    resources { # Limiti di CPU e memoria per istanza
22      limits = {
23        cpu = var.cpu_limit
24        memory = var.memory_limit # Es. "2Gi"
25      }
26    }
27  }
28}
```

```

24     }
25   }
26   # Probe per monitorare la salute dell'applicazione
27   startup_probe { http_get { path = "/" } # ... (configurazione dettagliata
    omessa per brevità) }
28   liveness_probe { http_get { path = "/" } # ... (configurazione dettagliata
    omessa per brevità) }
29 }
30 }
31 # ... (configurazione traffico omessa)
32 }
33

```

Listing 8: Estratto da terraform/modules/cloudrun_processor/main.tf

- **Metriche Basate su Log (LBMs):** Definite in terraform/main.tf per tracciare eventi chiave come *heartbeat* dello sniffer, upload di PCAP, errori di pubblicazione Pub/Sub, successo/fallimento delle conversioni *tshark*, pacchetti UDM processati/errori, e latenza di elaborazione. Queste metriche alimentano la dashboard operativa.
- **Dashboard Operativa:** Un file JSON (terraform/dashboards/main_operational_dashboard.json) definisce una dashboard completa in Cloud Monitoring utilizzando Monitoring Query Language (MQL). Questa dashboard visualizza le LBMs e le metriche standard dei servizi GCP, offrendo una visione centralizzata dello stato della pipeline. (Figure 7).

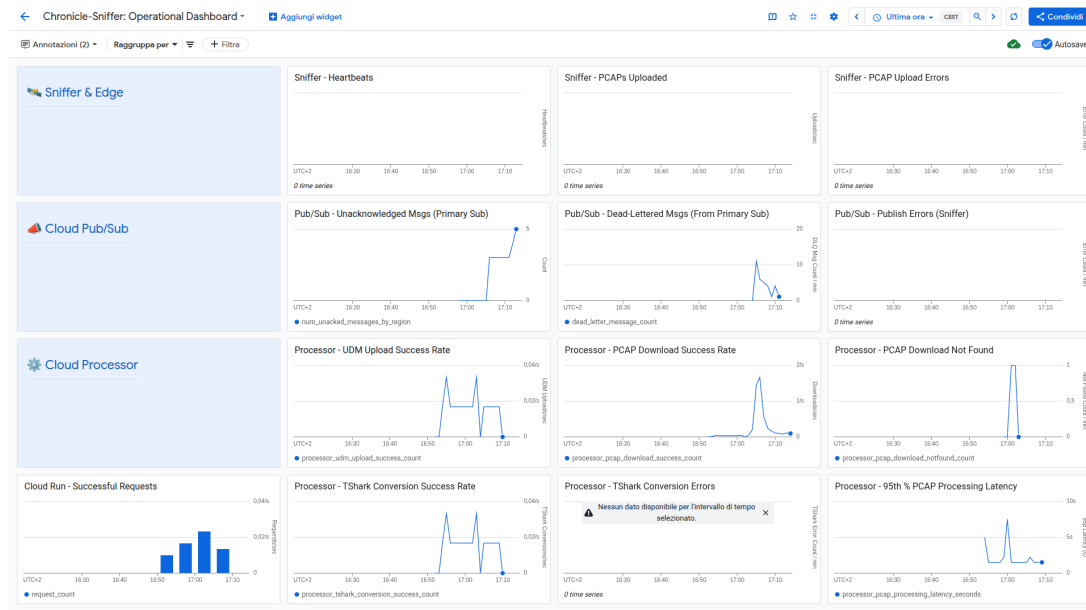


Figura 7: Dashboard Operativa in Cloud Monitoring.

- **Politica di Alert (Esempio):** Una politica di alert di esempio per gli sniffer inattivi (basata sulla metrica di *heartbeat*) dimostra come estendere l'osservabilità.
- **VM di Test (Opzionale):** Un modulo per creare una VM GCE che simula un ambiente on-prem, con uno script di avvio che prepara l'ambiente Docker Compose per lo sniffer.

7 Valore Didattico ed Evoluzione del Progetto

Questo progetto si distingue per la sua evoluzione da una soluzione locale a un'architettura cloud-nativa sofisticata. Il passaggio chiave è stato l'offloading dell'elaborazione intensiva al

cloud e l'adozione di un design *event-driven* e *serverless*. L'introduzione del parsing JSON in streaming (`ijson`) nello script `json2udm_cloud.py` è un esempio di ottimizzazione per sfruttare a pieno l'ambiente cloud, dovuta ai limiti di memoria di Cloud Run. La definizione dell'intera infrastruttura e della dashboard di monitoraggio come codice (Terraform) sottolinea le moderne pratiche DevOps per la gestione di sistemi scalabili e affidabili. Il progetto dimostra l'applicazione pratica di concetti come disaccoppiamento, resilienza (DLQ), sicurezza IAM e osservabilità end-to-end. Le competenze acquisite includono la gestione di infrastrutture complesse con IaC, il design di sistemi *serverless*, il debugging in ambienti distribuiti e l'ottimizzazione delle risorse cloud.

8 Aspetti Operativi

L'architettura progettata integra diversi meccanismi per garantire scalabilità, affidabilità, sicurezza e osservabilità, aspetti cruciali per un servizio operativo.

La **scalabilità** è gestita principalmente da Cloud Run, che adatta automaticamente il numero di istanze del Processor al carico di messaggi Pub/Sub, e dai servizi GCP sottostanti che scalano in modo trasparente. L'**affidabilità** si basa sul disaccoppiamento fornito da Pub/Sub, che garantisce la persistenza dei messaggi e gestisce i tentativi di consegna in caso di fallimenti temporanei del Processor. Cloud Run contribuisce ulteriormente monitorando la salute delle istanze tramite *health probes* e sostituendo quelle problematiche. La **sicurezza** è implementata attraverso il principio del *least privilege* per i Service Account, sull'autenticazione OIDC tra Pub/Sub e Cloud Run, la gestione delle chiavi SA on-premises e la configurazione sicura dei bucket GCS. Infine, l'**osservabilità** è un punto di forza di questa versione evoluta. Oltre ai log centralizzati in Cloud Logging, il sistema si avvale di strumenti di monitoraggio avanzati:

- **Metriche Basate su Log (LBMs):** Queste metriche, definite via Terraform, catturano eventi specifici dell'applicazione (es. `sniffer_heartbeat_count`, `pcap_files_uploaded_count`, `processor_pcap_processing_latency_seconds`, `processor_udm_packets_processed_count`). Sono essenziali per comprendere il comportamento interno e le prestazioni della pipeline.
- **Dashboard Operativa Personalizzata:** Anch'essa gestita come codice (`json`) e deployata da Terraform, utilizza MQL per visualizzare sia le LBMs sia le metriche standard dei servizi GCP. Offre una visione d'insieme dello stato della pipeline, delle prestazioni dei singoli componenti (sniffer, Pub/Sub, processor) e dei tassi di errore, cruciale per la manutenzione proattiva e il troubleshooting.

9 Test e Validazione

La validazione della pipeline è facilitata dalla Test Generator VM, un componente cruciale dell'infrastruttura definita tramite Terraform. Questa VM non è semplicemente un ambiente di test generico, ma è concepita specificamente per simulare un ambiente on-premises o edge in cui lo sniffer opererebbe tipicamente. Questo ambiente simulato è fondamentale perché permette di replicare le condizioni operative reali del componente Sniffer, che è l'unico elemento della pipeline destinato a operare al di fuori dell'infrastruttura GCP gestita. Lo script di avvio della VM (`startup_script_vm.sh`) automatizza la configurazione di un ambiente Docker Compose completo, pronto per eseguire il container dello sniffer. In questo ambiente simulato, è possibile utilizzare strumenti come `tcpreplay` per "riprodurre" file PCAP preesistenti (ad esempio, un file illustrativo come `synflood_capture.pcap` contenente traffico di un attacco simulato), immettendo traffico controllato nella rete della VM che verrà catturato dallo sniffer.

¹ *# Questo script, eseguito all'avvio della VM di test, automatizza la configurazione*
² *# dell'ambiente per lo sniffer. Recupera i parametri necessari dai metadati dell'istanza*

```

3 # (passati da Terraform) e crea i file di configurazione per Docker Compose.
4
5 # ... (installazione docker, pull immagine dello sniffer) ...
6
7 echo "Creating .env file for the sniffer in $SNIFFER_ENV_FILE"
8 # Il file .env conterra' le variabili d'ambiente specifiche per questa istanza
9 # dello sniffer sulla VM di test.
10 cat << EOF_ENV > "$SNIFFER_ENV_FILE"
11 GCP_PROJECT_ID=$VM_GCP_PROJECT_ID_FROM_METADATA
12 INCOMING_BUCKET=$VM_INCOMING_BUCKET_FROM_METADATA
13 PUBSUB_TOPIC_ID=$VM_PUBSUB_TOPIC_ID_FROM_METADATA
14 SNIFFER_ID=${VM_SNIFFER_ID_FROM_METADATA}
15 # GCP_KEY_FILE e' ${SNIFFER_GCP_KEY_CONTAINER_PATH}/key.json
16 GCP_KEY_FILE=${SNIFFER_GCP_KEY_CONTAINER_PATH}/key.json # Path interno al container
17 EOF_ENV
18
19 # ... (creazione di docker-compose.yml e docker-compose.override.yml specifici per la VM)
20 # L'override.yml gestisce il montaggio della chiave SA e della directory delle catture.

```

Listing 9: Estratto da `terraform/modules/test_generator_vm/startup_script_vm.sh`.

In aggiunta, per testare la reazione del sistema a scenari di traffico specifici o attacchi simulati, è possibile installare e utilizzare sulla VM strumenti come `hping3`. Ad esempio, un attacco SYN flood può essere simulato per osservare come i pacchetti anomali vengono catturati, processati e rappresentati nel formato UDM.

```

1 # Questo comando invia un gran numero di pacchetti TCP SYN alla porta 80
2 # dell'IP target, simulando l'inizio di molte connessioni senza completarle.
3 sudo hping3 -S --flood -p 80 <IP_TARGET_NELLA_RETE_DELLA_VM>

```

Listing 10: Esempio di comando `hping3` per simulare un SYN flood (da eseguire sulla VM di test)

Questa capacità di simulare un ambiente on-premises realistico, completo di strumenti per la generazione e la riproduzione di traffico, è fondamentale per validare end-to-end il funzionamento della pipeline e per osservarne il comportamento in condizioni controllate, inclusa la risposta a potenziali minacce di rete. File di esempio come un ipotetico `synflood_capture.pcap` e il corrispondente `synflood_capture.udm.json` servirebbero come riferimento per questa validazione.

10 Conclusioni e Sviluppi Futuri

Questo progetto ha raggiunto l'obiettivo di evolvere una semplice utility di conversione in una pipeline di analisi del traffico di rete robusta, scalabile, osservabile e gestita come codice su GCP. L'adozione di tecniche come il parsing in streaming, l'architettura *event-driven* con DLQ, e l'osservabilità tramite dashboard IaC dimostrano una comprensione matura dei principi per la creazione di servizi cloud affidabili.

Possibili sviluppi futuri includono:

- Implementazione di test automatici di integrazione e unitari.
- Integrazione diretta con l'API di ingestion di *Chronicle Security Operations* (Google SecOps).

Il sistema attuale, tuttavia, costituisce una solida piattaforma per l'analisi del traffico di rete su scala.

11 Bibliografia e Sitografia

Riferimenti bibliografici

- [1] Lucchesi, Filippo (2025). *chronicle-Sniffer*. GitHub Repository. <https://github.com/fillol/chronicle-sniffer> (Codice Sorgente)
- [2] Lucchesi, Filippo (2024). *Wireshark-to-Chronicle-Pipeline*. GitHub Repository. <https://github.com/fillol/Wireshark-to-Chronicle-Pipeline> (Progetto originario di Cybersecurity)
- [3] Immagine Sniffer su Dockerhub. <https://hub.docker.com/r/fillol/chronicle-sniffer>
- [4] Google Cloud. Cloud Run. <https://cloud.google.com/run/docs>
- [5] Google Cloud. Pub/Sub. <https://cloud.google.com/pubsub/docs>
- [6] Google Cloud. Cloud Storage. <https://cloud.google.com/storage/docs>
- [7] Google Cloud. Unified Data Model (UDM). <https://cloud.google.com/chronicle/docs/unified-data-model/udm-overview>
- [8] Google Cloud. Cloud Monitoring. <https://cloud.google.com/monitoring/docs>
- [9] Google Cloud. Identity and Access Management (IAM). <https://cloud.google.com/iam/docs>
- [10] Google SecOps. Chronicle <https://cloud.google.com/security/products/security-operations>
- [11] Terraform by HashiCorp. <https://developer.hashicorp.com/terraform/docs>
- [12] Docker. <https://docs.docker.com/>
- [13] TShark Manual Page. <https://www.wireshark.org/docs/man-pages/tshark.html>
- [14] Python Software Foundation. ijson <https://pypi.org/project/ijson/>
- [15] Hping3. <https://www.kali.org/tools/hping3/>
- [16] Tcpreplay Manual Page. <https://tcpreplay.appneta.com/wiki/tcpreplay-man.html>