

# Audit Log GCS Sink Service Operations Guide

- [1. Useful Links](#)
- [2. Connector APIs](#)
- [3. Logs and Health Metrics](#)
- [4. Troubleshooting the Issues](#)
- [5. Issues Observed in Prod](#)

## 1. Useful Links

List of all the frequently used system related documents link

| Name                              | Link  | Description  |
|-----------------------------------|---|--|
| ADT                               | <a href="#">API Logs for Product Metrics - ADT</a>      | Design Document Of Product Metrics.  |
| GCS Sink Service                  | <a href="#">Kafka Connect GCS Sink Service</a>          | GCS Sink Service Document.   |
| Integrating Audit Logging Service | <a href="#">Audit Logging Service Integration Guide</a> | Document to integrate the audit logging service.                           |
| Visualising Product Metrics       | <a href="#">Visualising Product Metrics</a>             | Data discovery access and accessing data stored in GCS via data discovery. |
| KCAAS DX overview                 | <a href="#">Kafka connect overview link</a>             | Kafka connect as a service overview link                                   |

## 2. Connector APIs

List the APIs exposed by the connector along with usage, sample payloads, and response structures.

| Curl    |   | Description             | Sample Response  |
|---------|---|-------------------------|--|
| Cluster | Curl  | Status of the connector | <pre>{   "connector": {     "state": "RUNNING",     "worker_id": "11.19.25.41:8083"   },   "name": "audit-log-gcs-sink-connector",   "type": "sink",   "tasks": [     {       "id": 0,       "state": "RUNNING",       "worker_id": "11.19.25.41:8083"     }   ] }</pre> |
| EUS     | <code>curl --location 'http://audit-log-gcs-sink.data-ventures-luminate-cperf.eus2-prod-a30.cluster.k8s.us.walmart.net/connectors/audit-log-gcs-sink-connector/status'</code> |                         |  |
| SCUS    | <code>curl --location 'http://audit-log-gcs-sink.data-ventures-luminate-cperf.scus-prod-a63.cluster.k8s.us.walmart.net/connectors/audit-log-gcs-sink-connector/status'</code> |                         |  |

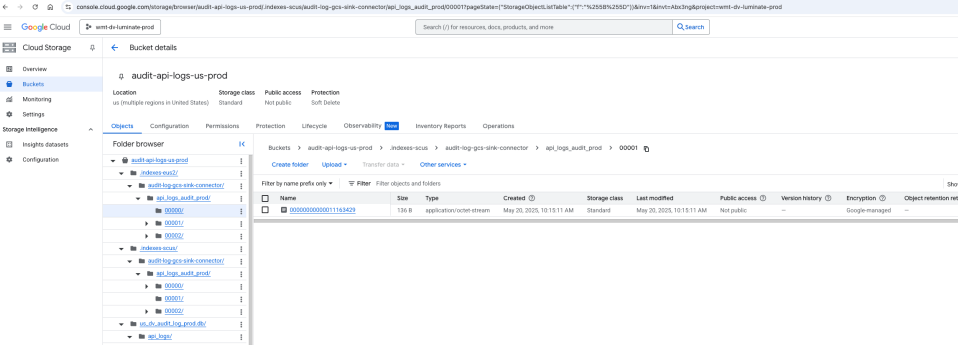
|         |  |                                 |  |
|---------|--|---------------------------------|--|
| Cluster | Curl   | Configurations of the connector | <pre>{   "connector.class": "io.lenses.streamreactor.connect.gcp.storage.sink.GCPStorageSinkConnector",   "errors.log.include.messages": "true",   "build.id": "0.0.289",   "transforms.InsertRollingRecordTimestamp.timezone": "GMT",   "tasks.max": "1",   "timezone": "GMT",   "topics": "api_logs_audit_prod",   "connect.gcpstorage.gcp.auth.mode": "FILE",   "connect.gcpstorage.gcp.file": "*****",   "transforms": "InsertRollingRecordTimestamp, FilterUS",   "locale": "en",   "errors.deadletterqueue.context.headers.enable": "true",   "transforms.InsertRollingRecordTimestamp.type": "io.lenses.connect.smt.header.InsertRollingRecordTimestampHeaders",   "connect.gcpstorage.gcp.project.id": "wmt-dv-luminate-prod",   "connect.gcpstorage.kcql": "INSERT INTO `audit-api-logs-us-prod:us_dv_audit_log_prod.db/api_logs` SELECT * FROM `api_logs_audit_prod` PARTITIONBY service_name, _header.date, endpoint_name STOREAS `PARQUET` PROPERTIES(\n 'flush.size'='50000000',\n 'flush.count'='5000',\n 'flush.interval'='600',\n 'key.suffix'='_eus2'\n );\n",   "transforms.FilterUS.type": "com.walmart.audit.log.sink.converter.AuditLogSinkUSFilter",   "errors.deadletterqueue.topic.name": "api_logs_audit_prod_DLQ",   "connect.gcpstorage.indexes.name": ".indexes-eus2",   "name": "audit-log-gcs-sink-connector",   "errors.tolerance": "all",   "transforms.InsertRollingRecordTimestamp.date.format": "yyyy-MM-dd",   "errors.log.enable": "true" }</pre> |
| EUS     | curl --location 'http://audit-log-gcs-sink.data-ventures-luminate-cperf.eus2-prod-a30.cluster.k8s.us.walmart.net/connectors/audit-log-gcs-sink-connector/config' |                                 |  |
| SCUS    | curl --location 'http://audit-log-gcs-sink.data-ventures-luminate-cperf.scus-prod-a63.cluster.k8s.us.walmart.net/connectors/audit-log-gcs-sink-connector/config' |                                 |  |
| Cluster | Curl   | Restarting the connectors       | <pre>{   "connector.class": "io.lenses.streamreactor.connect.gcp.storage.sink.GCPStorageSinkConnector",   "errors.log.include.messages": "true",   "build.id": "0.0.289",   "transforms.InsertRollingRecordTimestamp.timezone": "GMT",   "tasks.max": "1",   "timezone": "GMT",   "topics": "api_logs_audit_prod",   "connect.gcpstorage.gcp.auth.mode": "FILE",   "connect.gcpstorage.gcp.file": "*****",   "transforms": "InsertRollingRecordTimestamp, FilterUS",   "locale": "en",   "errors.deadletterqueue.context.headers.enable": "true",   "transforms.InsertRollingRecordTimestamp.type": "io.lenses.connect.smt.header.InsertRollingRecordTimestampHeaders",   "connect.gcpstorage.gcp.project.id": "wmt-dv-luminate-prod",   "connect.gcpstorage.kcql": "INSERT INTO `audit-api-logs-us-prod:us_dv_audit_log_prod.db/api_logs` SELECT * FROM `api_logs_audit_prod` PARTITIONBY service_name, _header.date, endpoint_name STOREAS `PARQUET` PROPERTIES(\n 'flush.size'='50000000',\n 'flush.count'='5000',\n 'flush.interval'='600',\n 'key.suffix'='_eus2'\n );\n",   "transforms.FilterUS.type": "com.walmart.audit.log.sink.converter.AuditLogSinkUSFilter",   "errors.deadletterqueue.topic.name": "api_logs_audit_prod_DLQ",   "connect.gcpstorage.indexes.name": ".indexes-eus2",   "name": "audit-log-gcs-sink-connector",   "errors.tolerance": "all",   "transforms.InsertRollingRecordTimestamp.date.format": "yyyy-MM-dd",   "errors.log.enable": "true" }</pre> |
| EUS     | curl --location 'http://audit-log-gcs-sink.data-ventures-luminate-cperf.eus2-prod-a30.cluster.k8s.us.walmart.net/connectors/restart_all'                         |                                 |  |
| SCUS    | curl --location 'http://audit-log-gcs-sink.data-ventures-luminate-cperf.scus-prod-a63.cluster.k8s.us.walmart.net/connectors/restart_all'                         |                                 |  |

| Cluster | Curl   | Update the logging level for the connectors.<br>Available levels - DEBUG, ERROR, INFO, WARN | <pre>[ "com.walmart.audit.log.sink.converter.AuditLogSinkUSFilter", "com.walmart.audit.log.sink.converter.CustomHeaderConverter", "io.confluent.connect.avro.AvroConverterConfig", "io.confluent.connect.avro.AvroData", "io.confluent.connect.avro.AvroDataConfig", "io.confluent.kafka.schemaregistry.rules.NoneAction", "io.confluent.kafka.serializers.AbstractKafkaSchemaSerDe", "io.confluent.kafka.serializers.KafkaAvroDeserializerConfig", "io.confluent.kafka.serializers.KafkaAvroSerializerConfig", "io.lenses.streamreactor.common.errors.ErrorPolicy\$", "io.lenses.streamreactor.common.errors.ThrowErrorPolicy", "io.lenses.streamreactor.common.util.AsciiArtPrinter", "io.lenses.streamreactor.connect.cloud.common.sink.seek.IndexManager", "io.lenses.streamreactor.connect.cloud.common.sink.transformers.TopicsTransformers\$", "io.lenses.streamreactor.connect.cloud.common.sink.writer.NoWriter", "io.lenses.streamreactor.connect.cloud.common.sink.writer.Uploading", "org.apache.kafka.common.utils.Utils", "org.apache.kafka.connect.cli.AbstractConnectCli", "org.apache.kafka.connect.connector.policy.AllConnectorClientConfigOverridePolicy", "org.apache.kafka.connect.connector.policy.NoneConnectorClientConfigOverridePolicy", "org.apache.kafka.connect.connector.policy.PrincipalConnectorClientConfigOverridePolicy", "org.apache.kafka.connect.data.Values\$ValueParser", "org.apache.kafka.connect.json.JsonConverterConfig", "org.apache.kafka.connect.mirror.MirrorCheckpointConnector", "org.apache.kafka.connect.mirror.MirrorSourceConnector", "org.apache.kafka.connect.rest.basic.auth.extension.BasicAuthSecurityRestExtension", "org.apache.kafka.connect.runtime.AbstractHerder", "org.apache.kafka.connect.runtime.WorkerConnector", "org.apache.kafka.connect.runtime.WorkerInfo", "org.apache.kafka.connect.runtime.WorkerSinkTask", "org.eclipse.jetty.util.log", "org.eclipse.jetty.util.thread.QueuedThreadPool", "org.eclipse.jetty.util.thread.ReservedThreadExecutor", "org.eclipse.jetty.util.thread.ShutdownThread", "org.eclipse.jetty.util.thread.ThreadPoolBudget", "org.eclipse.jetty.util.thread.strategy.EatWhatYouKill", "org.reflections", "org.reflections.Reflections", "root", "shaded.parquet.org.apache.thrift.transport.TIOStreamTransport" ]</pre> |
|---------|--|---|--|
| EUS     | curl --location 'http://audit-log-gcs-sink.data-ventures-luminate-cperf.eus2-prod-a30.cluster.k8s.us.walmart.net/connectors/log_all/TRACE' |   |  |
| SCUS    | curl --location 'http://audit-log-gcs-sink.data-ventures-luminate-cperf.scus-prod-a63.cluster.k8s.us.walmart.net/connectors/log_all/TRACE' |   |  |

### 3. Logs and Health Metrics

How to access logs (e.g., K8s ) and what metrics to monitor (e.g., consumer lag, memory utilization).

| Name    |  | Description  |
|---------|--|--|
| Topic   |  | <b>api_logs_audit_prod</b> is the topic in prod EUS and SCUS cluster. Separate topics are present in both cluster. |
| Cluster | Link                                       |  |
| EUS     | <a href="#">api_logs_audit_prod - EUS</a>  |  |
| SCUS    | <a href="#">api_logs_audit_prod - SCUS</a> |  |

| <p>Consumer group</p> <table border="1"> <thead> <tr> <th>Cluster</th><th>Link</th></tr> </thead> <tbody> <tr> <td>EUS</td><td><a href="#">connect-audit-log-gcs-sink-connector - EUS</a></td></tr> <tr> <td>SCUS</td><td><a href="#">connect-audit-log-gcs-sink-connector - SCUS</a></td></tr> </tbody> </table> | Cluster   | Link | EUS | <a href="#">connect-audit-log-gcs-sink-connector - EUS</a> | SCUS | <a href="#">connect-audit-log-gcs-sink-connector - SCUS</a> | <p><b>connect-audit-log-gcs-sink-connector</b> is the consumer group in prod EUS and SCUS cluster. Separate consumer groups are present in both cluster. We have alerts set up. When consumer lag is more than <b>50K we get WARNING alert</b>. When it is more than <b>75K we get critical alert</b>. Flush properties are as follows. Here once any of the threshold are met the connector will start writing the record on to GCS</p> <table border="1"> <thead> <tr> <th>Flush property</th><th>Threshold</th></tr> </thead> <tbody> <tr> <td>size</td><td>50 MB</td></tr> <tr> <td>interval</td><td>10 Mins</td></tr> <tr> <td>count</td><td>5000 records</td></tr> </tbody> </table> | Flush property | Threshold | size | 50 MB | interval | 10 Mins | count | 5000 records |
|---|---|------|-----|--|------|---|--|----------------|-----------|------|-------|----------|---------|-------|--------------|
| Cluster   | Link  |      |     |  |      |   |  |                |           |      |       |          |         |       |              |
| EUS   | <a href="#">connect-audit-log-gcs-sink-connector - EUS</a>  |      |     |  |      |   |  |                |           |      |       |          |         |       |              |
| SCUS  | <a href="#">connect-audit-log-gcs-sink-connector - SCUS</a>   |      |     |  |      |   |  |                |           |      |       |          |         |       |              |
| Flush property  | Threshold   |      |     |  |      |   |  |                |           |      |       |          |         |       |              |
| size  | 50 MB   |      |     |  |      |   |  |                |           |      |       |          |         |       |              |
| interval  | 10 Mins   |      |     |  |      |   |  |                |           |      |       |          |         |       |              |
| count   | 5000 records  |      |     |  |      |   |  |                |           |      |       |          |         |       |              |
| <p>K8s Dashboard</p> <table border="1"> <thead> <tr> <th>Cluster</th><th>Link</th></tr> </thead> <tbody> <tr> <td>EUS</td><td><a href="#">K8s Dashboard EUS - audit-log-gcs-sink</a></td></tr> <tr> <td>SCUS</td><td><a href="#">K8s Dashboard SCUS - audit-log-gcs-sink</a></td></tr> </tbody> </table>          | Cluster   | Link | EUS | <a href="#">K8s Dashboard EUS - audit-log-gcs-sink</a>     | SCUS | <a href="#">K8s Dashboard SCUS - audit-log-gcs-sink</a>     | <p>K8s dashboard in EUS and SCUS cluster. Service name is <b>audit-log-gcs-sink</b>. This can be used to monitor and download the logs. We can also monitor the pod health at that time. For a particular time range metrics Grafana is the go to dashboard.</p>   |                |           |      |       |          |         |       |              |
| Cluster   | Link  |      |     |  |      |   |  |                |           |      |       |          |         |       |              |
| EUS   | <a href="#">K8s Dashboard EUS - audit-log-gcs-sink</a>  |      |     |  |      |   |  |                |           |      |       |          |         |       |              |
| SCUS  | <a href="#">K8s Dashboard SCUS - audit-log-gcs-sink</a>   |      |     |  |      |   |  |                |           |      |       |          |         |       |              |
| <p>GCS Bucket US - <a href="#">audit-api-logs-us-prod</a></p>   | <p>Sink service stores the records in this GCS bucket. We can check here for latest updates on offsets commits based on <b>Last modified</b> time. Actual records will be stored in respective partitions</p> <p>EUS offsets maintained here by connector <b>audit-api-logs-us-prod/.indexes-eus2/audit-log-gcs-sink-connector/api_logs_audit_prod/00000</b></p> <p>SCUS offsets maintained here by connector <b>audit-api-logs-us-prod/.indexes-scus/audit-log-gcs-sink-connector/api_logs_audit_prod/00001</b></p>  |      |     |  |      |   |  |                |           |      |       |          |         |       |              |
| <p>Grafana Golden Signals - <a href="#">Grafana Golden Signals</a></p>  | <p>Application metrics can be found here i.e, health, latency, memory utilization etc.. Any deviation here then we need to take appropriate config updates in <b>kitt.yaml</b> .</p> <p>We can select the appropriate clusters in the dropdown.</p>   |      |     |  |      |   |  |                |           |      |       |          |         |       |              |
| <p>Kafka Broker Metrics - <a href="#">Managed Broker Metrics</a></p>  | <p>Metrics related to Kafka broker can be found here. This provides a very high level analysis at the broker level</p>  |      |     |  |      |   |  |                |           |      |       |          |         |       |              |

## 4. Troubleshooting the Issues

Common issues and how to debug them using logs, configurations, or metrics. There are chances the lag is due to rebalancing, connector being killed because of an error, GCS not responding etc.. Here are the standard steps to be followed in order to address any issue

| Step    | Description   | Action  |      |     |   |      |  |  |
|---------|---|---|------|-----|---|------|--|--|
| 1       | <div>Consumer groups analysis</div> <table><tr><th>Cluster</th><th>Link</th></tr><tr><td>EUS</td><td><a href="#">connect-audit-log-gcs-sink-conector - EUS</a></td></tr><tr><td>SCUS</td><td><a href="#">connect-audit-log-gcs-sink-conector - SCUS</a></td></tr></table> | Cluster   | Link | EUS | <a href="#">connect-audit-log-gcs-sink-conector - EUS</a> | SCUS | <a href="#">connect-audit-log-gcs-sink-conector - SCUS</a> | Check the logs in K8s for any exception, error, memory related errors. Analyse the consumer groups. Is it in rebalancing state ? Inactive state ? Some underlying infra changes too can cause consumer rebalancing. If the application is not able to heal from here then proceed to next steps. |
| Cluster | Link  |   |      |     |   |      |  |  |
| EUS     | <a href="#">connect-audit-log-gcs-sink-conector - EUS</a>   |   |      |     |   |      |  |  |
| SCUS    | <a href="#">connect-audit-log-gcs-sink-conector - SCUS</a>  |   |      |     |   |      |  |  |
| 2       | Check the connector status, config. Set of APIs to support this analysis<br>Refer pointer 2. APIs of the Connector  | Check the connector task status. Due to some internal errors the connector task can be killed and it may require restart as well. Even if the consumer is rebalancing for more than 10 minutes then better to restart the connector task. Enable the TRACE log level. On restarting wait for 10 minutes to see if consumer is able to rejoin the group and start consuming the message (Lag should be reduced / We can see latest offsets commit as GCS bucket). If it fails again then health related metrics to be checked. |      |     |   |      |  |  |
| 3       | Grafana golden signals analysis - <a href="#">Grafana Golden Signals</a>  | We can get to know the health metics of the connector here. Check for memory utilization, CPU utilization and take appropriate action updating configs in <b>kitt.yaml</b>  |      |     |   |      |  |  |
| 4       | Heap dump and Heap memory related analysis  | On analysing the heap dump / if in case of heap memory error try to update the <b>Xmx, Xms</b> in property <b>KAFKA_HEAP_OPTS</b> of <b>kitt.yaml</b>   |      |     |   |      |  |  |

## 5. Issues Observed in Prod

List of known or recurring issues observed in production with logs and mitigation steps.

| Issue                    | Description   | Logs  | Resolution  | Current config   |  |  |
|--------------------------|---|---|---|--|--|--|
| Consumer poll() time out | This can occur when consumer is failing to poll records in the stipulated interval. The four main properties for this to look for and update is <b>max.poll.records</b> , <b>max.poll.interval.ms</b> , <b>consumer.max.poll.records</b> , <b>consumer.max.poll.interval.ms</b> | [2025-05-14 15:02:15,417] WARN [audit-log-gcs-sink-conector task-1] [Consumer clientId=connector-consumer-audit-log-gcs-sink-conector-1, groupId=connect-audit-log-gcs-sink-conector] consumer poll timeout has expired. This means the time between subsequent calls to poll() was longer than the configured <b>max.poll.interval.ms</b> , which typically implies that the poll loop is spending too much time processing messages. You can address this either by increasing <b>max.poll.interval.ms</b> or by reducing the maximum size of batches returned in poll() with <b>max.poll.records</b> . | To increase the worker consumer's <b>max.poll.interval.ms</b> and decrease the <b>max.poll.records</b> . Keep the same values for connector's instance i.e, <b>consumer.max.poll.records</b> , <b>consumer.max.poll.interval.ms</b> | <p>Found in file -&gt; <b>kc_config.yaml</b></p> <p># For worker consumer</p> <pre>max.poll.records: 50 max.poll.interval.ms: 300000 heartbeat.interval.ms: 5000 session.timeout.ms: 15000 request.timeout.ms: 60000</pre> <p># For connector instance</p> <pre>consumer.max.poll.records: 50 consumer.max.poll.interval.ms: 300000 consumer.heartbeat.interval.ms: 5000 consumer.session.timeout.ms: 15000 consumer.request.timeout.ms: 60000</pre> |  |  |

|   |  |  |  |   |  |  |
|---|--|--|--|---|--|--|
| java.lang.<br>OutOfMemoryError<br>: Java<br>heap<br>space                 | When there is data which<br>can take more heap<br>memory / memory leaks<br>we can encounter this.<br>Due to this the pods can<br>be restarting intermittently.               | Failed to deserialize value for header 'x-edgeconnect-session-id' on topic 'api_logs_audit_prod',<br>so using byte array (org.apache.kafka.connect.storage.SimpleHeaderConverter:72) java.lang.<br>OutOfMemoryError: Java heap space   | Address if any memory leaks.<br>Increase the <b>Xmx, Xms</b> in<br>property <b>KAFKA_HEAP_OPTS</b><br>of <b>kitt.yaml</b>  | <b>KAFKA_HEAP_OPTS</b> : "-Xmx7g<br>-Xms5g -XX:<br>MetaspaceSize=96m -XX:<br>+UseG1GC -XX:<br>MaxGCPauseMillis=20 -XX:<br>InitiatingHeapOccupancyPe<br>rcent=35 -XX:<br>G1HeapRegionSize=16M -XX:<br>MinMetaspaceFreeRatio=50<br>-XX:<br>MaxMetaspaceFreeRatio=80<br>-XX:<br>+UseStringDeduplication" |  |  |
| GCS<br>storage<br>exception<br>/<br>Connector<br>killed for<br>any errors | There can be chances the<br>GCS is not responding /<br>timeout in getting<br>response from GCS. By<br>default the connector is<br>killed and it requires<br>manual restarts. | <p>[2025-05-16 10:50:14,613] ERROR [audit-log-gcs-sink-connector task-0] [audit-log-gcs-sink-connector - 1 of 1] Failed upload from data string (UploadableString(us_dv_audit_log_prod.db<br/>/api_logs/service_name=NRT/date=2025-05-16/endpoint_name=nrti_status/api_logs_audit_prod<br/>(0_000011006149)_eus2.parquet)) to Storage audit-api-logs-us-prod::indexes-eus2/audit-log-gcs-<br/>sink-connector/api_logs_audit_prod/00000/00000000000011006149 (io.lenses.streamreactor.<br/>connect.gcp.storage.storage.GCPStorageStorageInterface:165) com.google.cloud.storage.<br/>StorageException: Unexpected end of file from server at com.google.cloud.storage.<br/>StorageException.translate(StorageException.java:179) Caused by: <a href="#">java.net.SocketException</a>:<br/>Unexpected end of file from server at java.base/<a href="#">sun.net.www.http.HttpClient.parseHTTPHeader</a><br/>(Unknown Source)</p> <p>[2025-05-16 10:50:15,681] ERROR [audit-log-gcs-sink-connector task-0] WorkerSinkTask<br/>{id=audit-log-gcs-sink-connector-0} Task threw an uncaught and unrecoverable exception. Task<br/>is being killed and will not recover until manually restarted. Error: org.apache.kafka.connect.<br/>errors.ConnectException: fatal:</p> | Default behaviour of connector<br>is to not retry when there is<br>error and task will be killed.<br>Check for task status and just<br>restart the task. Currently we<br>have update the default task<br>behaviour to retry as mentioned<br>here. If no success in that time<br>then we need to restart the task<br>which gets killed. | <b>connect.gcpstorage.error.policy</b> :<br>RETRY<br><br><b>connect.gcpstorage.max.retries</b> : 5<br><br><b>connect.gcpstorage.retry.interval</b> :<br>5000<br><br><b>connect.gcpstorage.http.retry</b> .<br><b>timeout.multiplier</b> : 1.0   |  |  |