# Parquet|ORC: OutOfMemoryError when copying with parquet|orc format

Last updated by | Veena Pachauri | Mar 8, 2023 at 11:22 PM PST

202031- Parquet/ORC: OutOfMemoryError when copying
with parquet/orc format

Wednesday, November 14, 2018
9:17 PM

| SME | |
|---|---|
| **Symptoms** | Error message:<br>ErrorCode=UserErrorJavaInvocationException,'Type=Microsoft.DataTransfer.Common.Shared.HybridDeliveryException,Message=An error occurred when invoking java, message: java.lang.OutOfMemoryError:Java heap space / Direct memory |
| **Cause** | Default JVM heap size (1G) is not enough for JVM to do (de)serialization work in copying parquet/orc format data. |
| **Resolution** | 1. If parquet is used as sink, and error message contains "Java heap space" and "doubleCapacity" in call stack<br>    a. If customer is using selfhosted-IR and version < 3.20.7159.1, please upgrade to the latest version.<br>2. Please first try mitigation 1.<br>3. If not work, ask the customer several questions and try mitigate one by one:<br>    a. What's the total RAM of the host machine?  If < 8GB, please use more powerful machine.<br>    b. What's the concurrent job limit of the host node?  Please set lower number for parquet jobs.  For example: 3. (See mitigation 2 below).<br>    c. What's the schema of the parquet file?  It would be prone to OOM when:<br>        i. There're hundreds of columns inside.<br>        ii. The row group size (block) is very large .  Please reduce the row group size if the value >1GB.<br><br>Kusto Query to list the OOM jobs:<br><br>```\nlet startTime={startTime};\nlet endTime={endTime};\nJobInfo | where TIMESTAMP >= startTime and TIMESTAMP < endTime and  JobMode == "TransferJob"\n| join\n(  CustomLogEvent | where TIMESTAMP >= startTime and TIMESTAMP < endTime and TraceMessage ==\n"TransferServiceExecutorJobPayloadSplitted"\n    | extend payload = parse_json(substring(Message, 21, strlen(Message)-44))\n    | project\nJobId,jvmSize=payload.rule.maxJvmMemoryLimit,srcFormat=payload.source.format.type,sinkFormat=payload.sink.format.type\n    | join (\n    CustomLogEvent | where TIMESTAMP >= startTime and TIMESTAMP < endTime and TraceMessage == 'TransferManagerPullNotify'\n| where Message contains "OutOfMemory"\n       ) on JobId\n) on JobId\n| join (\n   CustomLogEvent | where TIMESTAMP >= startTime and TIMESTAMP < endTime and TraceMessage == "TranferServiceJobTelemetry"\n   | extend telemetry = parse_json(substring(Message, 21, strlen(Message)-44))\n   | project ActivityId,memUseRatio=telemetry['Memory.CommittedBytesInUseRatio'],\ninboundSize=telemetry['DataSizeInbound'],\n   outboundSize=telemetry['DataSizeOutbound']\n) on ActivityId\n| project TIMESTAMP, SubscriptionId, ActivityId, JobId, DataType, DestinationType,\nGatewayVersion,memUseRatio,inboundSize,outboundSize,jvmSize,srcFormat,sinkFormat,\nJVMOOM=(todouble(memUseRatio)<85 and sinkFormat in ('ParquetFormat','OrcFormat'))\n| where JVMOOM == 1\n```<br><br>• **For selfhosted IR:**<br><br>**First: check the total RAM of the host machine.  Please make sure total RAM is >= 8GM**<br><br>Add the following environment System variable in the machine that hosts the selfhosted IR:<br>_JAVA_OPTIONS "-Xms256m -Xmx16g"<br><br><br><br>**Then restart the IR.**<br>Note: this is only a sample value.  Customer can determine the min/max heap size by him/herself.<br><br>Also please do not set too many nodes on this machine if the parquet/orc work is heavy[]<br>.<br>You may change this value in ADF portal -> Connections -> Integration Runtimes -> Edit -> Nodes<br><br><br><br>• **For azure IR:**<br>Use a selfhosed IR instead. |
| **More Information** | Related incidents, reminder, best practice, how to avoid the issue, how to troubleshoot |
| **Tags** | Provide some tags that may help search |

| CSS Feedback | Please leave you feedback if you are using the article to help customers |
|---|---|

**How good have you found this content?**

😐 🙁