

TM



---

# Kafka Admin/Ops

Support around Cassandra  
and Kafka running in EC2

---



*Kafka growing*

---

# Kafka Admin, Ops, DevOps

---

- Kafka Admin
- Kafka Ops
- Kafka DevOps
- Production Systems

# Kafka Topic Creation Important for Operations



- ❖ Topics are added and modified using the topic tool

```
1  #!/usr/bin/env bash
2
3  cd ~/kafka-training
4
5  kafka/bin/kafka-topics.sh \
6      --create \
7      --zookeeper localhost:2181 \
8      --replication-factor 3 \
9      --partitions 23 \
10     --topic stock-prices \
11     --config min.insync.replicas=2
```

- ❖ **Replication factor** - replicas count amount of Kafka Brokers needed
  - ❖ use replication factor of at least 3 (or 2)
    - ❖ survive outages, head-room for upgrades and maintenance -ability to bounce servers
- ❖ **Partition count** - how much topic log will get sharded
  - ❖ determines broker count - if you have a partition count of 3, but have 5 servers, 2 not host topic log
  - ❖ consumers parallelism - active consumer count in consumer group



# Modifying Topics

- ❖ You can modify topic configuration
- ❖ You can add partitions
  - ❖ ***existing data partition don't change!***
  - ❖ Consumers semantics could break, data is not moved from existing partitions to new partitions
- ❖ You can use ***bin/kafka-topics.sh —alter*** to modify a topic
  - ❖ add partitions - ***you can't remove partitions!***
  - ❖ ***you can't change replication factor!***
  - ❖ modify config or ***delete*** it
- ❖ You can use ***bin/kafka-topics.sh —delete*** to delete a topic
  - ❖ Has to be enabled in Kafka Broker config - ***delete.topic.enable=true***

# Review of Kafka Topic Tools



## Create Topic

```
#!/usr/bin/env bash

cd ~/kafka-training

## Create a new Topic
kafka/bin/kafka-topics.sh \
  --create \
  --zookeeper localhost:2181 \
  --replication-factor 2 \
  --partitions 3 \
  --topic stock-prices \
  --config min.insync.replicas=1 \
  --config retention.ms=60000
```

## Describe Topic

```
#!/usr/bin/env bash
cd ~/kafka-training

# Describe existing topic
kafka/bin/kafka-topics.sh \
  --describe \
  --topic stock-prices \
  --zookeeper localhost:2181
```

## Delete Topic

```
#!/usr/bin/env bash
cd ~/kafka-training

# Delete Topic
kafka/bin/kafka-topics.sh \
  --delete \
  --zookeeper localhost:2181 \
  --topic stock-prices
```



# Alter Topic

```
alter-topic.sh x
1  #!/usr/bin/env bash
2
3  cd ~/kafka-training
4
5  ## Alter the topic
6  kafka/bin/kafka-topics.sh --alter \
7      --zookeeper localhost:2181 \
8      --partitions 13 \
9      --topic stock-prices \
10     --config min.insync.replicas=2 \
11     --delete-config retention.ms
12
```

- ❖ Changes *min.insync.replicas* from 1 to 2
- ❖ Changes partition count (*partitions*) from 3 to 13
- ❖ Use *—delete-config* to delete *retention.ms* configuration



# Modifying Topics with Alter

```
$ bin/delete-topic.sh
```

```
Topic stock-prices is marked for deletion.
```

```
$ bin/create-topic.sh
```

```
Created topic "stock-prices".
```

```
$ bin/describe-topic.sh
```

```
Topic:stock-prices      PartitionCount:3      ReplicationFactor:2      Configs:retention.ms=
      Topic: stock-prices      Partition: 0      Leader: 1      Replicas: 1,2      Isr: 1,2
      Topic: stock-prices      Partition: 1      Leader: 2      Replicas: 2,0      Isr: 2,0
      Topic: stock-prices      Partition: 2      Leader: 0      Replicas: 0,1      Isr: 0,1
```

```
$ bin/alter-topic.sh
```

```
Adding partitions succeeded!
```

```
$ bin/describe-topic.sh
```

```
Topic:stock-prices      PartitionCount:13      ReplicationFactor:2      Configs:min.insync.re
      Topic: stock-prices      Partition: 0      Leader: 1      Replicas: 1,2      Isr: 1,2
      Topic: stock-prices      Partition: 1      Leader: 2      Replicas: 2,0      Isr: 2,0
      Topic: stock-prices      Partition: 2      Leader: 0      Replicas: 0,1      Isr: 0,1
      ...
      Topic: stock-prices      Partition: 11      Leader: 0      Replicas: 0,1      Isr: 0,1
      Topic: stock-prices      Partition: 12      Leader: 1      Replicas: 1,0      Isr: 1,0
```



# Kafka Broker Graceful Shutdown

- ❖ Kafka Clustering detects Kafka broker shutdown or failure
  - ❖ Elects new partition leaders
  - ❖ For maintenance shutdowns Kafka supports graceful shutdown
- ❖ Graceful shutdown optimizations - ***controlled.shutdown.enable=true***
  - ❖ Topic logs data synced to disk = faster log recovery on restart by avoiding log recovery and checksum validation
  - ❖ Partitions are migrated to other Kafka brokers
  - ❖ Clean, fast leadership transfers, reduces partitions unavailability
  - ❖ Controlled shutdown fails if replicas on broker do not have in-sync replicas on another server





# Balancing Leadership

- ❖ When broker stops or **crashes** leadership moves to surviving brokers
    - ❖ crashed broker's partitions **transfers** to other replicas
    - ❖ If broker restarted becomes a **follower** for all its partitions
      - ❖ Recall only **leaders** read and write
- ```
bin/kafka-preferred-replica-election.sh \  
-zookeeper host:port
```
- ❖ kafka-preferred-replication.sh will rebalance leadership, OR
  - ❖ Kafka Broker Config: **auto.leader.rebalance.enable=true**
    - ❖ auto-balance leaders on change



# Kafka balancing across racks

- ❖ Kafka has rack awareness
  - ❖ spreads same partition replicas to different racks or AWS AZ (EC2 availability zones)
  - ❖ Survive single rack or single AZ outage
  - ❖ broker config: ***broker.rack=us-west-2a***
- ❖ During topic creation, rack constraint used to span replicas to as many racks as possible
  - ❖  $\min(\text{\#racks}, \text{replication-factor})$
- ❖ Assignment of replicas to brokers ensures leaders count per broker same, regardless rack distribution if racks have equal number of brokers
  - ❖ if rack has fewer brokers, then each broker in rack will get more replicas
  - ❖ keep broker count the same in each rack or AZ



# Checking Consumer Position

- ❖ Useful to see position of your consumers
  - ❖ Especially MirrorMaker consumers
- ❖ Tool to show consumer position
  - ❖ ***bin/kafka-consumer-groups.sh***
- ❖ Shows ***Topic*** and which ***Client*** (client id) and Consumer (consumer id) from consumer group is working with which Topic ***Partition***
  - ❖ GUID for Consumer ID based on ***client id plus GUID***
- ❖ ***Shows Lag between Consumer and Log***
- ❖ ***Shows Lag between Producer and what consumer can see (replicated vs non-replicated)***

# kafka-consumer-groups Describe



```
check-consumer-offsets.sh x
1  #!/usr/bin/env bash
2
3  cd ~/kafka-training
4
5  BOOTSTRAP_SERVERS="localhost:9092,localhost:9093"
6
7  kafka/bin/kafka-consumer-groups.sh --describe \
8      --bootstrap-server "$BOOTSTRAP_SERVERS" \
9      --group StockPriceConsumer
```

- ❖ Using `--describe`
- ❖ Specifies bootstrap server lists not ZooKeeper
- ❖ Specifies name of ConsumerGroup
- ❖ Will show lag, etc. for every consumer in group

# kafka-consumer-groups Describe Output



```
$ bin/check-consumer-offsets.sh
```

| TOPIC        | PARTITION | CURRENT-OFFSET | LOG-END-OFFSET | LAG | HOST       | CLIENT-ID |
|--------------|-----------|----------------|----------------|-----|------------|-----------|
| stock-prices | 5         | 910            | 910            | 0   | /10.0.1.11 | green-2   |
| stock-prices | 4         | 611            | 611            | 0   | /10.0.1.11 | green-1   |
| stock-prices | 2         | 949            | 949            | 0   | /10.0.1.11 | blue-2    |
| stock-prices | 6         | 39             | 39             | 0   | /10.0.1.11 | red-0     |
| stock-prices | 8         | 13             | 13             | 0   | /10.0.1.11 | red-2     |
| stock-prices | 1         | 13             | 13             | 0   | /10.0.1.11 | blue-1    |
| stock-prices | 3         | 1534           | 1534           | 0   | /10.0.1.11 | green-0   |
| stock-prices | 7         | -              | 0              | -   | /10.0.1.11 | red-1     |
| stock-prices | 0         | 611            | 611            | 0   | /10.0.1.11 | blue-0    |

- ❖ Shows **Topic** and which **Client** from the consumer group is working with which Topic **Partition** - Note also shows GUID for Consumer ID (not shown)
- ❖ **Current offset** is what is visible to Consumer (replicated to ISRs)
- ❖ **Log end** shows what the leader of has written

# kafka-consumer-groups Describe Output Lagging



```
$ bin/check-consumer-offsets.sh
```

| TOPIC        | PARTITION | CURRENT-OFFSET | LOG-END-OFFSET | LAG | HOST       | CLIENT-ID |
|--------------|-----------|----------------|----------------|-----|------------|-----------|
| stock-prices | 1         | 524            | 524            | 0   | /10.0.1.11 | blue-1    |
| stock-prices | 8         | 380            | 524            | 144 | /10.0.1.11 | red-2     |
| stock-prices | 7         | 0              | 0              | 0   | /10.0.1.11 | red-1     |
| stock-prices | 3         | 2959           | 3067           | 108 | /10.0.1.11 | green-0   |
| stock-prices | 0         | 909            | 1122           | 213 | /10.0.1.11 | blue-0    |
| stock-prices | 6         | 1464           | 1572           | 108 | /10.0.1.11 | red-0     |
| stock-prices | 5         | 1277           | 1421           | 144 | /10.0.1.11 | green-2   |
| stock-prices | 4         | 934            | 1122           | 188 | /10.0.1.11 | green-1   |
| stock-prices | 2         | 2464           | 2993           | 529 | /10.0.1.11 | blue-2    |

- ❖ Notice Partition 8, the replication is behind Current Offset is behind Log End
- ❖ Notice how partition 3 has 6x as many records as Partition 1
  - ❖ Could be an example of a hot spot!
- ❖ Notice how Partition 7 has no records so red-2 is idle!



# Managing Consumer Groups

- ❖ ConsumerGroupCommand - ***kafka-consumer-groups.sh***
  - ❖ you can also list, describe, or delete consumer groups
- ❖ Delete restriction -
  - ❖ Only works with older clients
  - ❖ No need for new client API because group is deleted automatically when last committed offset for group expires
  - ❖ If using older consumers that relied on ZooKeeper then you can use ***—delete***



# List Consumers

```
>_ list-consumers.sh x
1  #!/usr/bin/env bash
2
3  cd ~/kafka-training
4
5  BOOTSTRAP_SERVERS="localhost:9092,localhost:9093"
6
7  kafka/bin/kafka-consumer-groups.sh --list \
8      --bootstrap-server "$BOOTSTRAP_SERVERS"
9
```

- ❖ Use `—list` to get a list of consumers





# Expanding Kafka cluster

- ❖ Adding Kafka Brokers to cluster is simple
  - ❖ need unique broker id
  - ❖ new Kafka Brokers are not automatically assigned Topic partitions
  - ❖ You need to migrate partitions to it
- ❖ Migrating Topic Partitions is manually initiated
  - ❖ New Kafka Broker becomes followers of partitions
  - ❖ When it becomes ISR set member, then it gains leadership over partitions assigned to it
  - ❖ Once it becomes leader, existing replica will delete partition data if needed
- ❖ Kafka provides a partition reassignment tool



# Kafka Partition Reassignment Tool

- ❖ partition can be moved across brokers
- ❖ avoid hotspots, balance load on brokers
- ❖ you have to look at load on Kafka Broker
  - ❖ use ***kafka-consumer-groups.sh***
  - ❖ other admin tools to find hotspots (top, KPIs, etc.)
  - ❖ balance as needed

# Kafka Partition Reassignment Tool - Modes

## ❖ **GENERATE A PLAN —*generate***

- ❖ Inputs: Topics List, and Kafka Broker List
- ❖ Generates ***reassignment plan*** to move all topic partitions to new Kafka Brokers

## ❖ **EXECUTE A PLAN —*execute***

- ❖ Input: ***reassignment plan*** (***--reassignment-json-file***)
- ❖ Action: Does partition reassignment using plan

## ❖ **CHECK STATUS OF EXECUTE PLAN —*verify***

- ❖ Shows status of ***--execute***
- ❖ Outputs: Completed Successfully, Failed or In-Progress

# Generate Partition Reassignment Plan



```
reassign-partitions-generate-plan.sh x
1  #!/usr/bin/env bash
2  CONFIG=`pwd`/config
3  cd ~/kafka-training
4
5  # Generate Reassignment Plan
6  kafka/bin/kafka-reassign-partitions.sh --generate \
7      --broker-list 0,1,2,3 \
8      --topics-to-move-json-file "$CONFIG/move-topics.json" \
9      --zookeeper localhost:2181 > "$CONFIG/assignment-plan.json"
```

```
move-topics.json x
1  { "version": 1,
2    "topics": [
3      { "topic": "stocks" },
4      { "topic": "stock-prices" } ] }
```

- ❖ Added 4th Broker! Now we want it to have some partitions
- ❖ ***move-topics.json*** - list of topics to move in JSON format
- ❖ Generates assignment plan which needs to be edited

# Generated Partition Assignment Plan



```
assignment-plan.json x
partitions replicas
1 {
2   "version": 1,
3   "partitions": [
4     {
5       "topic": "stocks",
6       "partition": 7,
7       "replicas": [
8         0,
9         1,
10        2
11      ]
12    },
13    {
14      "topic": "stocks",
15      "partition": 2,
16      "replicas": [
17        3,
18        2,
19        0
20      ]
21    }
22  ]
23 }
```

- ❖ Assignment Plan
- ❖ List of Partitions
- ❖ List of Replicas
- ❖ Replicas might be moved to new Kafka Broker after plan executes
- ❖ Need to execute plan



# Execute Partition Reassignment Plan



```
reassign-partitions-execute-plan.sh x
1  #!/usr/bin/env bash
2  CONFIG=`pwd`/config
3  cd ~/kafka-training
4
5  # Execute reassignment plan
6  kafka/bin/kafka-reassign-partitions.sh --execute \
7      --reassignment-json-file "$CONFIG/assignment-plan.json" \
8      --throttle 100000 \
9      --zookeeper localhost:2181
```

- ❖ Executes reassignment plan
- ❖ Use generated plan or use modified generated plan
- ❖ Set throttle rate (optional) so it does not all happen at once
  - ❖ reduces load on Kafka Brokers

# Monitor Executing Partition Reassignment Plan



```
reassign-partitions-verify-plan.sh x
1  #!/usr/bin/env bash
2  CONFIG=`pwd`/config
3  cd ~/kafka-training
4
5  # Verify executing reassignment plan
6  kafka/bin/kafka-reassign-partitions.sh --verify \
7      --reassignment-json-file "$CONFIG/assignment-plan.json" \
8      --zookeeper localhost:2181
```

- ❖ Verify/Monitor reassignment plan
- ❖ Use generated plan or use modified generated plan that is already running
- ❖ Let's you know when the plan is done



# Decommissioning Kafka Brokers

- ❖ After we add a new broker,
  - ❖ add it to the `—broker-list`
  - ❖ Run generate plan
  - ❖ Execute plan
- ❖ To decommission Kafka Broker
  - ❖ Remove it from the `—broker-list`
  - ❖ Run generate plan, execute generate plan



# Generate Partition Reassignment Plan



```
reassign-partitions-generate-plan.sh x
1  #!/usr/bin/env bash
2  CONFIG=`pwd`/config
3  cd ~/kafka-training
4
5  # Generate Reassignment Plan
6  kafka/bin/kafka-reassign-partitions.sh --generate \
7      --broker-list 0,1,2 \
8      --topics-to-move-json-file "$CONFIG/move-topics.json" \
9      --zookeeper localhost:2181 > "$CONFIG/assignment-plan.json"
```

- ❖ Remove 4th Broker (3)! Now we want it reassign its partitions
- ❖ Generates assignment plan that moves partitions to 0,1,2



# Setting quotas

- ❖ You can configure quotas for client-id and user using ***kafka-configs.sh***
- ❖ Clients receive an unlimited quota
- ❖ You can set custom quotas for
  - ❖ (user, client-id) pair
  - ❖ user
  - ❖ client-id

# Setting quota for client-id, user Pair



quota.sh x

```
1  #!/usr/bin/env bash
2
3  cd ~/kafka-training
4
5  ## Add limit for stock_analyst user running as clientId stockConsumer
6  kafka/bin/kafka-configs.sh --alter \
7    --zookeeper localhost:2181 \
8    --add-config 'producer_byte_rate=1024,consumer_byte_rate=2048' \
9    --entity-type users \
10   --entity-name stock_analyst \
11   --entity-type clients \
12   --entity-name stockConsumer
```

- ❖ User stock\_analyst
- ❖ client id stockConsumer



# Quota Configuration

❖ Order of precedence for quota configuration is:

1. /config/users/<user>/clients/<client-id>
2. /config/users/<user>/clients/<default>
3. /config/users/<user>
4. /config/users/<default>/clients/<client-id>
5. /config/users/<default>/clients/<default>
6. /config/users/<default>
7. /config/clients/<client-id>
8. /config/clients/<default>



# Default Quota for Users

```
quota-default-user.sh x
1  #!/usr/bin/env bash
2
3  cd ~/kafka-training
4
5  ## Add limit to default user
6  kafka/bin/kafka-configs.sh --alter \
7    --zookeeper localhost:2181 \
8    --add-config 'producer_byte_rate=512,consumer_byte_rate=512' \
9    --entity-type users --entity-default
```

- ❖ Sets default quota for users



# Default Quota for Clients

```
quota-default-client.sh x
1  #!/usr/bin/env bash
2
3  cd ~/kafka-training
4
5  ## Add limit to default client
6  kafka/bin/kafka-configs.sh --alter \
7    --zookeeper localhost:2181 \
8    --add-config 'producer_byte_rate=512,consumer_byte_rate=512' \
9    --entity-type clients --entity-default
```

- ❖ Sets default quota for clients



# Describe a Quota

```
> quota-describe.sh x
1  #!/usr/bin/env bash
2
3  cd ~/kafka-training
4
5  ## Describe a quota
6  kafka/bin/kafka-configs.sh --describe \
7    --zookeeper localhost:2181 \
8    --entity-type users \
9    --entity-name stock_analyst \
10   --entity-type clients \
11   --entity-name stockConsumer
```

- ❖ You can see what quotas are set for a user

# Describe a Quota Output



```
$ bin/quota-describe.sh
```

```
Configs for user-principal 'stock_analyst', client-id 'stockConsumer'  
are producer_byte_rate=1024,consumer_byte_rate=2048
```

❖ Output from describe quota





# Multi-Datacenters Deploys

- ❖ Kafka may need to spans multiple datacenters or AWS regions
- ❖ Recommended approach deploy local Kafka cluster per datacenter
  - ❖ application and services using Kafka should be in same datacenter
  - ❖ Use mirroring between clusters in different datacenters
- ❖ Reduces latency from Kafka to application and services using Kafka avoid working over WAN
- ❖ Centralizes mirroring between data centers so it can be monitored
- ❖ If applications needs a global view of all data from all clusters
  - ❖ Use mirroring to provide clusters data from each cluster into one aggregate cluster
  - ❖ Aggregate clusters used by applications that require full data set
- ❖ Suggestion for most use cases

# If you need to cross WAN or DCs, ok



- ❖ Kafka batches and compresses records
  - ❖ Both producer and consumer can achieve high-throughput even over a high-latency connection
  - ❖ If needed increase the TCP socket buffer sizes for the producer, consumer, and broker
    - ❖ ***socket.send.buffer.bytes*** and ***socket.receive.buffer.bytes***
- ❖ Not a good idea to span DCs or regions
  - ❖ Really bad for ZooKeeper
  - ❖ More outages due to latency



# Important Client Configurations

- ❖ Producer configurations control
  - ❖ acks
  - ❖ compression
  - ❖ batch size
- ❖ Consumer Configuration
  - ❖ fetch size



# A Production Server Config

```
server-0.properties x
1  ## Increment by 1 for each broker
2  broker.id=0
3
4  # Kafka should have its own dedicated disk(s) or use SSD(s)
5  # To increase reads and writes, add more disks/log dirs JBOD.
6  log.dirs=./logs/kafka-0
7
8  ## Log config
9  default.replication.factor=3
10 num.partitions=8
11
12 ## Data must be replicated to at least two brokers
13 min.insync.replicas=2
14
15 ## Don't allow un-managed topics for production
16 auto.create.topics.enable=false
17
18 ## Run brokers spread over AZs or Racks
19 broker.rack=us-west2-a
20
21 ## Number of concurrent requests allowed
22 queued.max.requests=1000
23
24 ## Allow leaders to auto rebalance
25 auto.leader.rebalance.enable=true
```



# Java GC config

- Xmx6g
- Xms6g
- XX:MetaspaceSize=96m
- XX:+UseG1GC
- XX:MaxGCPauseMillis=20
- XX:InitiatingHeapOccupancyPercent=35
- XX:G1HeapRegionSize=16M
- XX:MinMetaspaceFreeRatio=50
- ❖ Use G1GC
- XX:MaxMetaspaceFreeRatio=80
- ❖ Heap Space should be 25% to 35% of available space for server
- ❖ Leave 50% for OS, Remember Kafka uses OS page cache
- ❖ Other tweaks for GC to limit overhead



# LinkedIn cluster

- ❖ One of LinkedIn's busiest clusters has:
  - ❖ 60 Kafka brokers
  - ❖ 50,000 partitions
  - ❖ Replication factor 2
  - ❖ Does 800k messages/sec in
  - ❖ 300 MB/sec inbound (writes/producers)
  - ❖ 1 GB/sec+ outbound (reads/consumers)
- ❖ 21 ms pause for 90% GC
- ❖ Less than 1 young GC per second



# Hardware and OS

- ❖ Dual quad-core Intel Xeon machines with 24GB of memory or higher
  - ❖ for production mission critical system
  - ❖ 24 GB total but only 25% of that for JVM (6 GB)
- ❖ Kafka Broker needs memory to buffer active readers and writers
  - ❖ to buffer for 30 seconds and memory needed is  $\text{write\_throughput} \times 30$
- ❖ Disk throughput is important
  - ❖ 8x7200 rpm SATA drives
  - ❖ Disk throughput is often performance bottleneck
  - ❖ JBOD - more disks is better

# OS



- ❖ Kafka production usually runs on Linux
- ❖ Ensure you have enough file descriptors
  - ❖ Kafka uses file descriptors for log segments and open connections
  - ❖  $(\text{number\_of\_partitions}) * (\text{partition\_size} / \text{segment\_size}) + \text{number\_of\_producer\_connections} + \text{number\_of\_consumer\_connections}$
  - ❖ Start with 100,000 or more file descriptors
- ❖ Max socket buffer size:
  - ❖ increased to enable high-performance data transfer between data centers
- ❖ Use JBOD instead of RAID, RAID ok, JBOD better
- ❖ Check flusher threads and PDF Flush but defaults should be ok
- ❖ Prefer filesystem XFS (largeio, nobarrier), EXT4 ok too (data=writeback, commit=num\_secs, nobh, delalloc)





# Monitoring

- ❖ Kafka uses Yammer Metrics
  - ❖ metrics reporting for Kafka Broke, Consumers and Producers
  - ❖ Reports stats using pluggable stats reporters
- ❖ Metrics exposed via JMX
- ❖ You can see what metrics are available with jconsole

# Kafka Broker Metrics -1 of 3



| DESCRIPTION                                         | JMX MBEAN NAME                                                                                          |
|-----------------------------------------------------|---------------------------------------------------------------------------------------------------------|
| <b>Message in rate</b>                              | kafka.server:type=BrokerTopicMetrics,name=MessagesInPerSec                                              |
| <b>Byte in rate</b>                                 | kafka.server:type=BrokerTopicMetrics,name=BytesInPerSec                                                 |
| <b>Request rate</b>                                 | kafka.network:type=RequestMetrics,name=RequestsPerSec,request={Produce FetchConsumer FetchFollower}     |
| <b>Byte out rate</b>                                | kafka.server:type=BrokerTopicMetrics,name=BytesOutPerSec                                                |
| <b>Log flush rate and time</b>                      | kafka.log:type=LogFlushStats,name=LogFlushRateAndTimeMs                                                 |
| <b>Time request waits in request queue</b>          | kafka.network:type=RequestMetrics,name=RequestQueueTimeMs,request={Produce FetchConsumer FetchFollower} |
| <b>Time request is processed at leader</b>          | kafka.network:type=RequestMetrics,name=LocalTimeMs,request={Produce FetchConsumer FetchFollower}        |
| <b>Messages count consumer lags behind producer</b> | kafka.consumer:type=consumer-fetch-manager-metrics,client-id={client-id} Attribute: records-lag-max     |

# Kafka Broker Metrics - 2 of 3



|                                          |                                                                         |                                                                                        |
|------------------------------------------|-------------------------------------------------------------------------|----------------------------------------------------------------------------------------|
| <b>Under replicated Count partitions</b> | kafka.server:type=ReplicaManager,name=UnderReplicatedPartitions         | 0                                                                                      |
| <b>Is controller active on broker?</b>   | kafka.controller:type=KafkaController,name=ActiveControllerCount        | Only 1 Kafka Broker is controller and has 1. All else should have 0.                   |
| <b>Leader election rate</b>              | kafka.controller:type=ControllerStats,name=LeaderElectionRateAndTimeMs  | >0 if failures                                                                         |
| <b>Unclean leader election rate</b>      | kafka.controller:type=ControllerStats,name=UncleanLeaderElectionsPerSec | 0                                                                                      |
| <b>Partition counts</b>                  | kafka.server:type=ReplicaManager,name=PartitionCount                    | mostly even across brokers                                                             |
| <b>Leader replica counts</b>             | kafka.server:type=ReplicaManager,name=LeaderCount                       | mostly even across brokers                                                             |
| <b>ISR shrink rate</b>                   | kafka.server:type=ReplicaManager,name=IsrShrinksPerSec                  | If a broker dies, ISR shrinks for some partitions. ISR expands when brokers come back. |
| <b>ISR expansion rate</b>                | kafka.server:type=ReplicaManager,name=IsrExpandsPerSec                  | Opposite of ISR shrink rate                                                            |

# Kafka Broker Metrics - 3 of 3



|                                               |                                                                                                            |                                                         |
|-----------------------------------------------|------------------------------------------------------------------------------------------------------------|---------------------------------------------------------|
| <b>Max follower lag</b>                       | kafka.server:type=ReplicaFetcherManager,name=MaxLag,clientId=Replica                                       | lag usually proportional to produce maximum batch size  |
| <b>Messages Lag per follower</b>              | kafka.server:type=FetcherLagMetrics,name=ConsumerLag,clientId=([-.\w]+),topic=([-.\w]+),partition=([0-9]+) | lag usually proportional to producer maximum batch size |
| <b>Requests waiting in producer purgatory</b> | kafka.server:type=DelayedOperationPurgatory,name=PurgatorySize,delayedOperation=Produce                    | >0 if ack=all is used                                   |
| <b>Requests waiting in fetch purgatory</b>    | kafka.server:type=DelayedOperationPurgatory,name=PurgatorySize,delayedOperation=Fetch                      | size depends on consumer config fetch.wait.max.ms       |
| <b>Request total time</b>                     | kafka.network:type=RequestMetrics,name=TotalTimeMs,request={Produce FetchConsumer FetchFollower}           | broken into queue, local, remote and response send time |
| <b>Leader replica counts</b>                  | kafka.server:type=ReplicaManager,name=LeaderCount                                                          | Should be even                                          |

# Common Metrics for Clients 1 of 2



| Metric                          | Description                                                                                                                                           |
|---------------------------------|-------------------------------------------------------------------------------------------------------------------------------------------------------|
| <b>connection-close-rate</b>    | Connections closed per second<br><br>JMX MBean Name<br>kafka.[producer consumer connect]:type=[producer consumer connect]-metrics,client-id=([-.\w]+) |
| <b>connection-creation-rate</b> | New connections established per second                                                                                                                |
| <b>network-io-rate</b>          | Average network operations count on all connections per second.                                                                                       |
| <b>outgoing-byte-rate</b>       | Average outgoing bytes count sent per second to all servers.                                                                                          |
| <b>request-rate</b>             | Average requests count sent per second.                                                                                                               |
| <b>request-size-avg</b>         | Average size of all requests                                                                                                                          |
| <b>request-size-max</b>         | Maximum size of any request                                                                                                                           |

# Common Metrics for Clients 2 of 2



| Metric                     | Description                                                                                                                                                                   |
|----------------------------|-------------------------------------------------------------------------------------------------------------------------------------------------------------------------------|
| <b>incoming-byte-rate</b>  | Average incoming byte count received by all sockets<br><br>JMX MBean Name<br>(kafka.[producer consumer connect]:type=[producer consumer connect]-metrics,client-id=([-.\w]+)) |
| <b>response-rate</b>       | Responses received sent per second.                                                                                                                                           |
| <b>select-rate</b>         | I/O layer checked for new I/O to perform per second count                                                                                                                     |
| <b>io-wait-time-ns-avg</b> | Average duration I/O thread spent waiting for a socket ready for reads/writes                                                                                                 |
| <b>io-wait-ratio</b>       | Fraction of time the I/O thread spent waiting.                                                                                                                                |
| <b>io-time-ns-avg</b>      | Average duration for I/O per select call in nanoseconds.                                                                                                                      |
| <b>io-ratio</b>            | Fraction of time I/O thread spent doing I/O.                                                                                                                                  |
| <b>connection-count</b>    | Current number of active connections.                                                                                                                                         |

# Per Kafka Broker Client Monitoring



| Metric                     | Description                                                                                                                                                                   |
|----------------------------|-------------------------------------------------------------------------------------------------------------------------------------------------------------------------------|
| <b>outgoing-byte-rate</b>  | Average outgoing byte count sent per second for node<br><br>JMX MBean Name: kafka.producer:type=[consumer producer connect]-node-metrics,client-id=([-.\w]+),node-id=([0-9]+) |
| <b>request-rate</b>        | Average requests count sent per second for a node.                                                                                                                            |
| <b>request-size-avg</b>    | Average size of all requests for node                                                                                                                                         |
| <b>request-size-max</b>    | Maximum size of any request sent for node                                                                                                                                     |
| <b>incoming-byte-rate</b>  | Average responses received count per second for node                                                                                                                          |
| <b>request-latency-avg</b> | Average request latency in ms for node                                                                                                                                        |
| <b>request-latency-max</b> | Maximum request latency in ms for node                                                                                                                                        |
| <b>response-rate</b>       | Responses received sent per second for node                                                                                                                                   |



# Kafka Producer Monitoring - 1 of 3



| Metric                        | Description                                                                                                                                                   |
|-------------------------------|---------------------------------------------------------------------------------------------------------------------------------------------------------------|
| <b>waiting-threads</b>        | User threads blocked count waiting for buffer memory to enqueue their records.<br><br>JMX MBean Name kafka.producer:type=producer-metrics,client-id=([-.\w]+) |
| <b>buffer-total-bytes</b>     | Maximum buffer memory size client can use                                                                                                                     |
| <b>buffer-available-bytes</b> | Total buffer memory size that is not being used                                                                                                               |
| <b>bufferpool-wait-time</b>   | Fraction of time an appender waits for space allocation                                                                                                       |
| <b>batch-size-avg</b>         | Average byte count sent per partition per-request.                                                                                                            |
| <b>batch-size-max</b>         | Max byte count sent per partition per-request.                                                                                                                |
| <b>compression-rate-avg</b>   | Average compression rate of record batches.                                                                                                                   |
| <b>record-queue-time-avg</b>  | Average time in ms record batches spent in record accumulator.                                                                                                |
| <b>record-queue-time-max</b>  | The maximum time in ms record batches spent in the record accumulator.                                                                                        |

# Kafka Producer Monitoring - 2 of 3



| Metric                         | Description                                                                                                   |
|--------------------------------|---------------------------------------------------------------------------------------------------------------|
| <b>request-latency-avg</b>     | Average request latency in ms.<br><br>JMX MBean Name kafka.producer:type=producer-metrics,client-id=([-.\w]+) |
| <b>request-latency-max</b>     | Maximum request latency in ms.                                                                                |
| <b>record-send-rate</b>        | Average record count sent per second                                                                          |
| <b>records-per-request-avg</b> | Average record count per request                                                                              |
| <b>record-retry-rate</b>       | Average per-second retried record send count                                                                  |
| <b>record-error-rate</b>       | Average per-second record send count that resulted in errors.                                                 |
| <b>record-size-max</b>         | Maximum record size.                                                                                          |
| <b>record-size-avg</b>         | Average record size.                                                                                          |
| <b>requests-in-flight</b>      | Current number of in-flight requests - waiting for a response.                                                |

# Kafka Producer Monitoring - 3 of 3



| Metric                           | Description                                                             |
|----------------------------------|-------------------------------------------------------------------------|
| <b>metadata-age</b>              | Age in seconds of current producer metadata being used                  |
| <b>record-send-rate</b>          | Average records sent count per second for topic                         |
| <b>byte-rate</b>                 | Average bytes sent count per second for topic                           |
| <b>compression-rate</b>          | Average record batches compression rate for topic                       |
| <b>record-retry-rate</b>         | Average per-second retried record send count for a topic                |
| <b>record-error-rate</b>         | Average per-second record sends that resulted in errors count for topic |
| <b>produce-throttle-time-max</b> | Maximum time in ms a request was throttled by a broker                  |
| <b>produce-throttle-time-avg</b> | Average time in ms a request was throttled by a broker                  |
| <b>requests-in-flight</b>        | Current number of in-flight requests - waiting for a response.          |

# Kafka Consumer Group Monitoring - 1 of 2

| Metric                             | Description                                                                                                                  |
|------------------------------------|------------------------------------------------------------------------------------------------------------------------------|
| <b>commit-latency-avg</b>          | Average duration for commit request<br><br><code>kafka.consumer:type=consumer-coordinator-metrics,client-id=([-.\w]+)</code> |
| <b>commit-latency-max</b>          | Max duration for a commit request                                                                                            |
| <b>commit-rate</b>                 | Commit call count per second                                                                                                 |
| <b>assigned-partitions</b>         | Partition count currently assigned to consumer                                                                               |
| <b>heartbeat-response-time-max</b> | Max duration for heartbeat request to receive response                                                                       |
| <b>heartbeat-rate</b>              | Average heartbeat count per second                                                                                           |
| <b>join-time-avg</b>               | Average duration for a group rejoin                                                                                          |
| <b>join-time-max</b>               | Max duration for a group rejoin                                                                                              |
| <b>join-rate</b>                   | Group join count per second                                                                                                  |

# Kafka Consumer Group Monitoring - 2 of 2

| Metric                            | Description                                  |
|-----------------------------------|----------------------------------------------|
| <b>sync-time-avg</b>              | Average duration for a group sync            |
| <b>sync-time-max</b>              | Max duration for a group sync                |
| <b>sync-rate</b>                  | Group sync count per second                  |
| <b>last-heartbeat-seconds-ago</b> | Second count since last controller heartbeat |

# Kafka Consumer Monitoring



| Metric                         | Description                               |
|--------------------------------|-------------------------------------------|
| <b>fetch-size-avg</b>          | Average byte size fetched per request     |
| <b>fetch-size-max</b>          | Maximum byte size fetched per request     |
| <b>bytes-consumed-rate</b>     | Average byte count consumed per second    |
| <b>records-per-request-avg</b> | Average record count in each request      |
| <b>records-consumed-rate</b>   | Average record count consumed per second  |
| <b>fetch-latency-avg</b>       | Average fetch request duration            |
| <b>fetch-latency-max</b>       | Max fetch request duration                |
| <b>fetch-rate</b>              | Fetch request count per second            |
| <b>records-lag-max</b>         | Max lag of record count for any partition |
| <b>fetch-throttle-time-avg</b> | Average throttle time in ms               |
| <b>fetch-throttle-time-max</b> | Maximum throttle time in ms               |

# Kafka Consumer Topic Fetch Monitoring



| Metric                         | Description                                                 |
|--------------------------------|-------------------------------------------------------------|
| <b>fetch-size-avg</b>          | Average byte size fetched per request for specific topic    |
| <b>fetch-size-max</b>          | Max byte size fetched per request for specific topic        |
| <b>bytes-consumed-rate</b>     | Average byte size consumed per second for specific topic    |
| <b>records-per-request-avg</b> | Average record count per request for specific topic         |
| <b>records-consumed-rate</b>   | Average record count consumed per second for specific topic |

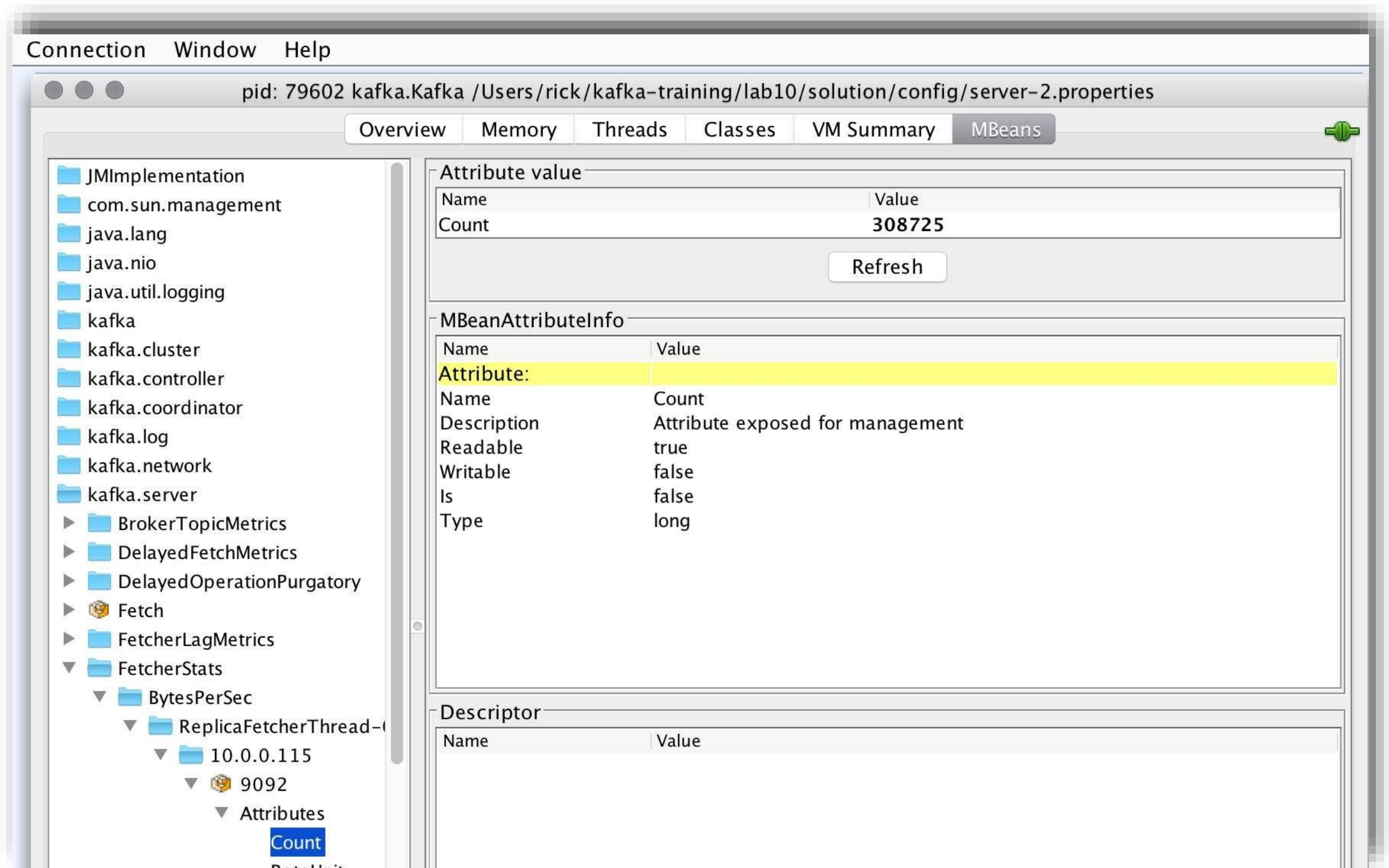


# Other Metrics

- ❖ Low level metrics
- ❖ Thread metrics
- ❖ Task Metrics
- ❖ Processor Node Metrics
  - ❖ Forwarding to other nodes
- ❖ State Store Metrics
- ❖ Good idea to monitor GC, JVM threads, etc.
- ❖ See metrics available with JConsole



# Kafka Broker Metrics via JConsole 1 of 2

Connection Window Help

pid: 79602 kafka.Kafka /Users/rick/kafka-training/lab10/solution/config/server-2.properties

Overview Memory Threads Classes VM Summary **MBeans**

**Attribute value**

| Name  | Value  |
|-------|--------|
| Count | 308725 |

Refresh

**MBeanAttributeInfo**

| Name              | Value                            |
|-------------------|----------------------------------|
| <b>Attribute:</b> |                                  |
| Name              | Count                            |
| Description       | Attribute exposed for management |
| Readable          | true                             |
| Writable          | false                            |
| Is                | false                            |
| Type              | long                             |

**Descriptor**

| Name | Value |
|------|-------|
|------|-------|

**Left Pane Tree:**

- JMImplementation
- com.sun.management
- java.lang
- java.nio
- java.util.logging
- kafka
  - kafka.cluster
  - kafka.controller
  - kafka.coordinator
  - kafka.log
  - kafka.network
  - kafka.server
    - BrokerTopicMetrics
    - DelayedFetchMetrics
    - DelayedOperationPurgatory
    - Fetch
    - FetcherLagMetrics
    - FetcherStats
      - BytesPerSec
        - ReplicaFetcherThread-0
          - 10.0.0.115
            - 9092
              - Attributes
                - Count

# Kafka Broker JConsole Metrics 2 of 2



pid: 79602 kafka.Kafka /Users/rick/kafka-training/lab10/solution/config/server-2.pro

Overview Memory Threads Classes VM Summary **MBeans**

▼ GroupMetadataManager

- ▼ NumGroups
  - ▼ Attributes
    - Value
  - ▼ Operations
    - objectName
- ▼ NumOffsets
  - ▼ Attributes
    - Value
  - Operations

▼ kafka.log

- ▼ Log
  - ▼ LogEndOffset
    - \_\_consumer\_offsets
    - ▼ stock-prices
      - 0
      - 3
      - 6
    - ▼ stocks
      - 0
      - ▼ 1
        - ▼ Attributes
          - Value

Attribute value

| Name  | Value |
|-------|-------|
| Value | 1374  |

Refresh

MBeanAttributeInfo

| Name        | Value                            |
|-------------|----------------------------------|
| Attribute:  |                                  |
| Name        | Value                            |
| Description | Attribute exposed for management |
| Readable    | true                             |
| Writable    | false                            |
| Is          | false                            |
| Type        | java.lang.Object                 |

Descriptor

| Name | Value |
|------|-------|
|------|-------|

# Kafka Producer Metrics JConsole



pid: 7852 com.intelij.rt.execution.application.AppMain com.cloudurable.kafka.producer.StockPriceP

Overview Memory Threads Classes VM Summary MBeans

com.sun.management  
java.lang  
java.nio  
java.util.logging  
kafka.producer  
  app-info  
  kafka-metrics-count  
    StockPriceProducerUtils  
      Attributes  
  producer-metrics  
    StockPriceProducerUtils  
      Attributes  
        buffer-exhausted-rate  
        **response-rate**

Attribute value

| Name          | Value              |
|---------------|--------------------|
| response-rate | 2.7192386131883075 |

Refresh

MBeanAttributeInfo

| Name        | Value                               |
|-------------|-------------------------------------|
| Attribute:  |                                     |
| Name        | response-rate                       |
| Description | Responses received sent per second. |
| Readable    | true                                |
| Writable    | false                               |
| Is          | false                               |
| Type        | double                              |

# Kafka Consumer JConsole Metrics



pid: 8131 com.intellij.rt.execution.application.AppMain com.cloudurable.kafka.consumer.ConsumerBlueMain

Overview Memory Threads Classes VM Summary **MBeans**

Tree structure:

- JMImplementation
- com.sun.management
- java.lang
- java.nio
- java.util.logging
- kafka.consumer
  - app-info
  - consumer-coordinator-metrics
    - blue
      - blue-0
        - Attributes
          - join-time-max
          - commit-latency-avg
          - sync-time-avg
          - join-rate
          - assigned-partitions
          - sync-rate
          - commit-rate**

Attribute value

| Name        | Value              |
|-------------|--------------------|
| commit-rate | 0.7556252098958917 |

Refresh

MBeanAttributeInfo

| Name        | Value                                 |
|-------------|---------------------------------------|
| Attribute:  |                                       |
| Name        | commit-rate                           |
| Description | The number of commit calls per second |
| Readable    | true                                  |
| Writable    | false                                 |
| Is          | false                                 |
| Type        | double                                |



# ZooKeeper Setup 1 of 3

- ❖ Don't put all ZooKeeper nodes in same same rack or in a single AWS availability Zones
- ❖ Decent hardware; don't use T2 Micro
- ❖ Use 5 to 7 servers for production tolerates 2 to 3 servers down
- ❖ For small deployment using 3 servers is ok (only 1 allowed down)
- ❖ Put transaction logs on dedicated disk group (***dataLogDir***)
- ❖ Put snapshots, message log, and OS on another disk/disk group (***dataDir***)
- ❖ Writes to transaction log are synchronous batches
  - ❖ Concurrent writes can significantly affect performance





# ZooKeeper Setup 2 of 3

- ❖ Use dedicated ZooKeeper cluster for Kafka
- ❖ ZooKeeper needs 3 to 5GB of heap with some room for OS (30% to 50% of System total)
- ❖ Monitoring ZooKeeper use JMX and or 4 letter words
- ❖ Keep ZooKeeper cluster small
  - ❖ Reduce quorums on the writes and subsequent cluster member updates
  - ❖ But don't go too small either
  - ❖ More ZooKeeper servers increases read capacity of ZooKeeper



# ZooKeeper Setup 3 of 3

- ❖ ZooKeeper requires little administration, but...
- ❖ ZooKeeper takes periodic snapshots of its data
  - ❖ snapshot plus log can rebuild ZooKeeper state
- ❖ ZooKeeper does not purge snapshots by default
  - ❖ Let's you back up snapshots
- ❖ You want to purge snapshots so disk does not fill up
  - ❖ `autopurge.snapRetainCount` (how many snapshots to keep)
  - ❖ `autopurge.purgeInterval` (duration in hours)
- ❖ Make sure you use rolling log files for logging