



- download
- introduction
- uses
- documentation
- quickstart
- performance
- clients
- ecosystem
- faq
- project
 - twitter
 - wiki
 - bugs
 - mailing lists
 - committees
 - powered by
 - papers & talks
- developers
 - code
 - projects
 - contributing
 - coding guide
 - protocol guide
 - unit tests

Kafka 0.9.0 Documentation

Prior releases: 0.7.2, 0.8.0, 0.8.1, 0.8.2, 0.8.3

- 1. Getting Started
 - 1.1 Introduction
 - 1.2 Use Cases
 - 1.3 Quick Start
 - 1.4 Ecosystem
 - 1.5 Upgrading
- 2. API
 - 2.1 Producer API
 - 2.2 Consumer API
 - 2.2.1 Old High Level Consumer API
 - 2.2.2 Old Simple Consumer API
 - 2.2.3 New Consumer API
- 3. Configuration
 - 3.1 Broker Configs
 - 3.2 Producer Configs
 - 3.3 Consumer Configs
 - 3.3.1 Old Consumer Configs
 - 3.3.2 New Consumer Configs
 - 3.4 Kafka Connect Configs
- 4. Design
 - 4.1 Motivation
 - 4.2 Persistence
 - 4.3 Idempotency
 - 4.4 The Producer
 - 4.5 The Consumer
 - 4.6 Message Delivery Semantics
 - 4.7 Replication
 - 4.8 Log Compaction
 - 4.9 Quotas
- 5. Implementation
 - 5.1 API Design
 - 5.2 Network Layer
 - 5.3 Messages
 - 5.4 Message Format
 - 5.5 Log
 - 5.6 Distribution
- 6. Operations
 - 6.1 Basic Kafka Operations
 - Adding and removing topics
 - Modifying topics
 - Graceful shutdown
 - Rebalancing leadership
 - Checking consumer position
 - Mirroring data between clusters
 - Sharding your cluster
 - Decommissioning brokers
 - Increasing replication factor
 - 6.2 Data centers
 - 6.2.1 Important Configs
 - Important Server Configs
 - Important Client Configs
 - Important Broker Configs
 - 6.2.4 Java Version
 - 6.2.5 Hardware and OS
 - OS
 - Disk and Filesystems
 - Application-Level File Management
 - Linux File Behavior
 - Buffer Pools
 - 6.2.6 Monitoring
 - 6.2.7 ZooKeeper
 - Stable Version
 - Operationalization
- 7. Security
 - 7.1 Security Overview
 - 7.2 Encryption and Authentication using SSL
 - 7.3 Authentication using SASL
 - 7.4 Authorization and ACLs
 - 7.5 ZooKeeper Authentication
 - New Clusters
 - Migrating Clusters
 - Migrating the ZooKeeper Ensemble
- 8. Kafka Connect
 - 8.1 Overview
 - 8.2 User Guide
 - 8.3 Connector Development Guide

1. Getting Started

1.1 Introduction

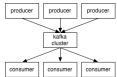
Kafka is a distributed, partitioned, replicated commit log service. It provides the functionality of a messaging system, but with a unique design.

What does all that mean?

First let's review some basic messaging terminology:

- Kafka maintains feeds of messages in categories called topics.
- We'll call processes that publish messages to a Kafka topic producers.
- We'll call processes that subscribe to topics and process the feed of published messages consumers.
- Kafka is run as a cluster comprised of one or more servers each of which is called a broker.

So, at a high level, producers send messages over the network to the Kafka cluster which in turn serves them up to consumers like this:

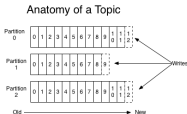


Communication between the clients and the servers is done with a simple, high-performance, language agnostic TCP protocol. We provide a Java client for Kafka, but clients are available in many languages.

Topics and Logs

Let's first dive into the high-level abstraction Kafka provides—the topic.

A topic is a category or feed name to which messages are published. For each topic, the Kafka cluster maintains a partitioned log that looks like this:



Each partition is an ordered, immutable sequence of messages that is continually appended to—a commit log. The messages in the partitions are each assigned a sequential id number called the offset that uniquely identifies each message within the partition.

The Kafka cluster retains all published messages—whether or not they have been consumed—for a configurable period of time. For example if the log retention is set to two days, then for the two days after a message is published it is available for consumption, after which it will be discarded to free up space. Kafka's performance is effectively constant with respect to data size so retaining lots of data is not a problem.

In fact the only metadata retained on a per-consumer basis is the position of the consumer in the log, called the "offset". This offset is controlled by the consumer: normally a consumer will advance its offset linearly as it reads messages, but in fact the position is controlled by the consumer and it can consume messages in any order it likes. For example a consumer can reset to an older offset to reprocess.

This combination of features means that Kafka consumers are very cheap—they can come and go without much impact on the cluster or on other consumers. For example, you can use our command line tools to "tail" the contents of any topic without changing what is consumed by any existing consumers.

The partitions in the log serve several purposes. First, they allow the log to scale beyond a size that will fit on a single server. Each individual partition must fit on the servers that host it, but a topic may have many partitions so it can handle an arbitrary amount of data. Second they act as the unit of parallelism—more on that in a bit.

Distribution

The partitions of the log are distributed over the servers in the Kafka cluster with each server handling data and requests for a share of the partitions. Each partition is replicated across a configurable number of servers for fault tolerance.

Each partition has one server which acts as the "leader" and zero or more servers which act as "followers". The leader handles all read and write requests for the partition while the followers passively replicate the leader. If the leader fails, one of the followers will automatically become the new leader. Each server acts as a leader for some of its partitions and a follower for others so load is well balanced within the cluster.

Producers

Producers publish data to the topics of their choice. The producer is responsible for choosing which message to assign to which partition within the topic. This can be done in a round-robin fashion simply to balance load or it can be done according to some semantic partition function (say based on some key in the message). More on the use of partitioning in a second.

Consumers

Messaging traditionally has two models: queueing and publish-subscribe. In a queue, a pool of consumers may read from a server and each message goes to one of them; in publish-subscribe the message is broadcast to all consumers. Kafka offers a single consumer abstraction that generalizes both of these—the consumer group.

Consumers label themselves with a consumer group name, and each message published to a topic is delivered to one consumer instance within each subscribing consumer group. Consumer instances can be in separate processes or on separate machines.

If all the consumer instances have the same consumer group, then this works just like a traditional queue balancing load over the consumers.

If all the consumer instances have different consumer groups, then this works like publish-subscribe and all messages are broadcast to all consumers.

More commonly, however, we have found that topics have a small number of consumer groups, one for each "logical subscriber". Each group is composed of many consumer instances for scalability and fault tolerance. This is nothing more than publish-subscribe semantics where the subscriber is cluster of consumers instead of a single process.

Kafka has stronger ordering guarantees than a traditional messaging system, too.

A traditional queue retains messages in order on the server, and if multiple consumers consume from the queue then the server hands out messages in the order they are stored. However, although the server hands out messages in order, the messages are delivered asynchronously to consumers, so they may arrive out of order on different consumers. This effectively means the ordering of the messages is lost in the presence of parallel consumption. Messaging systems often work around this by having a notion of "exclusive consumer" that allows only one process to consume from a queue, but of course this means that there is no parallelism in processing.

Kafka does it better. By having a notion of parallelism—the partition—within the topics, Kafka is able to provide both ordering guarantees and load balancing over a pool of consumer processes. This is achieved by assigning the partitions in the topic to the consumers in the consumer group so that each partition is consumed by exactly one consumer in the group. By doing this we ensure that the consumer is the only reader of that partition and consumes the data in order. Since there are many partitions this still balances the load over many consumer instances. Note however that there cannot be more consumer instances in a consumer group than partitions.

Kafka only provides a total order over messages within a partition, not between different partitions in a topic. Per-partition ordering combined with the ability to partition data by key is sufficient for most applications. However, if you require a total order over messages this can be achieved with a topic that has only one partition, though this will mean only one consumer process per consumer group.

At a high-level Kafka gives the following guarantee:

- Messages sent by a producer to a particular topic partition will be appended in the order they are sent. That is, if a message M1 is sent by the same producer as a message M2, and M1 is sent first, then M1 will have a lower offset than M2 and appear earlier in the log.
- A consumer instance sees messages in the order they are stored in the log.
- For a topic with replication factor N, we will tolerate up to N-1 server failures without losing any messages committed to the log.

More details on these guarantees are given in the design section of the documentation.

User Cases

Here is a description of a few of the popular use cases for Apache Kafka. For an overview of a number of these areas in action, see this blog post.

Messaging

Kafka works well as a replacement for a more traditional message broker. Message brokers are used for a variety of reasons (to decouple processing from data producers, to buffer unprocessed messages, etc). In comparison to most messaging systems Kafka has better throughput, built in partitioning, replication, and fault tolerance which makes it a good solution for large scale message processing applications.

In our experience messaging uses are often comparatively low throughput, but may require low end-to-end latency and often depend on the strong durability guarantees Kafka provides.

In this domain Kafka is comparable to traditional messaging systems such as [ActiveMQ](#) or [RabbitMQ](#).

[Website Activity Tracking](#)

The original use case for Kafka was to be able to rebuild a user activity tracking pipeline as a set of real-time publish-subscribe feeds. This means site activity (page views, searches, or other actions users may take) is published to central topics with one topic per activity type. These feeds are available for subscription for a range of use cases including real-time processing, real-time monitoring, and loading into Hadoop or offline data warehousing systems for offline processing and reporting.

Activity tracking is often very high volume as many activity messages are generated for each user page view.

Metrics

Kafka is often used for operational monitoring data. This involves aggregating statistics from distributed applications to produce centralized feeds of operational data.

[Log Aggregation](#)

Many people use Kafka as a replacement for a log aggregation solution. Log aggregation typically collects physical log files off servers and puts them in a central place (a file server or HDFS perhaps) for processing. Kafka abstracts away the details of files and gives a cleaner abstraction of log or event data as a stream of messages. This allows for lower-latency processing and easier support for multiple data sources and distributed data consumption. In comparison to log-centric systems like Scribe or Flume, Kafka offers equally good performance, stronger durability guarantees due to replication, and much lower end-to-end latency.

[Stream Processing](#)

Many users end up doing stage-wise processing of data where data is consumed from topics of raw data and then aggregated, enriched, or otherwise transformed into new Kafka topics for further consumption. For example a processing flow for article recommendation might crawl article content from RSS feeds and publish it to an "articles" topic; further processing might help normalize or deduplicate this content to a topic of cleaned article content; a final stage might attempt to match this content to users. This creates a graph of real-time data flow out of the individual topics. [Storm](#) and [Samza](#) are popular frameworks for implementing these kinds of transformations.

[Event Sourcing](#)

[Event sourcing](#) is a style of application design where state changes are logged as a time-ordered sequence of records. Kafka's support for very large stored log data makes it an excellent backend for an application built in this style.

[Commit Log](#)

Kafka can serve as a kind of external commit log for a distributed system. The log helps replicate data between nodes and acts as a re-syncing mechanism for failed nodes to restore their data. The [logcompaction](#) feature in Kafka helps support this usage. In this usage Kafka is similar to [Ripet](#) or [Bookkeeper](#) project.

[1.3 Quick Start](#)

This tutorial assumes you are starting fresh and have no existing Kafka or ZooKeeper data.

[Step 1: Download the code](#)

Download the 0.9.0.0 release and un-tar it.

```
> tar -xzf kafka_2.11-0.9.0.0.tgz
> cd kafka_2.11-0.9.0.0
```

[Step 2: Start the server](#)

Kafka uses ZooKeeper so you need to first start a ZooKeeper server if you don't already have one. You can use the convenience script packaged with kafka to get a quick and dirty single node ZooKeeper instance.

```
> bin/zookeeper-server-start.sh config/zookeeper.properties
[2013-04-22 15:01:37,495] INFO Reading configuration from: config/zookeeper.properties (org.apache.zookeeper.server.quorum.QuorumPeerConfig)
...
```

Now start the Kafka server:

```
> bin/kafka-server-start.sh config/server.properties
[2013-04-22 15:01:47,028] INFO Verifying properties (kafka.utils.VerifiableProperties)
[2013-04-22 15:01:47,031] INFO Property socket.send.buffer.bytes is overridden to 1048576 (kafka.utils.VerifiableProperties)
...
```

[Step 3: Create a topic](#)

Let's create a topic named "test" with a single partition and only one replica:

```
> bin/kafka-topics.sh --create --zookeeper localhost:2181 --replication-factor 1 --partitions 1 --topic test
```

We can now see that topic if we run the list topic command:

```
> bin/kafka-topics.sh --list --zookeeper localhost:2181
test
```

Alternatively, instead of manually creating topics you can also configure your brokers to auto-create topics when a non-existent topic is published to.

[Step 4: Send some messages](#)

Kafka comes with a command line client that will take input from a file or from standard input and send it out as messages to the Kafka cluster. By default each line will be sent as a separate message.

Run the producer and then type a few messages into the console to send to the server.

```
> bin/kafka-console-producer.sh --broker-list localhost:9092 --topic test
This is a message
This is another message
```

[Step 5: Start a consumer](#)

Kafka also has a command line consumer that will dump out messages to standard output.

```
> bin/kafka-console-consumer.sh --zookeeper localhost:2181 --topic test --from-beginning
This is a message
This is another message
```

If you have each of the above commands running in a different terminal then you should now be able to type messages into the producer terminal and see them appear in the consumer terminal.

All of the command line tools have additional options; running the command with no arguments will display usage information documenting them in more detail.

[Step 6: Setting up a multi-broker cluster](#)

So far we have been running against a single broker, but that's no fun. For Kafka, a single broker is just a cluster of size one, so nothing much changes other than starting a few more broker instances. But just to get feel for it, let's expand our cluster to three nodes (still all on our local machine).

First we make a config file for each of the brokers:

```
> cp config/server.properties config/server-1.properties
> cp config/server.properties config/server-2.properties
```

Now edit these new files and set the following properties:

```
config/server-1.properties:
broker.id=1
port=9093
log.dir=tmp/kafka-logs-1

config/server-2.properties:
broker.id=2
port=9094
log.dir=tmp/kafka-logs-2
```

The `broker.id` property is the unique and permanent name of each node in the cluster. We have to override the port and log directory only because we are running these all on the same machine and we want to keep the brokers from all trying to register on the same port or overwrite each others data.

We already have Zookeeper and our single node started, so we just need to start the two new nodes:

```
> bin/kafka-server-start.sh config/server-1.properties &
...
> bin/kafka-server-start.sh config/server-2.properties &
...
```

Now create a new topic with a replication factor of three:

```
> bin/kafka-topics.sh --create --zookeeper localhost:2181 --replication-factor 3 --partitions 1 --topic my-replicated-topic
```

Okay but now that we have a cluster how can we know which broker is doing what? To see that run the "describe topics" command:

```
> bin/kafka-topics.sh --describe --zookeeper localhost:2181 --topic my-replicated-topic
Topic: my-replicated-topic PartitionCount: 1 ReplicationFactor: 3 Configs:
Topic: my-replicated-topic Partition: 0 Leader: 1 Replicas: 1,2,0 Isr: 1,2,0
```

Here is an explanation of output. The first line gives a summary of all the partitions, each additional line gives information about one partition. Since we have only one partition for this topic there is only one line.

- "leader" is the node responsible for all reads and writes for the given partition. Each node will be the leader for a randomly selected portion of the partitions.
- "replicas" is the list of nodes that replicate the log for this partition regardless of whether they are the leader or even if they are currently alive.
- "isr" is the set of "in-sync" replicas. This is the subset of the replicas list that is currently alive and

```
caught up to the leader.

Note that in my example node 1 is the leader for the only partition of the topic.

We can run the same command on the original topic we created to see where it is:

> bin/kafka-topics.sh --describe --zookeeper localhost:2181 --topic test
Topic: test    PartitionCount: 1    ReplicationFactor: 1    Configs:
Topic: test    Partition: 0    Leader: 0    Replicas: 0    Isr: 0

So there is no surprise there—the original topic has no replicas and is on server 0, the only server in our cluster
when we created it.

Let's publish a few messages to our new topic:

> bin/kafka-console-producer.sh --broker-list localhost:9092 --topic my-replicated-topic
...
my test message 1
my test message 2
^C

Now let's consume these messages:

> bin/kafka-console-consumer.sh --zookeeper localhost:2181 --from-beginning --topic my-replicated-topic
...
My test message 1
My test message 2
^C

Now let's test out fault-tolerance. Broker 1 was acting as the leader so let's kill it:

> ps | grep server-1.properties
7564 tty002    0:15.91 /System/Library/Frameworks/JavaVM.framework/Versions/1.6/Home/bin/java...
> kill -9 7564

Leadership has switched to one of the slaves and node 1 is no longer in the in-sync replica set:

> bin/kafka-topics.sh --describe --zookeeper localhost:2181 --topic my-replicated-topic
Topic: my-replicated-topic    PartitionCount: 1    ReplicationFactor: 3    Configs:
Topic: my-replicated-topic    Partition: 0    Leader: 2    Replicas: 1,2,0    Isr: 2,0

But the messages are still be available for consumption even though the leader that took the writes originally
is down:

> bin/kafka-console-consumer.sh --zookeeper localhost:2181 --from-beginning --topic my-replicated-topic
...
My test message 1
My test message 2
^C
```

[Step 7: Use Kafka Connect to import/export data](#)

Writing data from the console and writing it back to the console is a convenient place to start, but you'll probably want to use data from other sources or export data from Kafka to other systems. For many systems, instead of writing custom integration code you can use Kafka Connect to import or export data. Kafka Connect is a tool included with Kafka that imports and exports data to Kafka. It is an extensible tool that runs connectors, which implement the custom logic for interacting with an external system. In this quickstart we'll see how to run Kafka Connect with simple connectors that import data from a file to a Kafka topic and export data from a Kafka topic to a file. First, we'll start by creating some seed data to test with:

```
> echo -e "foo\nbar" > test.txt

Next, we'll start two connectors running in standalone mode, which means they run in a single, local,
dedicated process. We provide three configuration files as parameters. The first is always the configuration for
the Kafka Connect process, containing common configuration such as the Kafka brokers to connect to and the
serialization format for data. The remaining configuration files each specify a connector to create. These files
include a unique connector name, the connector class to instantiate, and any other configuration required by
the connector.

> bin/connect-standalone.sh config/connect-standalone.properties config/connect-file-source.properties config/connect-file-sink.properties

These sample configuration files, included with Kafka, use the default local cluster configuration you started
earlier and create two connectors: the first is a source connector that reads lines from an input file and
produces each to a Kafka topic and the second is a sink connector that reads messages from a Kafka topic and
produces each as a line in an output file. During startup you'll see a number of log messages, including some
indicating that the connectors are being instantiated. Once the Kafka Connect process has started, the source
connector should start reading lines from

test.txt

and producing them to the topic:

connect-test

, and the sink connector should start reading messages from the topic:

connect-test

and write them to the file

test.sink.txt

. We can verify the data has been delivered through the entire pipeline by examining the contents of the output
file:

> cat test.sink.txt
foo
bar

Note that the data is being stored in the Kafka topic:

connect-test

, so we can also run a console consumer to see the data in the topic (or use custom consumer code to process
it):

> bin/kafka-console-consumer.sh --zookeeper localhost:2181 --topic connect-test --from-beginning
{"schema":{"type":"string","optional":false},"payload":"foo"}
{"schema":{"type":"string","optional":false},"payload":"bar"}
...

The connectors continue to process data, so we can add data to the file and see it move through the pipeline:

> echo "Another line" >> test.txt

You should see the line appear in the console consumer output and in the sink file.
```

[11.4 Ecosystem](#)

There are a plethora of tools that integrate with Kafka outside the main distribution. The ecosystem page lists many of these, including stream processing systems, Hadoop integration, monitoring, and deployment tools.

[11.5 Upgrading from Previous Versions](#)

[Upgrading from 0.8.0, 0.8.11, or 0.8.2 to 0.9.0.0](#)

0.9.0.0 has [potential breaking changes](#) (please review before upgrading) and an inter-broker protocol change from previous versions. This means that upgraded brokers and clients may not be compatible with older versions. It is important that you upgrade your Kafka cluster before upgrading your clients. If you are using MirrorMaker downstream clusters should be upgraded first as well.

Reconciling upgrades:

1. Update server.properties file on all brokers and add the following property:
inter.broker.protocol.version=0.8.2.X
2. Upgrade the brokers. This can be done a broker at a time by simply bringing it down, updating the code, and restarting it.
3. Once the entire cluster is upgraded, bump the protocol version by editing inter.broker.protocol.version and setting it to 0.9.0.0.
4. Restart the brokers one by one for the new protocol version to take effect

Note: If you are willing to accept downtime, you can simply take all the brokers down, update the code and start all of them. They will start with the new protocol by default.

Note: Bumping the protocol version and restarting can be done any time after the brokers were upgraded. It does not have to be immediately after.

[Potential breaking changes in 0.9.0.0](#)

- Java 1.6 is no longer supported.
- Scala 2.9 is no longer supported.
- Broker IDs above 1000 are now reserved by default to automatically assigned broker IDs. If your cluster has existing broker IDs above that threshold make sure to increase the reserved.broker.max.id broker configuration property accordingly.
- Configuration parameter replica.lag.max.messages was removed. Partition leaders will no longer consider the number of lagging messages when deciding which replicas are in sync.
- Configuration parameter replica.lag.time.max.ms now refers not just to the time passed since last fetch request from replica, but also to time since the replica last caught up. Replicas that are still fetching

- messages from leaders but did not catch up to the latest messages in replica.lag.time.max.ms will be considered out of sync.
- **Compacted topics no longer accept messages without key** and an exception is thrown by the producer if this is attempted. In 0.8.x, a message without key would cause the log compaction thread to subsequently complain and quit (and stop compacting all compacted topics).
 - **MirrorMaker no longer supports multiple target clusters.** As a result it will only accept a single `--consumer.config` parameter. To mirror multiple source clusters, you will need at least one MirrorMaker instance per source cluster, each with its own consumer configuration.
 - **Tools packaged under `org.apache.kafka.clients.tools.`** have been moved to `org.apache.kafka.tools.` All included scripts will still function as usual, only custom code directly importing these classes will be affected.
 - The default Kafka JVM performance options (`KAFKA_JVM_PERFORMANCE_OPTS`) have been changed in `kafka-run-class.sh`.
 - The `kafka-topics.sh` script (`kafka.admin.TopicCommand`) now exits with non-zero exit code on failure.
 - The `kafka-topics.sh` script (`kafka.admin.TopicCommand`) will now print a warning when topic names risk metric collisions due to the use of a " or " _ in the topic name, and error in the case of an actual collision.
 - The `kafka-console-producer.sh` script (`kafka.tools.ConsoleProducer`) will use the new producer instead of the old producer be default, and users have to specify 'old producer' to use the old producer.
 - By default all command line tools will print all logging messages to `stderr` instead of `stdout`.

Useful changes in 0.9.0.0

- The new broker id generation feature can be disable by setting `broker.id.generation.enable` to false.
- Configuration parameter `log.cleaner.enable` is now true by default. This means topics with a `cleanup.policy=compact` will now be compacted by default, and 128 MB of heap will be allocated to the cleaner process via `log.cleaner.dedupe.buffer.size`. You may want to review `log.cleaner.dedupe.buffer.size` and the other `log.cleaner` configuration values based on your usage of compacted topics.
- Default value of configuration parameter `fetch.min.bytes` for the new consumer is now 1 by default.

Deprecations in 0.9.0.0

- Altering topic configuration from the `kafka-topics.sh` script (`kafka.admin.TopicCommand`) has been deprecated. Going forward, please use the `kafka-configs.sh` script (`kafka.admin.ConfigCommand`) for this functionality.
- The `kafka-consumer-offset-checker.sh` (`kafka.tools.ConsumerOffsetChecker`) has been deprecated. Going forward, please use `kafka-consumer-groups.sh` (`kafka.admin.ConsumerGroupCommand`) for this functionality.
- The `kafka.tools.ProducerPerformance` class has been deprecated. Going forward, please use `org.apache.kafka.tools.ProducerPerformance` for this functionality (`kafka-producer-perf-test.sh` will also be changed to use the new class).

Upgrading from 0.8.1 to 0.8.2

0.8.2 is fully compatible with 0.8.1. The upgrade can be done one broker at a time by simply bringing it down, updating the code, and restarting it.

Upgrading from 0.8.0 to 0.8.1

0.8.1 is fully compatible with 0.8. The upgrade can be done one broker at a time by simply bringing it down, updating the code, and restarting it.

Upgrading from 0.7

Release 0.7 is incompatible with newer releases. Major changes were made to the API, Zookeeper data structures, and protocol, and configuration in order to add replication (Which was missing in 0.7). The upgrade from 0.7 to later versions requires a [special tutorial](#) for migration. This migration can be done without downtime.

2. API

Apache Kafka includes new java clients (in the `org.apache.kafka.clients` package). These are meant to supplant the older Scala clients, but for compatability they will co-exist for some time. These clients are available in a separate jar with minimal dependencies, while the old Scala clients remain packaged with the server.

2.1 Producer API

We encourage all new development to use the new Java producer. This client is production tested and generally both faster and more fully featured than the previous Scala client. You can use this client by adding a dependency on the client jar using the following example maven co-ordinates (you can change the version numbers with new releases):

```
<dependency>
  <groupId>org.apache.kafka</groupId>
  <artifactId>kafka-clients</artifactId>
  <version>0.9.0.0</version>
</dependency>
```

Examples showing how to use the producer are given in the [javadocs](#)

For those interested in the legacy Scala producer api, information can be found [here](#)

2.2 Consumer API

As of the 0.9.0 release we have added a new Java consumer to replace our existing high-level Zookeeper-based consumer and low-level consumer APIs. This client is considered beta quality. To ensure a smooth upgrade path for users, we still maintain the old 0.8 consumer clients that continue to work on an 0.9 Kafka cluster. In the following sections we introduce both the old 0.8 consumer APIs (both high-level `ConsumerConnector` and low level `SimpleConsumer`) and the new Java consumer API respectively.

2.2.1 HighLevelConsumer API

```
class Consumer {
  /**
   * Create a ConsumerConnector
   *
   * @param config at the minimum, need to specify the groupid of the consumer and the zookeeper
   *             connection string zookeeper.connect.
   */
  public static kafka.javaapi.consumer.ConsumerConnector createJavaConsumerConnector(ConsumerConfig config);
}

/**
 * V: Type of the message
 * K: Type of the optional key associated with the message
 */
public interface kafka.javaapi.consumer.ConsumerConnector {
  /**
   * Create a list of message streams of type T for each topic.
   *
   * @param topicCountMap a map of (topic, #streams) pair
   * @param decoder a decoder that converts from Message to T
   * @return a map of (topic, list of KafkaStream) pairs.
   *         The number of items in the list is #streams. Each stream supports
   *         an iterator over message/metadata pairs.
   */
  public <K,V> Map<String, List<KafkaStream<K,V>>> createMessageStreams(Map<String, Integer> topicCountMap, Decoder<K> keyDecoder, Decoder<V> valueDecoder);

  /**
   * Create a list of message streams of type T for each topic, using the default decoder.
   */
  public Map<String, List<KafkaStream<byte[], byte[]>>> createMessageStreams(Map<String, Integer> topicCountMap);

  /**
   * Create a list of message streams for topics matching a wildcard.
   *
   * @param topicFilter a TopicFilter that specifies which topics to
   *                   subscribe to (encapsulates a whitelist or a blacklist).
   * @param numStreams the number of message streams to return.
   * @param keyDecoder a decoder that decodes the message key
   * @param valueDecoder a decoder that decodes the message itself
   * @return a list of KafkaStream. Each stream supports an
   *         iterator over its MessageAndMetadata elements.
   */
  public <K,V> List<KafkaStream<K,V>> createMessageStreamsByFilter(TopicFilter topicFilter, int numStreams, Decoder<K> keyDecoder, Decoder<V> valueDecoder);

  /**
   * Create a list of message streams for topics matching a wildcard, using the default decoder.
   */
  public List<KafkaStream<byte[], byte[]>> createMessageStreamsByFilter(TopicFilter topicFilter, int numStreams);

  /**
   * Create a list of message streams for topics matching a wildcard, using the default decoder, with one stream.
   */
  public List<KafkaStream<byte[], byte[]>> createMessageStreamsByFilter(TopicFilter topicFilter);

  /**
   * Commit the offsets of all topic/partitions connected by this connector.
   */
  public void commitOffsets();

  /**
   * Shut down the connector
   */
  public void shutdown();
}
```

You can follow [this example](#) to learn how to use the high level consumer api.

2.2.2.10 SimpleConsumerAPI

```
class kafka.javaapi.consumer.SimpleConsumer {
    /**
     * Fetch a set of messages from a topic.
     *
     * @param request specifies the topic name, topic partition, starting byte offset, maximum bytes to be fetched.
     * @return a set of fetched messages
     */
    public FetchResponse fetch(kafka.javaapi.FetchRequest request);

    /**
     * Fetch metadata for a sequence of topics.
     *
     * @param request specifies the versionId, clientId, sequence of topics.
     * @return metadata for each topic in the request.
     */
    public kafka.javaapi.TopicMetadataResponse send(kafka.javaapi.TopicMetadataRequest request);

    /**
     * Get a list of valid offsets (up to maxSize) before the given time.
     *
     * @param request a [[kafka.javaapi.OffsetRequest]] object.
     * @return a [[kafka.javaapi.OffsetResponse]] object.
     */
    public kafka.javaapi.OffsetResponse getOffsetBefore(OffsetRequest request);

    /**
     * Close the SimpleConsumer.
     */
    public void close();
}
```

For most applications, the high level consumer Api is good enough. Some applications want features not exposed to the high level consumer yet (e.g., set initial offset when restarting the consumer). They can instead use our low level SimpleConsumer Api. The logic will be a bit more complicated and you can follow the [example in there](#)

2.2.2.11 NewConsumerAPI

This new unified consumer API removes the distinction between the 0.8 high-level and low-level consumer APIs. You can use this client by adding a dependency on the client jar using the following example maven co-ordinates (you can change the version numbers with new releases):

```
<dependency>
  <groupId>org.apache.kafka</groupId>
  <artifactId>kafka-clients</artifactId>
  <version>3.0.0</version>
</dependency>
```

Examples showing how to use the consumer are given in the [javautils](#).

3. Configuration

Kafka uses key-value pairs in the `property/file` format for configuration. These values can be supplied either from a file or programmatically.

3.1 BrokerConfig

The essential configurations are the following:

- broker.id
- log.dirs
- zookeeper.connect

Topic-level configurations and defaults are discussed in more detail [below](#).

Name	Description	Type	Default	Valid Values	Importance
zookeeper.connect	Zookeeper host string	string			high
advertised.host.name	Hostname to publish to ZooKeeper for clients to use. In IaaS environments, this may need to be different from the interface to which the broker binds. If this is not set, it will use the value for "host.name" if configured. Otherwise it will use the value returned from java.net.InetAddress.getCanonicalHostName().	string	null		high
advertised.listeners	Listeners to publish to ZooKeeper for clients to use. If different than the listeners above. In IaaS environments, this may need to be different from the interface to which the broker binds. If this is not set, the value for "listeners" will be used.	string	null		high
advertised.port	The port to publish to ZooKeeper for clients to use. In IaaS environments, this may need to be different from the port to which the broker binds. If this is not set, it will publish the same port that the broker binds to.	int	null		high
auto.create.topics.enable	Enable auto creation of topic on the server	boolean	true		high
auto.leader.rebalance.enable	Enables auto leader balancing. A background thread checks and triggers leader balance if required at regular intervals	boolean	true		high
background.threads	The number of threads to use for various background processing tasks	int	10	[1,...]	high
broker.id	The broker id for this server. To avoid conflicts between zookeeper generated brokerId and user's config.brokerId added MaxReservedBrokerId and zookeeper sequence starts from MaxReservedBrokerId + 1.	int	-1		high
compression.type	Specify the final compression type for a given topic. This configuration accepts the standard compression codecs ('gzip', 'snappy', 'lz4'). It additionally accepts 'uncompressed' which is equivalent to no compression; and 'producer' which means retain the original compression codec set by the producer.	string	producer		high
delete.topic.enable	Enables delete topic. Delete topic through the admin tool will have no effect if this config is turned off	boolean	false		high
host.name	hostname of broker. If this is set, it will only bind to this address. If this is not set, it will bind to all interfaces	string	''		high
leader.imbalance.check.interval.seconds	The frequency with which the partition rebalance check is triggered by the controller	long	300		high
leader.imbalance.per.broker.percentage	The ratio of leader imbalance allowed per broker. The controller would trigger a leader balance if it goes above this value per broker. The value is specified in percentage.	int	10		high
listeners	Listener List - Comma-separated list of URIs we will listen on and their protocols. Specify hostname as 0.0.0.0 to bind to all interfaces. Leave hostname empty to bind to default interface. Examples of legal listener lists: PLAINTEXT://myhost:9092,TRACE://9091 PLAINTEXT://0.0.0.0:9092,TRACE://localhost:9093	string	null		high
log.dir	The directory in which the log data is kept (supplemental for log.dirs property)	string	/tmp/kafka-logs		high
log.dirs	The directories in which the log data is kept. If not set, the value in log.dir is used	string	null		high
log.flush.interval.messages	The number of messages accumulated on a log partition before messages are flushed to disk	long	9223372036854775807	[1,...]	high
log.flush.interval.ms	The maximum time in ms that a message in any topic is kept in memory before flushed to disk. If not set, the value in	long	null		high

	log.flush.scheduler.interval.ms is used				
log.flush.offset.checkpoint.interval.ms	The frequency with which we update the persistent record of the last flush which acts as the log recovery point	int	60000	[0,...]	high
log.flush.scheduler.interval.ms	The frequency in ms that the log flusher checks whether any log needs to be flushed to disk	long	9223372036854775807		high
log.retention.bytes	The maximum size of the log before deleting it	long	-1		high
log.retention.hours	The number of hours to keep a log file before deleting it (in hours), tertiary to log.retention.ms property	int	168		high
log.retention.minutes	The number of minutes to keep a log file before deleting it (in minutes), secondary to log.retention.ms property. If not set, the value in log.retention.hours is used	int	null		high
log.retention.ms	The number of milliseconds to keep a log file before deleting it (in milliseconds), if not set, the value in log.retention.minutes is used	long	null		high
log.roll.hours	The maximum time before a new log segment is rolled out (in hours), secondary to log.roll.ms property	int	168	[1,...]	high
log.roll.jitter.hours	The maximum jitter to subtract from logRollTimeMills (in hours), secondary to log.roll.jitter.ms property	int	0	[0,...]	high
log.roll.jitter.ms	The maximum jitter to subtract from logRollTimeMills (in milliseconds), if not set, the value in log.roll.jitter.hours is used	long	null		high
log.roll.ms	The maximum time before a new log segment is rolled out (in milliseconds). If not set, the value in log.roll.hours is used	long	null		high
log.segment.bytes	The maximum size of a single log file	int	1073741824	[14,...]	high
log.segment.delete.delay.ms	The amount of time to wait before deleting a file from the filesystem	long	60000	[0,...]	high
message.max.bytes	The maximum size of message that the server can receive	int	1000012	[0,...]	high
min.insync.replicas	define the minimum number of replicas in ISR needed to satisfy a producer request with required.acks=-1 (or all)	int	1	[1,...]	high
num.io.threads	The number of io threads that the server uses for carrying out network requests	int	8	[1,...]	high
num.network.threads	the number of network threads that the server uses for handling network requests	int	3	[1,...]	high
num.recovery.threads.per.data.dir	The number of threads per data directory to be used for log recovery at startup and flushing at shutdown	int	1	[1,...]	high
num.replica.fetchers	Number of fetcher threads used to replicate messages from a source broker. Increasing this value can increase the degree of I/O parallelism in the follower broker.	int	1		high
offset.metadata.max.bytes	The maximum size for a metadata entry associated with an offset commit	int	4096		high
offsets.commit.required.acks	The required acks before the commit can be accepted. In general, the default (-1) should not be overridden	short	-1		high
offsets.commit.timeout.ms	Offset commit will be delayed until all replicas for the offsets topic receive the commit or this timeout is reached. This is similar to the producer request timeout.	int	5000	[1,...]	high
offsets.load.buffer.size	Batch size for reading from the offsets segments when loading offsets into the cache.	int	5242880	[1,...]	high
offsets.retention.check.interval.ms	Frequency at which to check for stale offsets	long	600000	[1,...]	high
offsets.retention.minutes	Log retention window in minutes for offsets topic	int	1440	[1,...]	high
offsets.topic.compression.codec	Compression codec for the offsets topic - compression may be used to achieve "atomic" commits	int	0		high
offsets.topic.num.partitions	The number of partitions for the offset commit topic (should not change after deployment)	int	50	[1,...]	high
offsets.topic.replication.factor	The replication factor for the offsets topic (set higher to ensure availability). To ensure that the effective replication factor of the offsets topic is the configured value, the number of alive brokers has to be at least the replication factor at the time of the first request for the offsets topic. If not, either the offsets topic creation will fail or it will get a replication factor of min(alive brokers, configured replication factor)	short	3	[1,...]	high
offsets.topic.segment.bytes	The offsets topic segment bytes should be kept relatively small in order to facilitate faster log compaction and cache loads	int	104857600	[1,...]	high
port	the port to listen and accept connections on	int	9092		high
queued.max.requests	The number of queued requests allowed before blocking the network threads	int	500	[1,...]	high
quota.consumer.default	Any consumer distinguished by clientId/consumer group will get throttled if it fetches more bytes than this value per second	long	9223372036854775807	[1,...]	high
quota.producer.default	Any producer distinguished by clientId will get throttled if it produces more bytes than this value per second	long	9223372036854775807	[1,...]	high
replica.fetch.max.bytes	The number of bytes of messages to attempt to fetch	int	1048576		high
replica.fetch.min.bytes	Minimum bytes expected for each fetch response. If not enough bytes, wait up to replica.fetch.wait.time.ms	int	1		high
replica.fetch.wait.max.ms	max wait time for each fetcher request issued by follower replicas. This value should always be less than the replica.lag.time.max.ms at all times to prevent frequent shrinking of ISR for low throughput topics	int	500		high
replica.high.watermark.checkpoint.interval.ms	The frequency with which the high watermark is saved out to disk	long	5000		high
replica.lag.time.max.ms	If a follower hasn't sent any fetch requests or hasn't consumed up to the leaders log end offset for at least this time, the leader will remove the follower from isr	long	10000		high
replica.socket.receive.buffer.bytes	The socket receive buffer for network requests	int	65536		high
replica.socket.timeout.ms	The socket timeout for network requests. Its value should be at least replica.fetch.wait.max.ms	int	30000		high
request.timeout.ms	The configuration controls the maximum amount of time the client will wait for the response of a request. If the response is not received before the timeout elapses the client will resend the request if necessary or fail the request if retries are exhausted.	int	30000		high
socket.receive.buffer.bytes	The SO_RCVBUF buffer of the socket server sockets	int	102400		high

socket.request.max.bytes	The maximum number of bytes in a socket request	int	104857600	[1,...]	high
socket.send.buffer.bytes	The <code>SO_SNDBUF</code> buffer of the socket server sockets	int	102400		high
unclean.leader.election.enable	Indicates whether to enable replicas not in the ISR set to be elected as leader as a last resort, even though doing so may result in data loss	boolean	true		high
zookeeper.connection.timeout.ms	The max time that the client waits to establish a connection to zookeeper. If not set, the value in <code>zookeeper.session.timeout.ms</code> is used	int	null		high
zookeeper.session.timeout.ms	Zookeeper session timeout	int	6000		high
zookeeper.set.acl	Set client to use secure ACLs	boolean	false		high
broker.id.generation.enable	Enable automatic broker id generation on the server? When enabled the value configured for <code>reserved.broker.max.id</code> should be reviewed.	boolean	true		medium
connections.max.idle.ms	Idle connections timeout: the server socket processor threads close the connections that idle more than this	long	600000		medium
controlled.shutdown.enable	Enable controlled shutdown of the server	boolean	true		medium
controlled.shutdown.max.retries	Controlled shutdown can fail for multiple reasons. This determines the number of retries when such failure happens	int	3		medium
controlled.shutdown.retry.backoff.ms	Before each retry, the system needs time to recover from the state that caused the previous failure (Controller fail over, replica lag etc). This config determines the amount of time to wait before retrying.	long	5000		medium
controller.socket.timeout.ms	The socket timeout for controller-to-broker channels	int	30000		medium
default.replication.factor	default replication factors for automatically created topics	int	1		medium
fetch.purgatory.purge.interval.requests	The purge interval (in number of requests) of the fetch request purgatory	int	1000		medium
group.max.session.timeout.ms	The maximum allowed session timeout for registered consumers	int	30000		medium
group.min.session.timeout.ms	The minimum allowed session timeout for registered consumers	int	6000		medium
inter.broker.protocol.version	Specify which version of the inter-broker protocol will be used. This is typically bumped after all brokers were upgraded to a new version. Example of some valid values are: 0.8.0, 0.8.1, 0.8.1.1, 0.8.2, 0.8.2.0, 0.8.2.1, 0.9.0.0, 0.9.0.0.1 Check AplVersion for the full list.	string	0.9.0.X		medium
log.cleaner.backoff.ms	The amount of time to sleep when there are no logs to clean	long	15000	[0,...]	medium
log.cleaner.dedupe.buffer.size	The total memory used for log deduplication across all cleaner threads	long	134217728		medium
log.cleaner.delete.retention.ms	How long are delete records retained?	long	86400000		medium
log.cleaner.enable	Enable the log cleaner process to run on the server? Should be enabled if using any topics with a cleanup.policy=compact including the internal offsets topic. If disabled those topics will not be compacted and continually grow in size.	boolean	true		medium
log.cleaner.io.buffer.load.factor	Log cleaner dedupe buffer load factor. The percentage full the dedupe buffer can become. A higher value will allow more log to be cleaned at once but will lead to more hash collisions	double	0.9		medium
log.cleaner.io.buffer.size	The total memory used for log cleaner I/O buffers across all cleaner threads	int	524288	[0,...]	medium
log.cleaner.io.max.bytes.per.second	The log cleaner will be throttled so that the sum of its read and write i/o will be less than this value on average	double	3.7976931348623151E308		medium
log.cleaner.min.cleanable.ratio	The minimum ratio of dirty log to total log for a log to be eligible for cleaning	double	0.5		medium
log.cleaner.threads	The number of background threads to use for log cleaning	int	1	[0,...]	medium
log.cleanup.policy	The default cleanup policy for segments beyond the retention window, must be either "delete" or "compact"	string	delete	[compact, delete]	medium
log.index.interval.bytes	The interval with which we add an entry to the offset index	int	4096	[0,...]	medium
log.index.size.max.bytes	The maximum size in bytes of the offset index	int	10485760	[4,...]	medium
log.preallocate	Should pre allocate file when create new segment? If you are using Kafka on Windows, you probably need to set it to true.	boolean	false		medium
log.retention.check.interval.ms	The frequency in milliseconds that the log cleaner checks whether any log is eligible for deletion	long	300000	[1,...]	medium
max.connections.per.ip	The maximum number of connections we allow from each ip address	int	2147483647	[1,...]	medium
max.connections.per.ip.overrides	Per-ip or hostname overrides to the default maximum number of connections	string	""		medium
num.partitions	The default number of log partitions per topic	int	1	[1,...]	medium
principal.builder.class	The fully qualified name of a class that implements the <code>PrincipalBuilder</code> interface, which is currently used to build the <code>Principal</code> for connections with the <code>SSL</code> <code>SecurityProtocol</code> .	class	class org.apache.kafka.common.security.auth.DefaultPrincipalBuilder		medium
producer.purgatory.purge.interval.requests	The purge interval (in number of requests) of the producer request purgatory	int	1000		medium
replica.fetch.backoff.ms	The amount of time to sleep when fetch partition error occurs	int	1000	[0,...]	medium
reserved.broker.max.id	Max number that can be used for a broker id	int	1000	[0,...]	medium
sasl.kerberos.kinit.cmd	Kerberos kinit command path.	string	/usr/bin/kinit		medium
sasl.kerberos.min.time.before.relogin	Login thread sleep time between refresh attempts.	long	60000		medium
sasl.kerberos.principal.to.local.rules	A list of rules for mapping from principal names to short names (typically operating system usernames). The rules are evaluated in order and the first rule that matches a principal name is used to map it to a short name. Any later rules in the list are ignored. By default, principal names of the form <code>{username}/{hostname}@{REALM}</code> are mapped to <code>{username}</code> . For more details on the format please see saslauthd(8) .	list	[DEFAULT]		medium
sasl.kerberos.service.name	The Kerberos principal name that Kafka runs as. This can be defined either in Kafka's JAAS config or in Kafka's config.	string	null		medium
sasl.kerberos.ticket.renew.jitter	Percentage of random jitter added to the renewal time.	double	0.05		medium
sasl.kerberos.ticket.renew.window.factor	Login thread will sleep until the specified window factor of time from last refresh to ticket's expiry has been reached, at which time it will try to renew the ticket.	double	0.8		medium
security.inter.broker.protocol	Security protocol used to communicate between brokers. Valid values are: PLAINTEXT, SSL, SASL_PLAINTEXT, SASL_SSL.	string	PLAINTEXT		medium

ssl.cipher.suites	A list of cipher suites. This is a named combination of authentication, encryption, MAC and key exchange algorithm used to negotiate the security settings for a network connection using TLS or SSL network protocol. By default all the available cipher suites are supported.	list	null		medium
ssl.client.auth	Configures kafka broker to request client authentication. The following settings are common: <ul style="list-style-type: none">• <code>ssl.client.auth=required</code> If set to required client authentication is required.• <code>ssl.client.auth=requested</code> This means client authentication is optional. unless requested, if this option is set client can choose not to provide authentication information about itself• <code>ssl.client.auth=none</code> This means client authentication is not needed.	string	none	(required, requested, none)	medium
ssl.enabled.protocols	The list of protocols enabled for SSL connections.	list	[TLSv1.2, TLSv1.1, TLSv1]		medium
ssl.key.password	The password of the private key in the key store file. This is optional for client.	password	null		medium
ssl.keymanager.algorithm	The algorithm used by key manager factory for SSL connections. Default value is the key manager factory algorithm configured for the Java Virtual Machine.	string	SunX509		medium
ssl.keystore.location	The location of the key store file. This is optional for client and can be used for two-way authentication for client.	string	null		medium
ssl.keystore.password	The store password for the key store file. This is optional for client and only needed if <code>ssl.keystore.location</code> is configured.	password	null		medium
ssl.keystore.type	The file format of the key store file. This is optional for client.	string	JKS		medium
ssl.protocol	The SSL protocol used to generate the SSLContext. Default setting is TLS, which is fine for most cases. Allowed values in recent JVMs are TLS, TLSv1.1 and TLSv1.2. SSL, SSLv2 and SSLv3 may be supported in older JVMs, but their usage is discouraged due to known security vulnerabilities.	string	TLS		medium
ssl.provider	The name of the security provider used for SSL connections. Default value is the default security provider of the JVM.	string	null		medium
ssl.trustmanager.algorithm	The algorithm used by trust manager factory for SSL connections. Default value is the trust manager factory algorithm configured for the Java Virtual Machine.	string	PKIX		medium
ssl.truststore.location	The location of the trust store file.	string	null		medium
ssl.truststore.password	The password for the trust store file.	password	null		medium
ssl.truststore.type	The file format of the trust store file.	string	JKS		medium
authorizer.class.name	The authorizer class that should be used for authorization	string	""		low
metric.reporters	A list of classes to use as metrics reporters. Implementing the <code>MetricReporter</code> interface allows plugging in classes that will be notified of new metric creation. The <code>JmxReporter</code> is always included to register JMX statistics.	list	[]		low
metrics.num.samples	The number of samples maintained to compute metrics.	int	2	[1,...]	low
metrics.sample.window.ms	The number of samples maintained to compute metrics.	long	30000	[1,...]	low
quota.window.num	The number of samples to retain in memory	int	11	[1,...]	low
quota.window.size.seconds	The time span of each sample	int	1	[1,...]	low
ssl.endpoint.identification.algorithm	The endpoint identification algorithm to validate server hostname using server certificate.	string	null		low
zookeeper.sync.time.ms	How far a ZK follower can be behind a ZK leader	int	2000		low

More details about broker configuration can be found in the scala class `kafka.server.KafkaConfig`.

Topic-level configuration Configurations pertinent to topics have both a global default as well as an optional per topic override. If no per topic configuration is given the global default is used. The override can be set at topic creation time by giving one or more `--config` options. This example creates a topic named `my-topic` with a custom max message size and flush rate:

```
> bin/kafka-topics.sh --zookeeper localhost:2181 --create --topic my-topic --partitions 1 --replication-factor 1 --config max.message.bytes=64000 --config flush.messages=1
```

Overrides can also be changed or set later using the alter topic command. This example updates the max message size for `my-topic`:

```
> bin/kafka-topics.sh --zookeeper localhost:2181 --alter --topic my-topic --config max.message.bytes=128000
```

To remove an override you can do

```
> bin/kafka-topics.sh --zookeeper localhost:2181 --alter --topic my-topic --deleteConfig max.message.bytes
```

The following are the topic-level configurations. The server's default configuration for this property is given under the Server Default Property heading, setting this default in the server config allows you to change the default given to topics that have no override specified.

#Property	Default	ServerDefault#Property	Description
cleanup.policy	delete	log.cleanup.policy	A string that is either "delete" or "compact". This string designates the retention policy to use on old log segments. The default policy ("delete") will discard old segments when their retention time or size limit has been reached. The "compact" setting will enable log compaction on the topic.
delete.retention.ms	8640000 (24 hours)	log.cleaner.delete.retention.ms	The amount of time to retain delete tombstone markers for log compacted topics. This setting also gives a bound on the time in which a consumer must complete a read if they begin from offset 0 to ensure that they get a valid snapshot of the final stage (otherwise delete tombstones may be collected before they complete their scan).
flush.messages	None	log.flush.interval.messages	This setting allows specifying an interval at which we will force an fsync of data written to the log. For example if this was set to 1 we would fsync after every message; if it were 5 we would fsync after every five messages. In general we recommend you not set this and use replication for durability and allow the operating system's background flush capabilities as it is more efficient. This setting can be overridden on a per-topic basis (see the per-topic configuration section).
flush.ms	None	log.flush.interval.ms	This setting allows specifying a time interval at which we will force an fsync of data written to the log. For example if this was set to 1000 we would fsync after 1000 ms had passed. In general we recommend you not set this and use replication for durability and allow the operating system's background flush capabilities as it is more efficient.
index.interval.bytes	4096	log.index.interval.bytes	This setting controls how frequently kafka adds an index entry to it's offset index. The default setting ensures that we index a message roughly every 4096 bytes. More indexing allows reads to jump closer to the exact position in the log but makes the index larger. You

			probably don't need to change this.
max.message.bytes	1,000,000	message.max.bytes	This is largest message size Kafka will allow to be appended to this topic. Note that if you increase this size you must also increase your consumer's fetch size so they can fetch messages this large.
min.cleanable.dirty.ratio	0.5	log.cleaner.min.cleanable.ratio	This configuration controls how frequently the log compactor will attempt to clean the log (assuming logcompaction is enabled). By default we will avoid cleaning a log where more than 50% of the log has been compacted. This ratio bounds the maximum space wasted in the log by duplicates (at 50% at most 50% of the log could be duplicates). A higher ratio will mean fewer, more efficient cleanings but will mean more wasted space in the log.
min.insync.replicas	1	min.insync.replicas	When a producer sets request.required.acks to -1, min.insync.replicas specifies the minimum number of replicas that must acknowledge a write for the write to be considered successful. If this minimum cannot be met, then the producer will raise an exception (either NotEnoughReplicas or NotEnoughReplicasAfterAppend). When used together, min.insync.replicas and request.required.acks allow you to enforce greater durability guarantees. A typical scenario would be to create a topic with a replication factor of 3, set min.insync.replicas to 2, and produce with request.required.acks of -1. This will ensure that the producer raises an exception if a majority of replicas do not receive a write.
retention.bytes	None	log.retention.bytes	This configuration controls the maximum size a log can grow to before we will discard old log segments to free up space if we are using the "delete" retention policy. By default there is no size limit only a time limit.
retention.ms	7 days	log.retention.minutes	This configuration controls the maximum time we will retain a log before we will discard old log segments to free up space if we are using the "delete" retention policy. This represents an SLA on how soon consumers must read their data.
segment.bytes	1 GB	log.segment.bytes	This configuration controls the segment file size for the log. Retention and cleaning is always done a file at a time so a larger segment size means fewer files but less granular control over retention.
segment.index.bytes	10 MB	log.index.size.max.bytes	This configuration controls the size of the index that maps offsets to file positions. We preallocate this index file and shrink it only after log rolls. You generally should not need to change this setting.
segment.ms	7 days	log.roll.hours	This configuration controls the period of time after which Kafka will force the log to roll even if the segment file isn't full to ensure that retention can delete or compact old data.
segment.jitter.ms	0	log.roll.jitter.ms.hours	The maximum jitter to subtract from log.roll.time.millis.

3.2.2 Producer Configs

Below is the configuration of the Java producer:

Name	Description	Type	Default	Valid Values	Importance
bootstrap.servers	A list of host/port pairs to use for establishing the initial connection to the Kafka cluster. The client will make use of all servers irrespective of which servers are specified here for bootstrapping—this list only impacts the initial hosts used to discover the full set of servers. This list should be in the form <code>host1:port1,host2:port2,...</code> . Since these servers are just used for the initial connection to discover the full cluster membership (which may change dynamically), this list need not contain the full set of servers (you may want more than one, though, in case a server is down).	list			high
key.serializer	Serializer class for key that implements the <code>Serializer</code> interface.	class			high
value.serializer	Serializer class for value that implements the <code>Serializer</code> interface.	class			high
acks	The number of acknowledgments the producer requires the leader to have received before considering a request complete. This controls the durability of records that are sent. The following settings are common: <ul style="list-style-type: none">• <code>acks=0</code> If set to zero then the producer will not wait for any acknowledgment from the server at all. The record will be immediately added to the socket buffer and considered sent. No guarantee can be made that the server has received the record in this case, and the <code>retry.backoff.config</code> configuration will not take effect (as the client won't generally know of any failures). The offset given back for each record will always be set to -1.• <code>acks=1</code> This will mean the leader will write the record to its local log but will respond without awaiting full acknowledgment from all followers. In this case should the leader fail immediately after acknowledging the record but before the followers have replicated it then the record will be lost.• <code>acks=all</code> This means the leader will wait for the full set of in-sync replicas to acknowledge the record. This guarantees that the record will not be lost as long as at least one in-sync replica remains alive. This is the strongest available guarantee.	string	1	[all, -1, 0, 1]	high
buffer.memory	The total bytes of memory the producer can use to buffer records waiting to be sent to the server. If records are sent faster than they can be delivered to the server the producer will either block or throw an exception based on the preference specified by <code>block.on.buffer.full</code> . This setting should correspond roughly to the total memory the producer will use, but is not a hard bound since not all memory the producer uses is used for buffering. Some additional memory will be used for compression (if compression is enabled) as well as for maintaining in-flight requests.	long	33554432	[0,...]	high

compression.type	The compression type for all data generated by the producer. The default is none (i.e. no compression). Valid values are none, gzip, snappy, or lz4. Compression is of full batches of data, so the efficacy of batching will also impact the compression ratio (more batching means better compression).	string	none		high
retries	Setting a value greater than zero will cause the client to resend any record whose send fails with a potentially transient error. Note that this retry is no different than if the client resent the record upon receiving the error. Allowing retries will potentially change the ordering of records because if two records are sent to a single partition, and the first fails and is retried but the second succeeds, then the second record may appear first.	int	0	[0,...,2147483647]	high
ssl.key.password	The password of the private key in the key store file. This is optional for client.	password	null		high
ssl.keystore.location	The location of the key store file. This is optional for client and can be used for two-way authentication for client.	string	null		high
ssl.keystore.password	The store password for the key store file. This is optional for client and only needed if ssl.keystore.location is configured.	password	null		high
ssl.truststore.location	The location of the trust store file.	string	null		high
ssl.truststore.password	The password for the trust store file.	password	null		high
batch.size	<p>The producer will attempt to batch records together into fewer requests whenever multiple records are being sent to the same partition. This helps performance on both the client and the server. This configuration controls the default batch size in bytes.</p> <p>No attempt will be made to batch records larger than this size.</p> <p>Requests sent to brokers will contain multiple batches, one for each partition with data available to be sent. A small batch size will make batching less common and may reduce throughput (a batch size of zero will disable batching entirely). A very large batch size may use memory a bit more wastefully as we will always allocate a buffer of the specified batch size in anticipation of additional records.</p>	int	16384	[0,...]	medium
client.id	An id string to pass to the server when making requests. The purpose of this is to be able to track the source of requests beyond just ip/port by allowing a logical application name to be included in server-side request logging.	string	**		medium
connections.max.idle.ms	Close idle connections after the number of milliseconds specified by this config.	long	540000		medium
linger.ms	<p>The producer groups together any records that arrive in between request transmissions into a single batched request. Normally this occurs only under load when records arrive faster than they can be sent out. However in some circumstances the client may want to reduce the number of requests even under moderate load. This setting accomplishes this by adding a small amount of artificial delay—that is, rather than immediately sending out a record the producer will wait for up to the given delay to allow other records to be sent so that the sends can be batched together. This can be thought of as analogous to Nagle's algorithm in TCP. This setting gives the upper bound on the delay for batching: once we get <code>batch.size</code> worth of records for a partition it will be sent immediately regardless of this setting, however if we have fewer than this many bytes accumulated for this partition we will "linger" for the specified time waiting for more records to show up. This setting defaults to 0 (i.e. no delay). Setting <code>linger.ms=5</code>, for example, would have the effect of reducing the number of requests sent but would add up to 5ms of latency to records sent in the absence of load.</p>	long	0	[0,...]	medium
max.block.ms	The configuration controls how long <code>(@link KafkaProducer.send())</code> and <code>(@link KafkaProducer.partitionsFor())</code> will block. These methods can be blocked either because the buffer is full or metadata unavailable. Blocking in the user-supplied serializers or partitioner will not be counted against this timeout.	long	60000	[0,...]	medium
max.request.size	The maximum size of a request. This is also effectively a cap on the maximum record size. Note that the server has its own cap on record size which may be different from this. This setting will limit the number of record batches the producer will send in a single request to avoid sending huge requests.	int	1048576	[0,...]	medium
partitioner.class	Partitioner class that implements the <code>Partitioner</code> interface.	class	class org.apache.kafka.clients.producer.internals.DefaultPartitioner		medium
receive.buffer.bytes	The size of the TCP receive buffer (SO_RECVBUF) to use when reading data.	int	32768	[0,...]	medium
request.timeout.ms	The configuration controls the maximum amount of time the client will wait for the response of a request. If the response is not received before the timeout elapses the client will	int	30000	[0,...]	medium

	resend the request if necessary or fail the request if retries are exhausted.				
<code>sasl.kerberos.service.name</code>	The Kerberos principal name that Kafka runs as. This can be defined either in Kafka's <code>JAAS</code> config or in Kafka's config.	string	null		medium
<code>security.protocol</code>	Protocol used to communicate with brokers. Valid values are: <code>PLAINTEXT</code> , <code>SSL</code> , <code>SASL_PLAINTEXT</code> , <code>SASL_SSL</code> .	string	<code>PLAINTEXT</code>		medium
<code>send.buffer.bytes</code>	The size of the TCP send buffer (<code>SO_SNDBUF</code>) to use when sending data.	int	131072	[0...]	medium
<code>ssl.enabled.protocols</code>	The list of protocols enabled for SSL connections.	list	[<code>TLSv1.2</code> , <code>TLSv1.1</code> , <code>TLSv1</code>]		medium
<code>ssl.keystore.type</code>	The file format of the key store file. This is optional for client.	string	<code>JKS</code>		medium
<code>ssl.protocol</code>	The SSL protocol used to generate the <code>SSLContext</code> . Default setting is <code>TLS</code> , which is fine for most cases. Allowed values in recent JVMs are <code>TLS</code> , <code>TLSv1.1</code> and <code>TLSv1.2</code> . <code>SSL</code> , <code>SSLv2</code> and <code>SSLv3</code> may be supported in older JVMs, but their usage is discouraged due to known security vulnerabilities.	string	<code>TLS</code>		medium
<code>ssl.provider</code>	The name of the security provider used for SSL connections. Default value is the default security provider of the JVM.	string	null		medium
<code>ssl.truststore.type</code>	The file format of the trust store file.	string	<code>JKS</code>		medium
<code>timeout.ms</code>	The configuration controls the maximum amount of time the server will wait for acknowledgments from followers to meet the acknowledgment requirements the producer has specified with the <code>acks</code> configuration. If the requested number of acknowledgments are not met when the timeout elapses an error will be returned. This timeout is measured on the server side and does not include the network latency of the request.	int	30000	[0...]	medium
<code>block.on.buffer.full</code>	When our memory buffer is exhausted we must either stop accepting new records (block) or throw errors. By default this setting is true and we block, however in some scenarios blocking is not desirable and it is better to immediately give an error. Setting this to <code>false</code> will accomplish that: the producer will throw a <code>BufferExhaustedException</code> if a record is sent and the buffer space is full.	boolean	<code>false</code>		low
<code>max.in.flight.requests.per.connection</code>	The maximum number of unacknowledged requests the client will send on a single connection before blocking. Note that if this setting is set to be greater than 1 and there are failed sends, there is a risk of message re-ordering due to retries (i.e., if retries are enabled).	int	5	[1...]	low
<code>metadata.fetch.timeout.ms</code>	The first time data is sent to a topic we must fetch metadata about that topic to know which servers host the topic's partitions. This fetch to succeed before throwing an exception back to the client.	long	60000	[0...]	low
<code>metadata.max.age.ms</code>	The period of time in milliseconds after which we force a refresh of metadata even if we haven't seen any partition leadership changes to proactively discover any new brokers or partitions.	long	300000	[0...]	low
<code>metric.reporters</code>	A list of classes to use as metrics reporters. Implementing the <code>MetricReporter</code> interface allows plugging in classes that will be notified of new metric creation. The <code>JmxReporter</code> is always included to register JMX statistics.	list	[]		low
<code>metrics.num.samples</code>	The number of samples maintained to compute metrics.	int	2	[1...]	low
<code>metrics.sample.window.ms</code>	The number of samples maintained to compute metrics.	long	3000	[0...]	low
<code>reconnect.backoff.ms</code>	The amount of time to wait before attempting to reconnect to a given host. This avoids repeatedly connecting to a host in a tight loop. This backoff applies to all requests sent by the consumer to the broker.	long	50	[0...]	low
<code>retry.backoff.ms</code>	The amount of time to wait before attempting to retry a failed fetch request to a given topic partition. This avoids repeated fetching and failing in a tight loop.	long	100	[0...]	low
<code>sasl.kerberos.kinit.cmd</code>	Kerberos kinit command path.	string	<code>/usr/bin/kinit</code>		low
<code>sasl.kerberos.min.time.before.relogin</code>	Login thread sleep time between refresh attempts.	long	60000		low
<code>sasl.kerberos.ticket.renew.jitter</code>	Percentage of random jitter added to the renewal time.	double	0.05		low
<code>sasl.kerberos.ticket.renew.window.factor</code>	Login thread will sleep until the specified window factor of time from last refresh to ticket's expiry has been reached, at which time it will try to renew the ticket.	double	0.8		low
<code>ssl.cipher.suites</code>	A list of cipher suites. This is a named combination of authentication, encryption, MAC and key exchange algorithm used to negotiate the security settings for a network connection using TLS or SSL network protocol. By default all the available cipher suites are supported.	list	null		low
<code>ssl.endpoint.identification.algorithm</code>	The endpoint identification algorithm to validate server hostname using server certificate.	string	null		low
<code>ssl.keymanager.algorithm</code>	The algorithm used by key manager factory for SSL connections. Default value is the key manager factory algorithm configured for the Java Virtual Machine.	string	<code>SunX509</code>		low
<code>ssl.trustmanager.algorithm</code>	The algorithm used by trust manager factory for SSL connections. Default value is the trust manager factory	string	<code>PKIX</code>		low

	algorithm configured for the Java Virtual Machine.				
--	--	--	--	--	--

For those interested in the legacy Scala producer configs, information can be found [here](#):

3.3.3ConsumerConfigs

We introduce both the old 0.8 consumer configs and the new consumer configs respectively below.

3.3.3.1OldConsumerConfigs

The essential old consumer configurations are the following:

- group.id
- zookeeper.connect

Property	Default	Description
group.id		A string that uniquely identifies the group of consumer processes to which this consumer belongs. By setting the same group id multiple processes indicate that they are all part of the same consumer group.
zookeeper.connect		Specifies the Zookeeper connection string in the form <code>hostname:port</code> where <code>host</code> and <code>port</code> are the host and port of a Zookeeper server. To allow connecting through other Zookeeper nodes when that Zookeeper machine is down you can also specify multiple hosts in the form <code>hostname1:port1,hostname2:port2,hostname3:port3</code> . The server may also have a Zookeeper chroot path as part of it's Zookeeper connection string which puts its data under some path in the global Zookeeper namespace. If so the consumer should use the same chroot path in its connection string. For example to give a chroot path of <code>/chroot/path</code> you would give the connection string as <code>hostname1:port1,hostname2:port2,hostname3:port3/chroot/path</code> .
consumer.id	null	Generated automatically if not set.
socket.timeout.ms	30 * 1000	The socket timeout for network requests. The actual timeout set will be <code>max.fetch.wait + socket.timeout.ms</code> .
socket.receive.buffer.bytes	64 * 1024	The socket receive buffer for network requests
fetch.message.max.bytes	1024 * 1024	The number of bytes of messages to attempt to fetch for each topic-partition in each fetch request. These bytes will be read into memory for each partition, so this helps control the memory used by the consumer. The fetch request size must be at least as large as the maximum message size the server allows or else it is possible for the producer to send messages larger than the consumer can fetch.
num.consumer.fetchers	1	The number fetcher threads used to fetch data.
auto.commit.enable	true	If true, periodically commit to Zookeeper the offset of messages already fetched by the consumer. This committed offset will be used when the process fails as the position from which the new consumer will begin.
auto.commit.interval.ms	60 * 1000	The frequency in ms that the consumer offsets are committed to zookeeper.
queued.max.message.chunks	2	Max number of message chunks buffered for consumption. Each chunk can be up to <code>fetch.message.max.bytes</code> .
rebalance.max.retries	4	When a new consumer joins a consumer group the set of consumers attempt to "rebalance" the load to assign partitions to each consumer. If the set of consumers changes while this assignment is taking place the rebalance will fail and retry. This setting controls the maximum number of attempts before giving up.
fetch.min.bytes	1	The minimum amount of data the server should return for a fetch request. If insufficient data is available the request will wait for that much data to accumulate before answering the request.
fetch.wait.max.ms	100	The maximum amount of time the server will block before answering the fetch request if there isn't sufficient data to immediately satisfy <code>fetch.min.bytes</code>
rebalance.backoff.ms	2000	Backoff time between retries during rebalance. If not set explicitly, the value in <code>zookeeper.sync.time.ms</code> is used.
refresh.leader.backoff.ms	200	Backoff time to wait before trying to determine the leader of a partition that has just lost its leader.
auto.offset.reset	largest	What to do when there is no initial offset in Zookeeper or if an offset is out of range: * <code>smallest</code> : automatically reset the offset to the smallest offset * <code>largest</code> : automatically reset the offset to the largest offset * anything else: throw exception to the consumer
consumer.timeout.ms	-1	Throw a timeout exception to the consumer if no message is available for consumption after the specified interval
exclude.internal.topics	true	Whether messages from internal topics (such as offsets) should be exposed to the consumer.
client.id	group id value	The client id is a user-specified string sent in each request to help trace calls. It should logically identify the application making the request.
zookeeper.session.timeout.ms	6000	Zookeeper session timeout. If the consumer fails to heartbeat to Zookeeper for this period of time it is considered dead and a rebalance will occur.
zookeeper.connection.timeout.ms	6000	The max time that the client waits while establishing a connection to zookeeper.
zookeeper.sync.time.ms	2000	How far a ZX follower can be behind a ZX leader
offsets.storage	zookeeper	Select where offsets should be stored (zookeeper or kafka).
offsets.channel.backoff.ms	1000	The backoff period when reconnecting the offsets channel or retrying failed offset fetch/commit requests.
offsets.channel.socket.timeout.ms	10000	Socket timeout when reading responses for offset fetch/commit requests. This timeout is also used for <code>ConsumerMetadata</code> requests that are used to query for the offset manager.
offsets.commit.max.retries	5	Retry the offset commit up to this many times on failure. This retry count only applies to offset commits during shut-down. It does not apply to commits originating from the auto-commit thread. It also does not apply to attempts to query for the offset coordinator before committing offsets, i.e., if a consumer metadata request fails for any reason, it will be retried and that retry does not count toward this limit.
dual.commit.enabled	true	If you are using "kafka" as offsets storage, you can dual commit offsets to Zookeeper (in addition to Kafka). This is required during migration from zookeeper based offset storage to kafka-based offset storage. With respect to any given consumer group, it is safe to turn this off after all instances within that group have been migrated to the new version that commits offsets to the broker (instead of directly to ZooKeeper).
partition.assignment.strategy	range	Select between the "range" or "roundrobin" strategy for assigning partitions to consumer streams. The round-robin partition assignor lays out all the available partitions and all the available consumer threads. It then proceeds to do a round-robin assignment from partition to consumer thread. If the subscriptions of all consumer instances are identical, then the partitions will be uniformly distributed. (i.e., the partition ownership counts will be within a delta of exactly one across all consumer threads.) Round-robin assignment is permitted only if: (a) Every topic has the same number of streams within a consumer instance (b) The set of subscribed topics is identical for every consumer instance within the group. Range partitioning works on a per-topic basis. For each topic, we lay out the available partitions in numeric order and the consumer threads in lexicographic order. We then divide the number of partitions by the total number of consumer streams (threads) to determine the number of partitions to assign to each consumer. If it does not evenly divide, then the first few consumers will have one extra partition.

More details about consumer configuration can be found in the scala class

`kafka.consumer.ConsumerConfig`.

3.3.3.2NewConsumerConfigs

Since 0.9.0.0 we have been working on a replacement for our existing simple and high-level consumers. The code is considered beta quality. Below is the configuration for the new consumer:

Name	Description	Type	Default	Valid Values	Importance
bootstrap.servers	A list of host/port pairs to use for establishing the initial connection to the Kafka cluster. The client will make use of all servers irrespective of which servers are specified here for bootstrapping--this list only impacts	list			high

	the initial hosts used to discover the full set of servers. This list should be in the form <code>host1:port1,host2:port2,...</code> Since these servers are just used for the initial connection to discover the full cluster membership (which may change dynamically), this list need not contain the full set of servers (you may want more than one, though, in case a server is down).				
<code>key.deserializer</code>	Deserializer class for key that implements the <code>Deserializer</code> interface.	class			high
<code>value.deserializer</code>	Deserializer class for value that implements the <code>Deserializer</code> interface.	class			high
<code>fetch.min.bytes</code>	The minimum amount of data the server should return for a fetch request. If insufficient data is available the request will wait for that much data to accumulate before answering the request. The default setting of 1 byte means that fetch requests are answered as soon as a single byte of data is available or the fetch request times out waiting for data to arrive. Setting this to something greater than 1 will cause the server to wait for larger amounts of data to accumulate which can improve server throughput a bit at the cost of some additional latency.	int	1	[0..]	high
<code>group.id</code>	A unique string that identifies the consumer group this consumer belongs to. This property is required if the consumer uses either the group management functionality by using <code>subscribe(topic)</code> or the Kafka-based offset management strategy.	string	""		high
<code>heartbeat.interval.ms</code>	The expected time between heartbeats to the consumer coordinator when using Kafka's group management facilities. Heartbeats are used to ensure that the consumer's session stays active and to facilitate rebalancing when new consumers join or leave the group. The value must be set lower than <code>session.timeout.ms</code> , but typically should be set no higher than 1/3 of that value. It can be adjusted even lower to control the expected time for normal rebalances.	int	3000		high
<code>max.partition.fetch.bytes</code>	The maximum amount of data per-partition the server will return. The maximum total memory used for a request will be <code>#partitions * max.partition.fetch.bytes</code> . This size must be at least as large as the maximum message size the server allows or else it is possible for the producer to send messages larger than the consumer can fetch. If that happens, the consumer can get stuck trying to fetch a large message on a certain partition.	int	1048576	[0..]	high
<code>session.timeout.ms</code>	The timeout used to detect failures when using Kafka's group management facilities.	int	30000		high
<code>ssl.key.password</code>	The password of the private key in the key store file. This is optional for client.	password	null		high
<code>ssl.keystore.location</code>	The location of the key store file. This is optional for client and can be used for two-way authentication for client.	string	null		high
<code>ssl.keystore.password</code>	The store password for the key store file. This is optional for client and only needed if ssl.keystore.location is configured.	password	null		high
<code>ssl.truststore.location</code>	The location of the trust store file.	string	null		high
<code>ssl.truststore.password</code>	The password for the trust store file.	password	null		high
<code>auto.offset.reset</code>	What to do when there is no initial offset in Kafka or if the current offset does not exist any more on the server (e.g. because that data has been deleted): <ul style="list-style-type: none"> earliest: automatically reset the offset to the earliest offset latest: automatically reset the offset to the latest offset none: throw exception to the consumer if no previous offset is found for the consumer's group anything else: throw exception to the consumer. 	string	latest	[latest, earliest, none]	medium
<code>connections.max.idle.ms</code>	Close idle connections after the number of milliseconds specified by this config.	long	540000		medium
<code>enable.auto.commit</code>	If true the consumer's offset will be periodically committed in the background.	boolean	true		medium
<code>partition.assignment.strategy</code>	The class name of the partition assignment strategy that the client will use to distribute partition ownership amongst consumer instances when group management is used.	list	[org.apache.kafka.clients.consumer.RangeAssignor]		medium
<code>receive.buffer.bytes</code>	The size of the TCP receive buffer (SO_RCVBUF) to use when reading data.	int	32768	[0..]	medium
<code>request.timeout.ms</code>	The configuration controls the maximum amount of time the client will wait for the response of a request. If the response is not received before the timeout elapses the client will resend the request if necessary or fail the request if retries are exhausted.	int	40000	[0..]	medium
<code>sasl.kerberos.service.name</code>	The Kerberos principal name that Kafka runs as. This can be defined either in Kafka's JAAS config or in Kafka's config.	string	null		medium
<code>security.protocol</code>	Protocol used to communicate with	string	PLAINTEXT		medium

	brokers. Valid values are: PLAINTEXT, SSL, SASL_PLAINTEXT, SASL_SSL.				
send.buffer.bytes	The size of the TCP send buffer (SO_SNDBUF) to use when sending data.	int	131072	[0...]	medium
ssl.enabled.protocols	The list of protocols enabled for SSL connections.	list	[TLSv1.2, TLSv1.1, TLSv1]		medium
ssl.keystore.type	The file format of the key store file. This is optional for client.	string	JKS		medium
ssl.protocol	The SSL protocol used to generate the SSLContext. Default setting is TLS, which is fine for most cases. Allowed values in recent JVMs are TLS, TLSv1.1 and TLSv1.2. SSL, SSLv2 and SSLv3 may be supported in older JVMs, but their usage is discouraged due to known security vulnerabilities.	string	TLS		medium
ssl.provider	The name of the security provider used for SSL connections. Default value is the default security provider of the JVM.	string	null		medium
ssl.truststore.type	The file format of the trust store file.	string	JKS		medium
auto.commit.interval.ms	The frequency in milliseconds that the consumer offsets are auto-committed to Kafka if <code>enable.auto.commit</code> is set to <code>true</code> .	long	5000	[0...]	low
check.crcs	Automatically check the CRC32 of the records consumed. This ensures no on-the-wire or on-disk corruption to the messages occurred. This check adds some overhead, so it may be disabled in cases seeking extreme performance.	boolean	true		low
client.id	An id string to pass to the server when making requests. The purpose of this is to be able to track the source of requests beyond just ip:port by allowing a logical application name to be included in server-side request logging.	string	""		low
fetch.max.wait.ms	The maximum amount of time the server will block before answering the fetch request if there isn't sufficient data to immediately satisfy the requirement given by <code>fetch.min.bytes</code> .	int	500	[0...]	low
metadata.max.age.ms	The period of time in milliseconds after which we force a refresh of metadata even if we haven't seen any partition leadership changes to proactively discover any new brokers or partitions.	long	300000	[0...]	low
metric.reporters	A list of classes to use as metrics reporters. Implementing the <code>MetricReporter</code> interface allows plugging in classes that will be notified of new metric creation. The <code>JmxReporter</code> is always included to register JMX statistics.	list	[]		low
metrics.num.samples	The number of samples maintained to compute metrics.	int	2	[1...]	low
metrics.sample.window.ms	The number of samples maintained to compute metrics.	long	30000	[0...]	low
reconnect.backoff.ms	The amount of time to wait before attempting to reconnect to a given host. This avoids repeatedly connecting to a host in a tight loop. This backoff applies to all requests sent by the consumer to the broker.	long	50	[0...]	low
retry.backoff.ms	The amount of time to wait before attempting to retry a failed fetch request to a given topic partition. This avoids repeated fetching and failing in a tight loop.	long	100	[0...]	low
ssl.kerberos.kinit.cmd	Kerberos kinit command path.	string	/usr/bin/kinit		low
ssl.kerberos.min.time.before.relogin	Login thread sleep time between refresh attempts.	long	60000		low
ssl.kerberos.ticket.renew.jitter	Percentage of random jitter added to the renewal time.	double	0.05		low
ssl.kerberos.ticket.renew.window.factor	Login thread will sleep until the specified window factor of time from last refresh to ticket's expiry has been reached, at which time it will try to renew the ticket.	double	0.8		low
ssl.cipher.suites	A list of cipher suites. This is a named combination of authentication, encryption, MAC and key exchange algorithm used to negotiate the security settings for a network connection using TLS or SSL network protocol. By default all the available cipher suites are supported.	list	null		low
ssl.endpoint.identification.algorithm	The endpoint identification algorithm to validate server hostname using server certificate.	string	null		low
ssl.keymanager.algorithm	The algorithm used by key manager factory for SSL connections. Default value is the key manager factory algorithm configured for the Java Virtual Machine.	string	SunX509		low
ssl.trustmanager.algorithm	The algorithm used by trust manager factory for SSL connections. Default value is the trust manager factory algorithm configured for the Java Virtual Machine.	string	PKIX		low

3.4.4 Kafka Connect Configs

Name	Description	Type	Default	Stable Value	Importance
group.id	A unique string that identifies the Connect cluster group this worker belongs to.	string			high
internal.key.converter	Converter class for internal key Connect data that implements the <code>Converter</code> interface. Used for converting data like offsets and configs.	class			high
internal.value.converter	Converter class for offset value Connect data that implements the <code>Converter</code> interface. Used for converting data like offsets and configs.	class			high
key.converter	Converter class for key Connect data that implements the <code>Converter</code> interface.	class			high
value.converter	Converter class for value Connect data that implements the <code>Converter</code> interface.	class			high

	implements the <code>Converter</code> interface.				
bootstrap.servers	A list of host/port pairs to use for establishing the initial connection to the Kafka cluster. The client will make use of all servers irrespective of which servers are specified here for bootstrapping—this list only impacts the initial hosts used to discover the full set of servers. This list should be in the form <code>host1:port1,host2:port2,...</code> . Since these servers are just used for the initial connection to discover the full cluster membership (which may change dynamically), this list need not contain the full set of servers (you may want more than one, though, in case a server is down).	list	[localhost:9092]	high	
cluster	ID for this cluster, which is used to provide a namespace so multiple Kafka Connect clusters or instances may co-exist while sharing a single Kafka cluster.	string	connect	high	
heartbeat.interval.ms	The expected time between heartbeats to the group coordinator when using Kafka's group management facilities. Heartbeats are used to ensure that the worker's session stays active and to facilitate rebalancing when new members join or leave the group. The value must be set lower than <code>session.timeout.ms</code> , but typically should be set no higher than 1/3 of that value. It can be adjusted even lower to control the expected time for normal rebalances.	int	3000	high	
session.timeout.ms	The timeout used to detect failures when using Kafka's group management facilities.	int	30000	high	
ssl.key.password	The password of the private key in the key store file. This is optional for client.	password	null	high	
ssl.keystore.location	The location of the key store file. This is optional for client and can be used for two-way authentication for client.	string	null	high	
ssl.keystore.password	The store password for the key store file. This is optional for client and only needed if <code>ssl.keystore.location</code> is configured.	password	null	high	
ssl.truststore.location	The location of the trust store file.	string	null	high	
ssl.truststore.password	The password for the trust store file.	password	null	high	
connections.max.idle.ms	Close idle connections after the number of milliseconds specified by this config.	long	540000	medium	
receive.buffer.bytes	The size of the TCP receive buffer (<code>SO_RCVBUF</code>) to use when reading data.	int	32768	[0..]	medium
request.timeout.ms	The configuration controls the maximum amount of time the client will wait for the response of a request. If the response is not received before the timeout elapses the client will resend the request if necessary or fail the request if retries are exhausted.	int	40000	[0..]	medium
ssl.kerberos.service.name	The Kerberos principal name that Kafka runs as. This can be defined either in Kafka's JAAS config or in Kafka's config.	string	null	medium	
security.protocol	Protocol used to communicate with brokers. Valid values are: PLAINTEXT, SSL, SASL_PLAINTEXT, SASL_SSL.	string	PLAINTEXT	medium	
send.buffer.bytes	The size of the TCP send buffer (<code>SO_SNDBUF</code>) to use when sending data.	int	131072	[0..]	medium
ssl.enabled.protocols	The list of protocols enabled for SSL connections.	list	[TLSv1.2, TLSv1.1, TLSv1]	medium	
ssl.keystore.type	The file format of the key store file. This is optional for client.	string	JKS	medium	
ssl.protocol	The SSL protocol used to generate the <code>SSLContext</code> . Default setting is TLS, which is fine for most cases. Allowed values in recent JVMs are TLS, TLSv1.1 and TLSv1.2. SSL, SSLv2 and SSLv3 may be supported in older JVMs, but their usage is discouraged due to known security vulnerabilities.	string	TLS	medium	
ssl.provider	The name of the security provider used for SSL connections. Default value is the default security provider of the JVM.	string	null	medium	
ssl.truststore.type	The file format of the trust store file.	string	JKS	medium	
worker.sync.timeout.ms	When the worker is out of sync with other workers and needs to resynchronize configurations, wait up to this amount of time before giving up, leaving the group, and waiting a backoff period before rejoining.	int	3000	medium	
worker.unsync.backoff.ms	When the worker is out of sync with other workers and fails to catch up within <code>worker.sync.timeout.ms</code> , leave the Connect cluster for this long before rejoining.	int	300000	medium	
client.id	An id string to pass to the server when making requests. The purpose of this is to be able to track the source of requests beyond just ip/port by allowing a logical application name to be included in server-side request logging.	string	""	low	
metadata.max.age.ms	The period of time in milliseconds after which we force a refresh of metadata even if we haven't seen any partition leadership changes to proactively discover any new brokers or partitions.	long	300000	[0..]	low
metric.reporters	A list of classes to use as metrics reporters. Implementing the <code>MetricReporter</code> interface allows plugging in classes that will be notified of new metric creation. The <code>JmxReporter</code> is always included to register JMX statistics.	list	[]	low	
metrics.num.samples	The number of samples maintained to compute metrics.	int	2	[1..]	low
metrics.sample.window.ms	The number of samples maintained to compute metrics.	long	30000	[0..]	low
offset.flush.interval.ms	Interval at which to try committing offsets for tasks.	long	60000	low	
offset.flush.timeout.ms	Maximum number of milliseconds to wait for records to flush and partition offset data to be committed to offset storage before cancelling the process and restoring the offset data to be committed in a future attempt.	long	5000	low	
reconnect.backoff.ms	The amount of time to wait before attempting to reconnect to a given host. This avoids repeatedly connecting to a host in a tight loop. This backoff applies to all requests sent by the consumer to the broker.	long	50	[0..]	low
rest.advertised.host.name	If this is set, this is the hostname that will be given out to other workers to connect to.	string	null	low	
rest.advertised.port	If this is set, this is the port that will be given out to other workers to connect to.	int	null	low	

rest.host.name	Hostname for the REST API. If this is set, it will only bind to this interface.	string	null		low
rest.port	Port for the REST API to listen on.	int	8083		low
retry.backoff.ms	The amount of time to wait before attempting to retry a failed fetch request to a given topic partition. This avoids repeated fetching, and falling in a tight loop.	long	100	[0,...]	low
saai.kerberos.kinit.cmd	Kerberos kinit command path.	string	/usr/bin/kinit		low
saai.kerberos.min.time.before.relogin	Login thread sleep time between refresh attempts.	long	60000		low
saai.kerberos.ticket.renew.jitter	Percentage of random jitter added to the renewal time.	double	0.05		low
saai.kerberos.ticket.renew.window.factor	Login thread will sleep until the specified window factor of time from last refresh to ticket's expiry has been reached, at which time it will try to renew the ticket.	double	0.8		low
ssl.cipher.suites	A list of cipher suites. This is a named combination of authentication, encryption, MAC, and key exchange algorithm used to negotiate the security settings for a network connection using TLS or SSL network protocol. By default all the available cipher suites are supported.	list	null		low
ssl.endpoint.identification.algorithm	The endpoint identification algorithm to validate server hostname using server certificate.	string	null		low
ssl.keymanager.algorithm	The algorithm used by key manager factory for SSL connections. Default value is the key manager factory algorithm configured for the Java Virtual Machine.	string	SunX509		low
ssl.trustmanager.algorithm	The algorithm used by trust manager factory for SSL connections. Default value is the trust manager factory algorithm configured for the Java Virtual Machine.	string	PKIX		low
task.shutdown.graceful.timeout.ms	Amount of time to wait for tasks to shutdown gracefully. This is the total amount of time, not per task. All task have shutdown triggered, then they are waited on sequentially.	long	5000		low

4. Design

4.1 Motivation

We designed Kafka to be able to act as a unified platform for handling all the real-time data feeds [a large company might have](#). To do this we had to think through a fairly broad set of use cases.

It would have to have high throughput to support high volume event streams such as real time log aggregation.

It would need to deal gracefully with large data backlogs to be able to support periodic data loads from offline systems.

It also meant the system would have to handle low-latency delivery to handle more traditional messaging use-cases.

We wanted to support partitioned, distributed, real-time processing of these feeds to create new, derived feeds. This motivated our partitioning and consumer model.

Finally in cases where the stream is fed into other data systems for serving, we knew the system would have to be able to guarantee fault tolerance in the presence of machine failures.

Supporting these uses led us to a design with a number of unique elements, more akin to a database log than a traditional messaging system. We will outline some elements of the design in the following sections.

4.2 Persistence

Don't leave the filesystem!

Kafka relies heavily on the filesystem for storing and caching messages. There is a general perception that "disks are slow" which makes people skeptical that a persistent structure can offer competitive performance. In fact disks are both much slower and much faster than people expect depending on how they are used; and a properly designed disk structure can often be as fast as the network.

The key fact about disk performance is that the throughput of hard drives has been diverging from the latency of a disk seek for the last decade. As a result the performance of linear writes on a [JBOD](#) configuration with six 7200rpm SATA RAID-5 array is about 600MB/sec but the performance of random writes is only about 100k/sec—a difference of over 600X. These linear reads and writes are the most predictable of all usage patterns, and are heavily optimized by the operating system. A modern operating system provides read-ahead and write-behind techniques that prefetch data in large block multiples and group smaller logical writes into large physical writes. A further discussion of this issue can be found in this [H3M@Amazon](#) article; they actually find that *sequential disk access can in some cases be faster than random memory access!*

To compensate for this performance divergence modern operating systems have become increasingly aggressive in their use of main memory for disk caching. A modern OS will happily divert *all* free memory to disk caching with little performance penalty when the memory is reclaimed. All disk reads and writes will go through this unified cache. This feature cannot easily be turned off without using direct I/O, so even if a process maintains an in-process cache of the data, this data will likely be duplicated in OS pagecache, effectively storing everything twice.

Furthermore we are building on top of the JVM, and anyone who has spent any time with Java memory usage knows two things:

1. The memory overhead of objects is very high, often doubling the size of the data stored (or worse).
2. Java garbage collection becomes increasingly fiddly and slow as the in-heap data increases.

As a result of these factors using the filesystem and relying on pagecache is superior to maintaining an in-memory cache or other structure—we at least double the available cache by having automatic access to all free memory, and likely double again by storing a compact byte structure rather than individual objects. Doing so will result in a cache of up to 28-30GB on a 32GB machine without GC penalties. Furthermore this cache will stay warm even if the service is restarted, whereas the in-process cache will need to be rebuilt in memory (which for a 10GB cache may take 10 minutes) or else it will need to start with a completely cold cache (which likely means terrible initial performance). This also greatly simplifies the code as all logic for maintaining coherency between the cache and filesystem is now in the OS, which tends to do so more efficiently and more correctly than one-off in-process attempts. If your disk usage favors linear reads then read-ahead is effectively pre-populating this cache with useful data on each disk read.

This suggests a design which is very simple: rather than maintain as much as possible in-memory and flush it all out to the filesystem in a panic when we run out of space, we invert that. All data is immediately written to a persistent log on the filesystem without necessarily flushing to disk. In effect this just means that it is transferred into the kernel's pagecache.

This style of pagecache-centric design is described in an [article](#) on the design of Varnish here (along with a healthy dose of arrogance).

Constant Time Suffices

The persistent data structure used in messaging systems are often a per-consumer queue with an associated BTree or other general-purpose random access data structures to maintain metadata about messages. BTrees are the most versatile data structure available, and make it possible to support a wide variety of transactional and non-transactional semantics in the messaging system. They do come with a fairly high cost, though: Btree operations are $O(\log N)$. Normally $O(\log N)$ is considered essentially equivalent to constant time, but this is not true for disk operations. Disk seeks come at 10 ms a pop, and each disk can do only one seek at a time so parallelism is limited. Hence even a handful of disk seeks leads to very high overhead. Since storage systems mix very fast cached operations with very slow physical disk operations, the observed performance of tree structures is often superlinear as data increases with fixed cache—i.e. doubling your data makes things much worse than twice as slow.

Intuitively a persistent queue could be built on simple reads and appends to files as is commonly the case with logging solutions. This structure has the advantage that all operations are $O(1)$ and reads do not block writes or each other. This has obvious performance advantages since the performance is completely decoupled from the data size—one server can now take full advantage of a number of cheap, low-rotational speed 1-TB SATA drives. Though they have poor seek performance, these drives have acceptable performance for large reads

and writes and come at 1/3 the price and 3x the capacity.

Having access to virtually unlimited disk space without any performance penalty means that we can provide some features not usually found in a messaging system. For example, in Kafka, instead of attempting to delete messages as soon as they are consumed, we can retain messages for a relative long period (say a week). This leads to a great deal of flexibility for consumers, as we will describe.

4.3 Efficiency

We have put significant effort into efficiency. One of our primary use cases is handling web activity data, which is very high volume: each page view may generate dozens of writes. Furthermore we assume each message published is read by at least one consumer (often many), hence we strive to make consumption as cheap as possible.

We have also found, from experience building and running a number of similar systems, that efficiency is a key to effective multi-tenant operations. If the downstream infrastructure service can easily become a bottleneck due to a small bump in usage by the application, such small changes will often create problems. By being very fast we help ensure that the application will tip-over under load before the infrastructure. This is particularly important when trying to run a centralized service that supports dozens or hundreds of applications on a centralized cluster as changes in usage patterns are a near-daily occurrence.

We discussed disk efficiency in the previous section. Once poor disk access patterns have been eliminated, there are two common causes of inefficiency in this type of system: too many small I/O operations, and excessive byte copying.

The small I/O problem happens both between the client and the server and in the server's own persistent operations.

To avoid this, our protocol is built around a "message set" abstraction that naturally groups messages together. This allows network requests to group messages together and amortize the overhead of the network roundtrip rather than sending a single message at a time. The server in turn appends chunks of messages to its log in one go, and the consumer fetches large linear chunks at a time.

This simple optimization produces orders of magnitude speed up. Batching leads to larger network packets, larger sequential disk operations, contiguous memory blocks, and so on, all of which allows Kafka to turn a bursty stream of random message writes into linear writes that flow to the consumers.

The other inefficiency is in byte copying. At low message rates this is not an issue, but under load the impact is significant. To avoid this we employ a standardized binary message format that is shared by the producer, the broker, and the consumer (so data chunks can be transferred without modification between them).

The message log maintained by the broker is itself just a directory of files, each populated by a sequence of message sets that have been written to disk in the same format used by the producer and consumer. Maintaining this common format allows optimization of the most important operation: network transfer of persistent log chunks. Modern unix operating systems offer a highly optimized code path for transferring data out of pagecache to a socket; in Linux this is done with the `sendfile` system call.

To understand the impact of `sendfile`, it is important to understand the common data path for transfer of data from file to socket:

1. The operating system reads data from the disk into pagecache in kernel space
2. The application reads the data from kernel space into a user-space buffer
3. The application writes the data back into kernel space into a socket buffer
4. The operating system copies the data from the socket buffer to the NIC buffer where it is sent over the network

This is clearly inefficient, there are four copies and two system calls. Using `sendfile`, this re-copying is avoided by allowing the OS to send the data from pagecache to the network directly. So in this optimized path, only the final copy to the NIC buffer is needed.

We expect a common use case to be multiple consumers on a topic. Using the zero-copy optimization above, data is copied into pagecache exactly once and reused on each consumption instead of being stored in memory and copied out to kernel space every time it is read. This allows messages to be consumed at a rate that approaches the limit of the network connection.

This combination of pagecache and `sendfile` means that on a Kafka cluster where the consumers are mostly caught up you will see no read activity on the disks whatsoever as they will be serving data entirely from cache.

For more background on the `sendfile` and zero-copy support in Java, see [this article](#)

End-to-end Batch Compression

In some cases the bottleneck is actually not CPU or disk but network bandwidth. This is particularly true for a data pipeline that needs to send messages between data centers over a wide-area network. Of course the user can always compress its messages one at a time without any support needed from Kafka, but this can lead to very poor compression ratios as much of the redundancy is due to repetition between messages of the same type (e.g. field names in JSON or user agents in web logs or common string values). Efficient compression requires compressing multiple messages together rather than compressing each message individually.

Kafka supports this by allowing recursive message sets. A batch of messages can be clumped together and compressed and sent to the server in this form. This batch of messages will be written in compressed form and will remain compressed in the log and will only be decompressed by the consumer.

Kafka supports GZIP and Snappy compression protocols. More details on compression can be found [here](#).

4.4 The Producer

Load Balancing

The producer sends data directly to the broker that is the leader for the partition without any intervening routing tier. To help the producer do this all Kafka nodes can answer a request for metadata about which servers are alive and where the leaders for the partitions of a topic are at any given time to allow the producer to appropriately direct its requests.

The client controls which partition it publishes messages to. This can be done at random, implementing a kind of random load balancing, or it can be done by some semantic partitioning function. We expose the interface for semantic partitioning by allowing the user to specify a key to partition by and using this to hash to a partition (there is also an option to override the partition function if need be). For example if the key chosen was a user id then all data for a given user would be sent to the same partition. This in turn will allow consumers to make locality assumptions about their consumption. This style of partitioning is explicitly designed to allow locality-sensitive processing in consumers.

Asynchronous Send

Batching is one of the big drivers of efficiency, and to enable batching the Kafka producer will attempt to accumulate data in memory and to send out larger batches in a single request. The batching can be configured to accumulate no more than a fixed number of messages and to wait no longer than some fixed latency bound (say 64k or 10 ms). This allows the accumulation of more bytes to send, and few larger I/O operations on the servers. This buffering is configurable and gives a mechanism to trade off a small amount of additional latency for better throughput.

Details on `configuration` and `api` for the producer can be found elsewhere in the documentation.

4.5 The Consumer

The Kafka consumer works by issuing "fetch" requests to the brokers leading the partitions it wants to consume. The consumer specifies its offset in the log with each request and receives back a chunk of log beginning from that position. The consumer thus has significant control over this position and can rewind it to re-consume data if need be.

Push vs. Pull

An initial question we considered is whether consumers should pull data from brokers or brokers should push data to the consumer. In this respect Kafka follows a more traditional design, shared by most messaging systems, where data is pushed to the broker from the producer and pulled from the broker by the consumer.

Some logging-centric systems, such as `Kerberos` and `Apache Flume` follow a very different push based path where data is pushed downstream. There are pros and cons to both approaches. However a push-based system has difficulty dealing with diverse consumers as the broker controls the rate at which data is transferred. The goal is generally for the consumer to be able to consume at the maximum possible rate; unfortunately in a push system this means the consumer tends to be overwhelmed when its rate of consumption falls below the rate of production (a denial of service attack, in essence). A pull-based system has the nicer property that the consumer simply falls behind and catches up when it can. This can be mitigated with some kind of backoff protocol by which the consumer can indicate it is overwhelmed, but getting the rate of transfer to fully utilize (but never over-utilize) the consumer is trickier than it seems. Previous attempts at building systems in this fashion led us to go with a more traditional pull model.

Another advantage of a pull-based system is that it lends itself to aggressive batching of data sent to the

consumer. A push-based system must choose to either send a request immediately or accumulate more data and then send it later without knowledge of whether the downstream consumer will be able to immediately process it. If tuned for low latency this will result in sending a single message at a time only for the transfer to end up being buffered anyway, which is wasteful. A pull-based design fixes this as the consumer always pulls all available messages after its current position in the log (or up to some configurable max size). So one gets optimal batching without introducing unnecessary latency.

The deficiency of a naive pull-based system is that if the broker has no data the consumer may end up polling in a tight loop, effectively busy-waiting for data to arrive. To avoid this we have parameters in our poll request that allow the consumer request to block in a "long poll" waiting until data arrives (and optionally waiting until a given number of bytes is available to ensure large transfer sizes).

You could imagine other possible designs which would be only pull, end-to-end. The producer would locally write to a local log, and brokers would pull from that with consumers pulling from them. A similar type of "store-and-forward" producer is often proposed. This is intriguing but we felt not very suitable for our target use cases which have thousands of producers. Our experience running persistent data systems at scale led us to feel that involving thousands of disks in the system across many applications would not actually make things more reliable and would be a nightmare to operate. And in practice we have found that we can run a pipeline with strong SLAs at large scale without a need for producer persistence.

Consumer Position

Keeping track of what has been consumed, is, surprisingly, one of the key performance points of a messaging system.

Most messaging systems keep metadata about what messages have been consumed on the broker. That is, as a message is handed out to a consumer, the broker either records that fact locally immediately or it may wait for acknowledgement from the consumer. This is a fairly intuitive choice, and indeed for a single machine server it is not clear where else this state could go. Since the data structure used for storage in many messaging systems scale poorly, this is also a pragmatic choice—since the broker knows what is consumed it can immediately delete it, keeping the data size small.

What is perhaps not obvious, is that getting the broker and consumer to come into agreement about what has been consumed is not a trivial problem. If the broker records a message as **consumed** immediately every time it is handed out over the network, then if the consumer fails to process the message (say because it crashes or the request times out or whatever) that message will be lost. To solve this problem, many messaging systems add an acknowledgement feature which means that messages are only marked as **not consumed** when they are sent; the broker waits for a specific acknowledgement from the consumer to record the message as **consumed**. This strategy fixes the problem of losing messages, but creates new problems. First of all, if the consumer processes the message but fails before it can send an acknowledgement then the message will be consumed twice. The second problem is around performance, now the broker must keep multiple states about every single message (first to lock it so it is not given out a second time, and then to mark it as permanently consumed so that it can be removed). Tricky problems must be dealt with, like what to do with messages that are sent but never acknowledged.

Kafka handles this differently. Our topic is divided into a set of totally ordered partitions, each of which is consumed by one consumer at any given time. This means that the position of consumer in each partition is just a single integer, the offset of the next message to consume. This makes the state about what has been consumed very small, just one number for each partition. This state can be periodically checkpointed. This makes the equivalent of message acknowledgements very cheap.

There is a side benefit of this decision. A consumer can deliberately *rewind* back to an old offset and re-consume data. This violates the common contract of a queue, but turns out to be an essential feature for many consumers. For example, if the consumer code has a bug and is discovered after some messages are consumed, the consumer can re-consume those messages once the bug is fixed.

Offline Data Count

Scalable persistence allows for the possibility of consumers that only periodically consume such as batch data loads that periodically bulk-load data into an offline system such as Hadoop or a relational data warehouse.

In the case of Hadoop we parallelize the data load by splitting the load over individual map tasks, one for each node/topic/partition combination, allowing full parallelism in the loading. Hadoop provides the task management, and tasks which fail can restart without danger of duplicate data—they simply restart from their original position.

4.5 Message Delivery Semantics

Now that we understand a little about how producers and consumers work, let's discuss the semantic guarantees Kafka provides between producer and consumer. Clearly there are multiple possible message delivery guarantees that could be provided:

- *At most once*—Messages may be lost but are never redelivered.
- *At least once*—Messages are never lost but may be redelivered.
- *Exactly once*—this is what people actually want, each message is delivered once and only once.

It's worth noting that this breaks down into two problems: the *durability* guarantees for publishing a message and the *guarantees* when consuming a message.

Many systems claim to provide "exactly once" delivery semantics, but it is important to read the fine print, most of these claims are misleading (i.e. they don't translate to the case where consumers or producers can fail, or cases where there are multiple consumer processes, or cases where data written to disk can be lost).

Kafka's semantics are straight forward. When publishing a message we have a notion of the message being "committed" to the log. Once a published message is committed it will not be lost as long as one broker that replicates the partition to which this message was written remains "alive". The definition of alive as well as a description of which types of failures we attempt to handle will be described in more detail in the next section. For now let's assume a perfect, lossless broker and try to understand the guarantees to the producer and consumer. If a producer attempts to publish a message and experiences a network error it cannot be sure if this error happened before or after the message was committed. This is similar to the semantics of inserting into a database table with an autogenerated key.

These are not the strongest possible semantics for publishers. Although we cannot be sure of what happened in the case of a network error, it is possible to allow the producer to generate a sort of "primary key" that makes retrying the produce request idempotent. This feature is not trivial for a replicated system because of course it must work even (or especially) in the case of a server failure. With this feature it would suffice for the producer to retry until it receives acknowledgement of a successfully committed message at which point we would guarantee the message had been published exactly once. We hope to add this in a future Kafka version.

Not all use cases require such strong guarantees. For uses which are latency sensitive we allow the producer to specify the durability level it desires. If the producer specifies that it wants to wait on the message being committed this can take on the order of 10 ms. However the producer can also specify that it wants to perform the send completely asynchronously or that it wants to wait only until the leader (but not necessarily the followers) have the message.

Now let's describe the semantics from the point-of-view of the consumer. All replicas have the exact same log with the same offsets. The consumer controls its position in this log. If the consumer never crashed it could just store this position in memory, but if the consumer fails and we want this topic partition to be taken over by another process the new process will need to choose an appropriate position from which to start processing. Let's say the consumer reads some messages - it has several options for processing the messages and updating its position.

1. It can read the messages, then save its position in the log, and finally process the messages. In this case there is a possibility that the consumer process crashes after saving its position but before saving the output of its message processing. In this case the process that took over processing would start at the saved position even though a few messages prior to that position had not been processed. This corresponds to "at-most-once" semantics as in the case of a consumer failure messages may not be processed.
2. It can read the messages, process the messages, and finally save its position. In this case there is a possibility that the consumer process crashes after processing messages but before saving its position. In this case when the new process takes over the first few messages it receives will already have been processed. This corresponds to the "at-least-once" semantics in the case of consumer failure. In many cases messages have a primary key and so the updates are idempotent (receiving the same message twice just overwrites a record with another copy of itself).
3. So what about exactly once semantics (i.e. the thing you actually want)? The limitation here is not actually a feature of the messaging system but rather the need to co-ordinate the consumer's position with what is actually stored as output. The classic way of achieving this would be to introduce a two-phase commit between the storage for the consumer position and the storage of the consumers output. But this can be handled more simply and generally by simply letting the consumer store its offset in the same place as its output. This is better because many of the output systems a consumer might want to write to will not support a two-phase commit. As an example of this, our Hadoop ETl that populates data in HDFS stores its offsets in HDFS with the data it reads so that it is guaranteed that either data and offsets are both updated or neither is. We follow similar patterns for many other data systems which require these stronger semantics and for which the messages do not have a primary key to allow

for deduplication.

So effectively Kafka guarantees at-least-once delivery by default and allows the user to implement at most once delivery by disabling retries on the producer and committing its offset prior to processing a batch of messages. Exactly-once delivery requires co-operation with the destination storage system but Kafka provides the offset which makes implementing this straight forward.

4.7.7 Replication

Kafka replicates the log for each topic's partitions across a configurable number of servers (you can set this replication factor on a topic-by-topic basis). This allows automatic failover to these replicas when a server in the cluster fails so messages remain available in the presence of failures.

Other messaging systems provide some replication related features, but, in our (totally biased) opinion, this appears to be a tacked on thing, not heavily used, and with large downsides: slaves are inactive, throughput is heavily impacted, it requires fiddly manual configuration, etc. Kafka is meant to be used with replication by default—in fact we implement un-replicated topics as replicated topics where the replication factor is one.

The unit of replication is the topic partition. Under non-failure conditions, each partition in Kafka has a single leader and zero or more followers. The total number of replicas including the leader constitute the replication factor. All reads and writes go to the leader of the partition. Typically, there are many more partitions than brokers and the leaders are evenly distributed among brokers. The logs on the followers are identical to the leader's log—all have the same offsets and messages in the same order (though, of course, at any given time the leader may have a few as-yet-unreplicated messages at the end of its log).

Followers consume messages from the leader just as a normal Kafka consumer would and apply them to their own log. Having the followers pull from the leader has the nice property of allowing the follower to naturally batch together log entries they are applying to their log.

As with most distributed systems automatically handling failures requires having a precise definition of what it means for a node to be "alive". For Kafka node liveness has two conditions

- 1. A node must be able to maintain its session with ZooKeeper (via ZooKeeper's heartbeat mechanism)
- 2. If it is a slave it must replicate the writes happening on the leader and not fall "too far" behind

We refer to nodes satisfying these two conditions as being "in sync" to avoid the vagueness of "alive" or "failed". The leader keeps track of the set of "in sync" nodes. If a follower dies, gets stuck, or falls behind, the leader will remove it from the list of in sync replicas. The determination of stuck and lagging replicas is controlled by the `replica.lag.time.max.ms` configuration.

In distributed systems terminology we only attempt to handle a "fail/recover" model of failures where nodes suddenly cease working and then later recover (perhaps without knowing that they have died). Kafka does not handle so-called "Byzantine" failures in which nodes produce arbitrary or malicious responses (perhaps due to bugs or foul play).

A message is considered "committed" when all in sync replicas for that partition have applied it to their log. Only committed messages are ever given out to the consumer. This means that the consumer need not worry about potentially seeing a message that could be lost if the leader fails. Producers, on the other hand, have the option of either waiting for the message to be committed or not, depending on their preference for tradeoff between latency and durability. This preference is controlled by the `request.required.acks` setting that the producer uses.

The guarantee that Kafka offers is that a committed message will not be lost, as long as there is at least one in sync replica alive, at all times.

Kafka will remain available in the presence of node failures after a short fail-over period, but may not remain available in the presence of network partitions.

`ReplicatedLoggs: Quorums, ISRs, and State Machines (Dhrmaji)`

At its heart a Kafka partition is a replicated log. The replicated log is one of the most basic primitives in distributed data systems, and there are many approaches for implementing one. A replicated log can be used by other systems as a primitive for implementing other distributed systems in the `state-machine` style.

A replicated log models the process of coming into consensus on the order of a series of values (generally numbering the log entries 0, 1, 2, ...). There are many ways to implement this, but the simplest and fastest is with a leader who chooses the ordering of values provided to it. As long as the leader remains alive, all followers need to only copy the values and ordering the leader chooses.

Of course if leaders didn't fail we wouldn't need followers! When the leader does die we need to choose a new leader from among the followers. But followers themselves may fall behind or crash so we must ensure we choose an up-to-date follower. The fundamental guarantee a log replication algorithm must provide is that if we tell the client a message is committed, and the leader fails, the new leader we elect must also have that message. This yields a tradeoff: if the leader waits for more followers to acknowledge a message before declaring it committed then there will be more potentially electable leaders.

If you choose the number of acknowledgements required and the number of logs that must be compared to elect a leader such that there is guaranteed to be an overlap, then this is called a Quorum.

A common approach to this tradeoff is to use a majority vote for both the commit decision and the leader election. This is not what Kafka does, but let's explore it anyway to understand the tradeoffs. Let's say we have $2f+1$ replicas. If $f+1$ replicas must receive a message prior to a commit being declared by the leader, and if we elect a new leader by electing the follower with the most complete log from at least $f+1$ replicas, then, with no more than f failures, the leader is guaranteed to have all committed messages. This is because among any $f+1$ replicas, there must be at least one replica that contains all committed messages. That replica's log will be the most complete and therefore will be selected as the new leader. There are many remaining details that each algorithm must handle (such as precisely defined what makes a log more complete, ensuring log consistency during leader failure or changing the set of servers in the replica set) but we will ignore these for now.

This majority vote approach has a very nice property: the latency is dependent on only the fastest servers. That is, if the replication factor is three, the latency is determined by the faster slave not the slower one.

There are a rich variety of algorithms in this family including ZooKeeper's `2Safe`, `Haft`, and `Witness` (see `Replication`). The most similar academic publication we are aware of to Kafka's actual implementation is `Rea` (that's from Microsoft).

The downside of majority vote is that it doesn't take many failures to leave you with no electable leaders. To tolerate one failure requires three copies of the data, and to tolerate two failures requires five copies of the data. In our experience having only enough redundancy to tolerate a single failure is not enough for a practical system, but doing every write five times, with 5x the disk space requirements and 1/5th the throughput, is not very practical for large volume data problems. This is likely why quorum algorithms more commonly appear for shared cluster configuration such as ZooKeeper but are less common for primary data storage. For example in HDFS the namenode's high-availability feature is built on a `majority-vote-no-leader` (journal) but this more expensive approach is not used for the data itself.

Kafka takes a slightly different approach to choosing its quorum set. Instead of majority vote, Kafka dynamically maintains a set of in-sync replicas (ISR) that are caught-up to the leader. Only members of this set are eligible for election as leader. A write to a Kafka partition is not considered committed until all in-sync replicas have received the write. This ISR set is persisted to ZooKeeper whenever it changes. Because of this, any replica in the ISR is eligible to be elected leader. This is an important factor for Kafka's usage model where there are many partitions and ensuring leadership balance is important. With this ISR model and $f+1$ replicas, a Kafka topic can tolerate f failures without losing committed messages.

For most use cases we hope to handle, we think this tradeoff is a reasonable one. In practice, to tolerate f failures, both the majority vote and the ISR approach will wait for the same number of replicas to acknowledge before committing a message (e.g. to survive one failure a majority quorum needs three replicas and one acknowledgement and the ISR approach requires two replicas and one acknowledgement). The ability to commit without the slowest servers is an advantage of the majority vote approach. However, we think it is ameliorated by allowing the client to choose whether they block on the message commit or not, and the additional throughput and disk space due to the lower required replication factor is worth it.

Another important design distinction is that Kafka does not require that crashed nodes recover with all their data intact. It is not uncommon for replication algorithms in this space to depend on the existence of "stable storage" that cannot be lost in any failure-recovery scenario without potential consistency violations. There are two primary problems with this assumption. First, disk errors are the most common problem we observe in real operation of persistent data systems and they often do not leave data intact. Secondly, even if this were not a problem, we do not want to require the use of fsync on every write for our consistency guarantees as this can reduce performance by two to three orders of magnitude. Our protocol for allowing a replica to rejoin the ISR ensures that before rejoining, it must fully re-sync again even if it lost unflushed data in its crash.

`Unordered leader election: What if they all die?`

Note that Kafka's guarantee with respect to data loss is predicated on at least one replica remaining in sync. If all the nodes replicating a partition die, this guarantee no longer holds.

However a practical system needs to do something reasonable when all the replicas die. If you are unlucky enough to have this occur, it is important to consider what will happen. There are two behaviors that could be

- implemented:
1. Wait for a replica in the ISR to come back to life and choose this replica as the leader (hopefully it still has all its data).
 2. Choose the first replica (not necessarily in the ISR) that comes back to life as the leader.

This is a simple tradeoff between availability and consistency: if we wait for replicas in the ISR, then we will remain unavailable as long as those replicas are down. If such replicas were destroyed or their data was lost, then we are permanently down. If, on the other hand, a non-in-sync replica comes back to life and we allow it to become leader, then its log becomes the source of truth even though it is not guaranteed to have every committed message. In our current release we choose the second strategy and favor choosing a potentially inconsistent replica when all replicas in the ISR are dead. In the future, we would like to make this configurable to better support use cases where downtime is preferable to inconsistency.

This dilemma is not specific to Kafka. It exists in any quorum-based scheme. For example in a majority voting scheme, if a majority of servers suffer a permanent failure, then you must either choose to lose 100% of your data or violate consistency by taking what remains on an existing server as your new source of truth.

Availability and Durability Guarantees

When writing to Kafka, producers can choose whether they wait for the message to be acknowledged by 0,1 or all (-1) replicas. Note that "acknowledgement by all replicas" does not guarantee that the full set of assigned replicas have received the message. By default, when request.required.acks=-1, acknowledgment happens as soon as all the current in-sync replicas have received the message. For example, if a topic is configured with only two replicas and one fails (i.e. only one in-sync replica remains), then writes that specify request.required.acks=-1 will succeed. However, these writes could be lost if the remaining replica also fails. Although this ensures maximum availability of the partition, this behavior may be undesirable to some users who prefer durability over availability. Therefore, we provide two topic-level configurations that can be used to prefer message durability over availability:

1. Disable unclean leader election - if all replicas become unavailable, then the partition will remain unavailable until the most recent leader becomes available again. This effectively prefers unavailability over the risk of message loss. See the previous section on Unclean Leader Election for clarification.
2. Specify a minimum ISR size - the partition will only accept writes if the size of the ISR is above a certain minimum, in order to prevent the loss of messages that were written to just a single replica, which subsequently becomes unavailable. This setting only takes effect if the producer uses request.required.acks=-1 and guarantees that the message will be acknowledged by at least this many in-sync replicas. This setting offers a trade-off between consistency and availability. A higher setting for minimum ISR size guarantees better consistency since the message is guaranteed to be written to more replicas which reduces the probability that it will be lost. However, it reduces availability since the partition will be unavailable for writes if the number of in-sync replicas drops below the minimum threshold.

Replica Management

The above discussion on replicated logs really covers only a single log, i.e. one topic partition. However a Kafka cluster will manage hundreds or thousands of these partitions. We attempt to balance partitions within a cluster in a round-robin fashion to avoid clustering all partitions for high-volume topics on a small number of nodes. Likewise we try to balance leadership so that each node is the leader for a proportional share of its partitions.

It is also important to optimize the leadership election process as that is the critical window of unavailability. A naive implementation of leader election would end up running an election per partition for all partitions a node hosted when that node failed. Instead, we elect one of the brokers as the "controller". This controller detects failures at the broker level and is responsible for changing the leader of all affected partitions in a failed broker. The result is that we are able to batch together many of the required leadership change notifications which makes the election process far cheaper and faster for a large number of partitions. If the controller fails, one of the surviving brokers will become the new controller.

Log Compaction

Log compaction ensures that Kafka will always retain at least the last known value for each message key within the log of data for a single topic partition. It addresses use cases and scenarios such as restoring state after application crashes or system failure, or reloading caches after application restarts during operational maintenance. Let's dive into these use cases in more detail and then describe how compaction works.

So far we have described only the simpler approach to data retention where old log data is discarded after a fixed period of time or when the log reaches some predetermined size. This works well for temporal event data such as logging where each record stands alone. However an important class of data streams are the log of changes to keyed, mutable data (for example, the changes to a database table).

Let's discuss a concrete example of such a stream. Say we have a topic containing user email addresses; every time a user updates their email address we send a message to this topic using their user id as the primary key. Now say we send the following messages over some time period for a user with id 123, each message corresponding to a change in email address (messages for other ids are omitted):

```
123 => bill@microsoft.com
      .
      .
123 => bill@gatesfoundation.org
      .
      .
123 => bill@gmail.com
```

Log compaction gives us a more granular retention mechanism so that we are guaranteed to retain at least the last update for each primary key (e.g. bill@gmail.com). By doing this we guarantee that the log contains a full snapshot of the final value for every key not just keys that changed recently. This means downstream consumers can restore their own state off this topic without us having to retain a complete log of all changes.

Let's start by looking at a few use cases where this is useful, then we'll see how it can be used.

1. Database change subscription. It is often necessary to have a data set in multiple data systems, and often one of these systems is a database of some kind (either a RDBMS or perhaps a new-fangled key-value store). For example you might have a database, a cache, a search cluster, and a Hadoop cluster. Each change to the database will need to be reflected in the cache, the search cluster, and eventually in Hadoop. In the case that one is only handling the real-time updates you only need recent log. But if you want to be able to reload the cache or restore a failed search node you may need a complete data set.
2. Event sourcing. This is a style of application design which co-locates query processing with application design and uses a log of changes as the primary store for the application.
3. Journaling for high-availability. A process that does local computation can be made fault-tolerant by logging out changes that it makes to its local state so another process can reload these changes and carry on if it should fail. A concrete example of this is handling counts, aggregations, and other "group by"-like processing in a stream query system. Samza, a real-time stream-processing framework, uses this feature for exactly this purpose.

In each of these cases one needs primarily to handle the real-time feed of changes, but occasionally, when a machine crashes or data needs to be re-loaded or re-processed, one needs to do a full load. Log compaction allows feeding both of these use cases off the same backing topic. This style of usage of a log is described in more detail in this blog post:

The general idea is quite simple. If we had infinite log retention, and we logged each change in the above cases, then we would have captured the state of the system at each time from when it first began. Using this complete log we could restore to any point in time by replaying the first N records in the log. This hypothetical complete log is not very practical for systems that update a single record many times as the log will grow without bound even for a stable dataset. The simple log retention mechanism which throws away old updates will bound space but the log is no longer a way to restore the current state--now restoring from the beginning of the log no longer recreates the current state as old updates may not be captured at all.

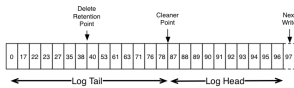
Log compaction is a mechanism to give finer-grained per-record retention, rather than the coarser-grained time-based retention. The idea is to selectively remove records where we have a more recent update with the same primary key. This way the log is guaranteed to have at least the last state for each key.

This retention policy can be set per-topic, so a single cluster can have some topics where retention is enforced by size or time and other topics where retention is enforced by compaction.

This functionality is inspired by one of LinkedIn's oldest and most successful pieces of infrastructure—a database changelog caching service called Databus. Unlike most log-structured storage systems Kafka is built for subscription and organizes data for fast linear reads and writes. Unlike Databus, Kafka acts as a source-of-truth store so it is useful even in situations where the upstream data source would not otherwise be replayable.

Log Compaction Basics

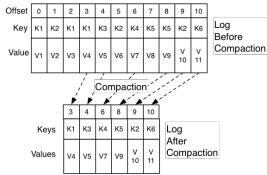
Here is a high-level picture that shows the logical structure of a Kafka log with the offset for each message.



The head of the log is identical to a traditional Kafka log. It has dense, sequential offsets and retains all messages. Log compaction adds an option for handling the tail of the log. The picture above shows a log with a compacted tail. Note that the messages in the tail of the log retain the original offset assigned when they were first written—that never changes. Note also that all offsets remain valid positions in the log, even if the message with that offset has been compacted away; in this case this position is indistinguishable from the next highest offset that does appear in the log. For example, in the picture above the offsets 36, 37, and 38 are all equivalent positions and a read beginning at any of these offsets would return a message set beginning with 38.

Compaction also allows for deletes. A message with a key and a null payload will be treated as a delete from the log. This delete marker will cause any prior message with that key to be removed (as would any new message with that key), but delete markers are special in that they will themselves be cleaned out of the log after a period of time to free up space. The point in time at which deletes are no longer retained is marked as the "delete retention point" in the above diagram.

The compaction is done in the background by periodically recycling log segments. Cleaning does not block reads and can be throttled to use no more than a configurable amount of I/O throughput to avoid impacting producers and consumers. The actual process of compacting a log segment looks something like this:



What guarantees does log compaction provide?

Log compaction guarantees the following:

- 1. Any consumer that stays caught up to within the head of the log will see every message that is written; those messages will have sequential offsets.
- 2. Ordering of messages is always maintained. Compaction will never re-order messages, just remove some.
- 3. The offset for a message never changes. It is the permanent identifier for a position in the log.
- 4. Any read progressing from offset 0 will see at least the final state of all records in the order they were written. All delete markers for deleted records will be seen provided the reader reaches the head of the log in a time period less than the topic's `delete.retention.ms` setting (the default is 24 hours). This is important as delete marker removal happens concurrently with read (and thus it is important that we not remove any delete marker prior to the reader seeing it).
- 5. Any consumer progressing from the start of the log will see at least the final state of all records in the order they were written. All delete markers for deleted records will be seen provided the consumer reaches the head of the log in a time period less than the topic's `delete.retention.ms` setting (the default is 24 hours). This is important as delete marker removal happens concurrently with read, and thus it is important that we do not remove any delete marker prior to the consumer seeing it.

Log Compaction Details

Log compaction is handled by the log cleaner, a pool of background threads that recopy log segment files, removing records whose key appears in the head of the log. Each compactor thread works as follows:

- 1. It chooses the log that has the highest ratio of log head to log tail
- 2. It creates a succinct summary of the last offset for each key in the head of the log
- 3. It recopies the log from beginning to end removing keys which have a later occurrence in the log. New, clean segments are swapped into the log immediately so the additional disk space required is just one additional log segment (not a fully copy of the log).
- 4. The summary of the log head is essentially just a space-compact hash table. It uses exactly 24 bytes per entry. As a result with 8GB of cleaner buffer one cleaner iteration can clean around 366GB of log head (assuming 1k messages).

Configuring The Log Cleaner

The log cleaner is disabled by default. To enable it set the server config

```
log.cleaner.enable=true
```

This will start the pool of cleaner threads. To enable log cleaning on a particular topic you can add the log-specific property

```
log.cleanup.policy=compact
```

This can be done either at topic creation time or using the alter topic command.

Further cleaner configurations are described [here](#):

Log Compaction Limitations

- 1. You cannot configure yet how much log is retained without compaction (the "head" of the log). Currently all segments are eligible except for the last segment, i.e. the one currently being written to.
- 2. Log compaction is not yet compatible with compressed topics.

4.9 Quotas

Starting in 0.9, the Kafka cluster has the ability to enforce quotas on produce and fetch requests. Quotas are basically byte-rate thresholds defined per client id. A client id logically identifies an application making a request. Hence a single client id can span multiple producer and consumer instances and the quota will apply for all of them as a single entity i.e. if client id="test-client" has a produce quota of 10MB/sec, this is shared across all instances with that same id.

Why are quotas necessary?

It is possible for producers and consumers to produce/consume very high volumes of data and thus monopolize broker resources, cause network saturation and generally DOS other clients and the brokers themselves. Having quotas protects against these issues and is all the more important in large multi-tenant clusters where a small set of badly behaved clients can degrade user experience for the well behaved ones. In fact, when running Kafka as a service this even makes it possible to enforce API limits according to an agreed upon contract.

Enforcement

By default, each unique client id receives a fixed quota in bytes/sec as configured by the cluster (`quota.producer.default`, `quota.consumer.default`). This quota is defined on a per-broker basis. Each client can publish/fetch a maximum of `K` bytes/sec per broker before it gets throttled. We decided that defining these quotas per broker is much better than having a fixed cluster-wide bandwidth per client because that would require a mechanism to share client quota usage among all the brokers. This can be harder to get right than the quota implementation itself.

How does a broker react when it detects a quota violation? In our solution, the broker does not return an error rather it attempts to slow down a client exceeding its quota. It computes the amount of delay needed to bring a guilty client under (it's quota and delays the response for that time. This approach keeps the quota violation transparent to clients (outside of client-side metrics). This also keeps them from having to implement any special backoff and retry behavior which can get tricky. In fact, bad client behavior (retry without backoff) can exacerbate the very problem quotas are trying to solve.

Client byte rate is measured over multiple small windows (for e.g. 30 windows of 1 second each) in order to detect and correct quota violations quickly. Typically, having large measurement windows (for e.g. 10 windows of 30 seconds each) leads to large bursts of traffic followed by long delays which is not great in terms of user experience.

@author:erikgr

It is possible to override the default quota for client ids that need a higher (or even lower) quota. The mechanism is similar to the per-topic log config overrides. Client-id overrides are written to ZooKeeper under `/configs/clients`. These overrides are read by all brokers and are effective immediately. This lets us change quotas without having to do a rolling restart of the entire cluster. See [here](#) for details.

5. Implementation

5.1.1 API Design

Producer APIs

The Producer API that wraps the 2 low-level producers - `kafka.producer.SyncProducer` and `kafka.producer.async.AsyncProducer`.

```
class Producer {  
    /*  
     * Sends the data, partitioned by key to the topic using either the */  
     * synchronous or the asynchronous producer */  
    public void send(kafka.javaapi.producer.ProducerData<K,V> producerData);  
    /*  
     * Sends a list of data, partitioned by key to the topic using either */  
     * the synchronous or the asynchronous producer */  
    public void send(java.util.List<kafka.javaapi.producer.ProducerData<K,V>> producerData);  
    /*  
     * Closes the producer and cleans up */  
    public void close();  
}
```

The goal is to expose all the producer functionality through a single API to the client. The new producer -

- can handle queuing/buffering of multiple producer requests and asynchronous dispatch of the batched data -

`kafka.producer.Producer` provides the ability to batch multiple produce requests (`@producer.type=async`), before serializing and dispatching them to the appropriate kafka broker partition. The size of the batch can be controlled by a few config parameters. As events enter a queue, they are buffered in a queue, until either `queue.time` or `batch.size` is reached. A background thread (`kafka.producer.async.ProducerSendThread`) dequeues the batch of data and lets the `kafka.producer.EventHandler` serialize and send the data to the appropriate kafka broker partition. A custom event handler can be plugged in through the `event.handler` config parameter. At various stages of this producer queue pipeline, it is helpful to be able to inject callbacks, either for plugging in custom logging/tracing code or custom monitoring logic. This is possible by implementing the `kafka.producer.async.CallbackHandler` interface and setting `callback.handler` config parameter to that class.

- handles the serialization of data through a user-specified `Encoder` -

```
interface Encoder<T> {  
    public Message toMessage(T data);  
}
```

The default is the no-op `kafka.serializer.DefaultEncoder`

- provides software load balancing through an optionally user-specified `Partitioner` -

The routing decision is influenced by the `kafka.producer.Partitioner`.

```
interface Partitioner<T> {  
    int partition(T key, int numPartitions);  
}
```

The partition API uses the key and the number of available broker partitions to return a partition id. This id is used as an index into a sorted list of broker_ids and partitions to pick a broker partition for the producer request. The default partitioning strategy is `hash(key)%numPartitions`. If the key is null, then a random broker partition is picked. A custom partitioning strategy can also be plugged in using the `partitioner.class` config parameter.

Consumer APIs

We have 2 levels of consumer APIs. The low-level "simple" API maintains a connection to a single broker and has a close correspondence to the network requests sent to the server. This API is completely stateless, with the offset being passed in on every request, allowing the user to maintain this metadata however they choose.

The high-level API hides the details of brokers from the consumer and allows consuming off the cluster of machines without concern for the underlying topology. It also maintains the state of what has been consumed. The high-level API also provides the ability to subscribe to topics that match a filter expression (i.e., either a `whitelist` or a `blacklist` regular expression).

Low-level API

```
class SimpleConsumer {  
    /*  
     * Send fetch request to a broker and get back a set of messages. */  
    public ByteBufferMessageSet fetch(FetchRequest request);  
    /*  
     * Send a list of fetch requests to a broker and get back a response set. */  
    public MultiFetchResponse multiFetch(List<FetchRequest> fetches);  
    /**  
     * Set a list of valid offsets (up to maxTime) before the given time.  
     * The result is a list of offsets, in descending order.  
     * @param time: time in millisecs,  
     *             if set to OffsetRequests.MODULES.LATEST_TIME(), get from the latest offset available.  
     *             if set to OffsetRequests.MODULES.EARLIEST_TIME(), get from the earliest offset available.  
     */  
    public long[] getOffsetsBefore(String topic, int partition, long time, int maxNumOffsets);  
}
```

The low-level API is used to implement the high-level API as well as being used directly for some of our offline consumers which have particular requirements around maintaining state.

High-level API

```
/*  
 * create a connection to the cluster */  
ConsumerConnector connector = Consumer.create(consumerConfig);  
  
interface ConsumerConnector {  
    /**  
     * This method is used to get a list of KafkaStreams, which are iterators over  
     * MessageAndMetadata objects from which you can obtain messages and their  
     * associated metadata (currently only topic).  
     * Input: a map of <topic, <streams>  
     * Output: a Map of <topic, list of Message streams>  
     */  
    public Map<String, List<KafkaStream>> createMessageStreams(Map<String, Int> topicCountMap);  
    /**  
     * You can also obtain a list of KafkaStreams, that iterate over messages  
     * from topics that match a TopicFilter. (A TopicFilter encapsulates a  
     * whitelist or a blacklist which is a standard Java regex.)  
     */  
    public List<KafkaStream> createMessageStreamsByFilter(  
        TopicFilter topicFilter, int numStreams);  
    /* Commit the offsets of all messages consumed so far. */  
    public commitOffsets();  
    /* Shut down the connector */  
    public shutdown();  
}
```

This API is centered around iterators, implemented by the `KafkaStream` class. Each `KafkaStream` represents the stream of messages from one or more partitions on one or more servers. Each stream is used for single threaded processing, so the client can provide the number of desired streams in the create call. Thus a stream may represent the merging of multiple server partitions (to correspond to the number of processing threads), but each partition only goes to one stream.

The `createMessageStreams` call registers the consumer for the topic, which results in rebalancing the consumer/broker assignment. The API encourages creating many topic streams in a single call in order to minimize this rebalancing. The `createMessageStreamsByFilter` call (additionally) registers watchers to discover new topics that match its filter. Note that each stream that `createMessageStreamsByFilter` returns may iterate over messages from multiple topics (i.e., if multiple topics are allowed by the filter).

5.2 Network Layer

The network layer is a fairly straight-forward NIO server, and will not be described in great detail. The sendfile implementation is done by giving the `MessageSet` interface a `writeTo` method. This allows the file-backed message set to use the more efficient `transferTo` implementation instead of an in-process buffered write. The threading model is a single acceptor thread and `N` processor threads which handle a fixed number of connections each. This design has been pretty thoroughly tested [at Lawrence](#) and found to be simple to implement and fast. The protocol is kept quite simple to allow for future implementation of clients in other languages.

5.3 Messages

Messages consist of a fixed size header and variable length opaque byte array payload. The header contains a format version and a CRC32 checksum to detect corruption or truncation. Leaving the payload opaque is the right decision: there is a great deal of progress being made on serialization libraries right now, and any particular choice is unlikely to be right for all uses. Needless to say a particular application using Kafka would likely mandate a particular serialization type as part of its usage. The `MessageSet` interface is simply an iterator over messages with specialized methods for bulk reading and writing to an NIO Channel.

5.4 Message Format

```
/**
 * A message. The format of an N byte message is the following:
 *
 * * If magic byte is 0
 * *
 * * 1. 1 byte "magic" identifier to allow format changes
 * *
 * * 2. 4 byte CRC32 of the payload
 * *
 * * 3. N - 5 byte payload
 * *
 * * If magic byte is 1
 * *
 * * 1. 1 byte "magic" identifier to allow format changes
 * *
 * * 2. 1 byte "attributes" identifier to allow annotations on the message independent of the version (e.g. compression enabled, type of codec used)
 * *
 * * 3. 4 byte CRC32 of the payload
 * *
 * * 4. N - 6 byte payload
 * *
 */
```

5.5 Log

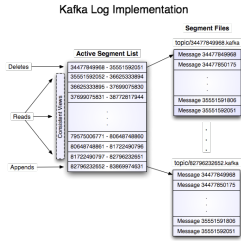
A log for a topic named "my_topic" with two partitions consists of two directories (namely `my_topic_0` and `my_topic_1`) populated with data files containing the messages for that topic. The format of the log files is a sequence of "log entries"; each log entry is a 4 byte integer *N* storing the message length which is followed by the *N* message bytes. Each message is uniquely identified by a 64-bit integer offset giving the byte position of the start of this message in the stream of all messages ever sent to that topic on that partition. The on-disk format of each message is given below. Each log file is named with the offset of the first message it contains. So the first file created will be 00000000000.kafka, and each additional file will have an integer name roughly 5 bytes from the previous file where 5 is the max log file size given in the configuration.

The exact binary format for messages is versioned and maintained as a standard interface so message sets can be transferred between producer, broker, and client without recopying or conversion when desirable. This format is as follows:

On-disk format of a message

```
message length : 4 bytes (value: 1-4*n)
"magic" value : 1 byte
crc32 : 4 bytes
payload : n bytes
```

The use of the message offset as the message id is unusual. Our original idea was to use a GUID generated by the producer, and maintain a mapping from GUID to offset on each broker. But since a consumer must maintain an ID for each server, the global uniqueness of the GUID provides no value. Furthermore the complexity of maintaining the mapping from a random id to an offset requires a heavy weight index structure which must be synchronized with disk, essentially requiring a full persistent random-access data structure. Thus to simplify the lookup structure we decided to use a simple per-partition atomic counter which could be coupled with the partition id and node id to uniquely identify a message; this makes the lookup structure simpler, though multiple seeks per consumer request are still likely. However once we settled on a counter, the jump to directly using the offset seemed natural - both after all are monotonically increasing integers unique to a partition. Since the offset is hidden from the consumer API this decision is ultimately an implementation detail and we went with the more efficient approach.



Writes:

The log allows serial appends which always go to the last file. This file is rolled over to a fresh file when it reaches a configurable size (say 1GB). The log takes two configuration parameter *N* which gives the number of messages to write before forcing the OS to flush the file to disk, and *S* which gives a number of seconds after which a flush is forced. This gives a durability guarantee of losing at most *N* messages or *S* seconds of data in the event of a system crash.

Reads:

Reads are done by giving the 64-bit logical offset of a message and an *S*-byte max chunk size. This will return an iterator over the messages contained in the *S*-byte buffer. *S* is intended to be larger than any single message, but in the event of an abnormally large message, the read can be retried multiple times, each time doubling the buffer size, until the message is read successfully. A maximum message and buffer size can be specified to make the server reject messages larger than some size, and to give a bound to the client on the maximum it need ever read to get a complete message. It is likely that the read buffer ends with a partial message, this is easily detected by the size delimiting.

The actual process of reading from an offset requires first locating the log segment file in which the data is stored, calculating the file-specific offset from the global offset value, and then reading from that file offset. The search is done as a simple binary search variation against an in-memory range maintained for each file.

The log provides the capability of getting the most recently written message to allow clients to start subscribing as of "right now". This is also useful in the case the consumer fails to consume its data within its SLA-specified number of days. In this case when the client attempts to consume a non-existent offset it is given an `OutOfRangeException` and can either reset itself or fail as appropriate to the use case.

The following is the format of the results sent to the consumer.

```
MessageSetSend (fetch result)

total length : 4 bytes
error code : 2 bytes
message 1 : x bytes
...
message n : x bytes

MultiMessageSetSend (multiFetch result)

total length : 4 bytes
error code : 2 bytes
MessageSetSend 1
...
MessageSetSend n
```

Deletes:

Data is deleted one log segment at a time. The log manager allows pluggable delete policies to choose which files are eligible for deletion. The current policy deletes any log with a modification time of more than *N* days ago, though a policy which retained the last *N* GB could also be useful. To avoid locking reads while still allowing deletes that modify the segment list we use a copy-on-write style segment list implementation that provides consistent views to allow a binary search to proceed on an immutable static snapshot view of the log segments while deletes are progressing.

Guarantees:

The log provides a configuration parameter *N* which controls the maximum number of messages that are written before forcing a flush to disk. On startup a log recovery process is run that iterates over all messages in the newest log segment and verifies that each message entry is valid. A message entry is valid if the sum of its size and offset are less than the length of the file AND the CRC32 of the message payload matches the CRC stored with the message. In the event corruption is detected the log is truncated to the last valid offset.

Note that two kinds of corruption must be handled: truncation in which an unwritten block is lost due to a crash, and corruption in which a nonsense block is ADDED to the file. The reason for this is that in general the OS makes no guarantee of the write order between the file inode and the actual block data so in addition to losing written data the file can gain nonsense data if the inode is updated with a new size but a crash occurs before the block containing that data is not written. The CRC detects this corner case, and prevents it from corrupting the log (though the unwritten messages are, of course, lost).

5.6 Distribution

ConsumerOffsetTracking

The high-level consumer tracks the maximum offset it has consumed in each partition and periodically commits its offset vector so that it can resume from those offsets in the event of a restart. Kafka provides the option to store all the offsets for a given consumer group in a designated broker (for that group) called the *offset manager*, i.e., any consumer instance in that consumer group should send its offset commits and fetches to that offset manager (broker). The high-level consumer handles this automatically; if you use the simple consumer you will need to manage offsets manually. This is currently unsupported in the Java simple consumer which can only commit or fetch offsets in ZooKeeper. If you use the Scala simple consumer you can discover the offset manager and explicitly commit or fetch offsets to the offset manager. A consumer can look up its offset manager by issuing a `GroupCoordinatorRequest` to any Kafka broker and reading the `GroupCoordinatorResponse` which will contain the offset manager. The consumer can then proceed to commit or fetch offsets from the offsets manager broker. In case the offset manager moves, the consumer will need to rediscover the offset manager: if you wish to manage your offsets manually, you can take a look at these [code samples](#) that explain how to issue `OffsetCommitRequest` and `OffsetFetchRequest`.

When the offset manager receives an `OffsetCommitRequest`, it appends the request to a special `compacted` Kafka topic named `__consumer_offsets`. The offset manager sends a successful offset commit response to the consumer only after all the replicas of the offsets topic receive the offsets. In case the offsets fail to replicate within a configurable timeout, the offset commit will fail and the consumer may retry the commit after backing off. (This is done automatically by the high-level consumer.) The brokers periodically compact the offsets topic since it only needs to maintain the most recent offset commit per partition. The offset manager also caches the offsets in an in-memory table in order to serve offset fetches quickly.

When the offset manager receives an offset fetch request, it simply returns the last committed offset vector from the offsets cache. In case the offset manager was just started or if it just became the offset manager for a new set of consumer groups (by becoming a leader for a partition of the offsets topic), it may need to load the offsets topic partition into the cache. In this case, the offset fetch will fail with an `OffsetLoadInProgress` exception and the consumer may retry the `OffsetFetchRequest` after backing off. (This is done automatically by the high-level consumer.)

Migrating Offsets from ZooKeeper to Kafka

Kafka consumers in earlier releases store their offsets by default in ZooKeeper. It is possible to migrate these consumers to commit offsets into Kafka by following these steps:

- 1. Set `offsets.storage=kafka` and `dual.commit.enabled=true` in your consumer config.
- 2. Do a rolling bounce of your consumers and then verify that your consumers are healthy.
- 3. Set `dual.commit.enabled=false` in your consumer config.
- 4. Do a rolling bounce of your consumers and then verify that your consumers are healthy.

A roll-back (i.e., migrating from Kafka back to ZooKeeper) can also be performed using the above steps if you set `offsets.storage=zookeeper`.

ZooKeeper Directories

The following gives the ZooKeeper structures and algorithms used for co-ordination between consumers and brokers.

Notation

When an element in a path is denoted `[xyz]`, that means that the value of `xyz` is not fixed and there is in fact a ZooKeeper `znode` for each possible value of `xyz`. For example `/topics/[topic]` would be a directory named `/topics` containing a sub-directory for each topic name. Numerical ranges are also given such as `[0...5]` to indicate the subdirectories 0, 1, 2, 3, 4. An arrow \rightarrow is used to indicate the contents of a `znode`. For example `/hello \rightarrow world` would indicate a `znode` `/hello` containing the value "world".

Broker Node Registry

`/brokers/ids/[0...N] \rightarrow host:port` (ephemeral node)

This is a list of all present broker nodes, each of which provides a unique logical broker id which identifies it to consumers (which must be given as part of its configuration). On startup, a broker node registers itself by creating a `znode` with the logical broker id under `/brokers/ids`. The purpose of the logical broker id is to allow a broker to be moved to a different physical machine without affecting consumers. An attempt to register a broker id that is already in use (say because two servers are configured with the same broker id) is an error.

Since the broker registers itself in ZooKeeper using ephemeral `znodes`, this registration is dynamic and will disappear if the broker is shutdown or dies (thus notifying consumers it is no longer available).

Broker Topic Registry

`/brokers/topics/[topic]/[0...N] \rightarrow nPartitions` (ephemeral node)

Each broker registers itself under the topics it maintains and stores the number of partitions for that topic.

Consumers and Consumer Groups

Consumers of topics also register themselves in ZooKeeper, in order to coordinate with each other and balance the consumption of data. Consumers can also store their offsets in ZooKeeper by setting `offsets.storage=zookeeper`. However, this offset storage mechanism will be deprecated in a future release. Therefore, it is recommended to [migrate offsets to storage in Kafka](#).

Multiple consumers can form a group and jointly consume a single topic. Each consumer in the same group is given a shared `group_id`. For example if one consumer is your foobar process, which is run across three machines, then you might assign this group of consumers the id "foobar". This group id is provided in the configuration of the consumer, and is your way to tell the consumer which group it belongs to.

The consumers in a group divide up the partitions as fairly as possible, each partition is consumed by exactly one consumer in a consumer group.

ConsumerIdRegistry

In addition to the `group_id` which is shared by all consumers in a group, each consumer is given a transient, unique `consumer_id` (of the form `hostname:uid`) for identification purposes. Consumer ids are registered in the following directory.

`/consumers/[group_id]/ids/[consumer_id] \rightarrow {"topic1": #streams, ..., "topicN": #streams}` (ephemeral node)

Each of the consumers in the group registers under its group and creates a `znode` with its `consumer_id`. The value of the `znode` contains a map of `<topic, #streams>`. This id is simply used to identify each of the consumers which is currently active within a group. This is an ephemeral node so it will disappear if the consumer process dies.

ConsumerOffset

Consumers track the maximum offset they have consumed in each partition. This value is stored in a ZooKeeper directory if `offsets.storage=zookeeper`. This value is stored in a ZooKeeper directory.

`/consumers/[group_id]/offsets/[topic]/[broker_id-partition_id] \rightarrow offset_counter_value` (persistent node)

PartitionOwnerRegistry

Each broker partition is consumed by a single consumer within a given consumer group. The consumer must establish its ownership of a given partition before any consumption can begin. To establish its ownership, a consumer writes its own id in an ephemeral node under the particular broker partition it is claiming.

`/consumers/[group_id]/owners/[topic]/[broker_id-partition_id] \rightarrow consumer_node_id` (ephemeral node)

Broker node registration

The broker nodes are basically independent, so they only publish information about what they have. When a broker joins, it registers itself under the broker node registry directory and writes information about its host name and port. The broker also register the list of existing topics and their logical partitions in the broker topic registry. New topics are registered dynamically when they are created on the broker.

Consumer registration algorithm

When a consumer starts, it does the following:

1. Register itself in the consumer id registry under its group.
2. Register a watch on changes (new consumers joining or any existing consumers leaving) under the consumer id registry. (Each change triggers rebalancing among all consumers within the group to which the changed consumer belongs.)
3. Register a watch on changes (new brokers joining or any existing brokers leaving) under the broker id registry. (Each change triggers rebalancing among all consumers in all consumer groups.)
4. If the consumer creates a message stream using a topic filter, it also registers a watch on changes (new topics being added) under the broker topic registry. (Each change will trigger re-evaluation of the available topics to determine which topics are allowed by the topic filter. A new allowed topic will trigger rebalancing among all consumers within the consumer group.)
5. Force itself to rebalance within its consumer group.

Consumer rebalancing algorithm

The consumer rebalancing algorithms allows all the consumers in a group to come into consensus on which consumer is consuming which partitions. Consumer rebalancing is triggered on each addition or removal of both broker nodes and other consumers within the same group. For a given topic and a given consumer group, broker partitions are divided evenly among consumers within the group. A partition is always consumed by a single consumer. This design simplifies the implementation. Had we allowed a partition to be concurrently consumed by multiple consumers, there would be contention on the partition and some kind of locking would be required. If there are more consumers than partitions, some consumers won't get any data at all. During rebalancing, we try to assign partitions to consumers in such a way that reduces the number of broker nodes each consumer has to connect to.

Each consumer does the following during rebalancing:

1. For each topic T that C_i subscribes to
2. let P_T be all partitions producing topic T
3. let C_T be all consumers in the same group as C_i that consume topic T
4. sort P_T (so partitions on the same broker are clustered together)
5. sort C_T
6. let i be the index position of C_i in C_T and let N = size(P_T)/size(C_T)
7. assign partitions from i*N to (i+1)*N - 1 to consumer C_i
8. remove current entries owned by C_i from the partition owner registry
9. add newly assigned partitions to the partition owner registry (we may need to retry this until the original partition owner releases its ownership)

When rebalancing is triggered at one consumer, rebalancing should be triggered in other consumers within the same group about the same time.

6. Operations

Here is some information on actually running Kafka as a production system based on usage and experience at LinkedIn. Please send us any additional tips you know of.

6.1 Basic Kafka Operations

This section will review the most common operations you will perform on your Kafka cluster. All of the tools reviewed in this section are available under the `bin/` directory of the Kafka distribution and each tool will print details on all possible commandline options if it is run with no arguments.

Adding and removing topics

You have the option of either adding topics manually or having them be created automatically when data is first published to a non-existent topic. If topics are auto-created then you may want to tune the default `topic.configurations` used for auto-created topics.

Topics are added and modified using the topic tool:

```
> bin/kafka-topics.sh --zookeeper zk_host:port/chroot --create --topic my_topic_name
--partitions 20 --replication-factor 3 --config x=y
```

The replication factor controls how many servers will replicate each message that is written. If you have a replication factor of 3 then up to 2 servers can fail before you will lose access to your data. We recommend you use a replication factor of 2 or 3 so that you can transparently bounce machines without interrupting data consumption.

The partition count controls how many logs the topic will be sharded into. There are several impacts of the partition count. First each partition must fit entirely on a single server. So if you have 20 partitions the full data set (and read and write load) will be handled by no more than 20 servers (no counting replicas). Finally the partition count impacts the maximum parallelism of your consumers. This is discussed in greater detail in the [concepts section](#).

The configurations added on the command line override the default settings the server has for things like the length of time data should be retained. The complete set of per-topic configurations is documented [here](#).

Modifying topics

You can change the configuration or partitioning of a topic using the same topic tool.

To add partitions you can do

```
> bin/kafka-topics.sh --zookeeper zk_host:port/chroot --alter --topic my_topic_name
--partitions 40
```

Be aware that one use case for partitions is to semantically partition data, and adding partitions doesn't change the partitioning of existing data so this may disturb consumers if they rely on that partition. That is if data is partitioned by `hash(key) % number_of_partitions` then this partitioning will potentially be shuffled by adding partitions but Kafka will not attempt to automatically redistribute data in any way.

To add configs:

```
> bin/kafka-topics.sh --zookeeper zk_host:port/chroot --alter --topic my_topic_name --config x=y
```

To remove a config:

```
> bin/kafka-topics.sh --zookeeper zk_host:port/chroot --alter --topic my_topic_name --delete-config x
```

And finally deleting a topic:

```
> bin/kafka-topics.sh --zookeeper zk_host:port/chroot --delete --topic my_topic_name
```

Topic deletion option is disabled by default. To enable it set the server config

```
delete.topic.enable=true
```

Kafka does not currently support reducing the number of partitions for a topic.

Instructions for changing the replication factor of a topic can be found [here](#)

Graceful shutdown

The Kafka cluster will automatically detect any broker shutdown or failure and elect new leaders for the partitions on that machine. This will occur whether a server fails or it is brought down intentionally for maintenance or configuration changes. For the later cases Kafka supports a more graceful mechanism for stopping a server then just killing it. When a server is stopped gracefully it has two optimizations it will take advantage of:

1. It will sync all its logs to disk to avoid needing to do any log recovery when it restarts (i.e. validating the checksum for all messages in the tail of the log). Log recovery takes time so this speeds up intentional restarts.
2. It will migrate any partitions the server is the leader for to other replicas prior to shutting down. This will make the leadership transfer faster and minimize the time each partition is unavailable to a few milliseconds.

Syncing the logs will happen automatically happen whenever the server is stopped other than by a hard kill, but the controlled leadership migration requires using a special setting:

```
controlled.shutdown.enable=true
```

Note that controlled shutdown will only succeed if all the partitions hosted on the broker have replicas (i.e. the replication factor is greater than 1 and at least one of these replicas is alive). This is generally what you want since shutting down the last replica would make that topic partition unavailable.

Rebalancing leadership

Whenever a broker stops or crashes leadership for that broker's partitions transfers to other replicas. This means that by default when the broker is restarted it will only be a follower for all its partitions, meaning it will not be used for client reads and writes.

To avoid this imbalance, Kafka has a notion of preferred replicas. If the list of replicas for a partition is 1,5,9 then node 1 is preferred as the leader to either node 5 or 9 because it is earlier in the replica list. You can have the Kafka cluster try to restore leadership to the restored replicas by running the command:

```
> bin/kafka-preferred-replica-election.sh --zookeeper zk_host:port/chroot
```

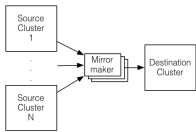
Since running this command can be tedious you can also configure Kafka to do this automatically by setting

the following configuration:

```
auto.leader.rebalance.enable=true
```

Mirroring data between clusters

We refer to the process of replicating data between Kafka clusters "mirroring" to avoid confusion with the replication that happens amongst the nodes in a single cluster. Kafka comes with a tool for mirroring data between Kafka clusters. The tool reads from a source cluster and writes to a destination cluster, like this:



A common use case for this kind of mirroring is to provide a replica in another datacenter. This scenario will be discussed in more detail in the next section.

You can run many such mirroring processes to increase throughput and for fault-tolerance (if one process dies, the others will take over the additional load).

Data will be read from topics in the source cluster and written to a topic with the same name in the destination cluster. In fact the mirror maker is little more than a Kafka consumer and producer hooked together.

The source and destination clusters are completely independent entities; they can have different numbers of partitions and the offsets will not be the same. For this reason the mirror cluster is not really intended as a fault-tolerance mechanism (as the consumer position will be different); for that we recommend using normal in-cluster replication. The mirror maker process will, however, retain and use the message key for partitioning so order is preserved on a per-key basis.

Here is an example showing how to mirror a single topic (named my-topic) from two input clusters:

```
> bin/kafka-run-class.sh kafka.tools.MirrorMaker
--consumer.config consumer-1.properties --consumer.config consumer-2.properties
--producer.config producer.properties --whitelist my-topic
```

Note that we specify the list of topics with the `--whitelist` option. This option allows any regular expression using [Java-style regular expressions](#). So you could mirror two topics named A and B using `--whitelist "A|B"`. Or you could mirror *all* topics using `--whitelist ""`. Make sure to quote any regular expression to ensure the shell doesn't try to expand it as a file path. For convenience we allow the use of `of` instead of `|` to specify a list of topics.

Sometimes it is easier to say what it is that you don't want. Instead of using `--whitelist` to say what you want to mirror you can use `--blacklist` to say what to exclude. This also takes a regular expression argument.

Combining mirroring with the configuration `auto.create.topics.enable=true` makes it possible to have a replica cluster that will automatically create and replicate all data in a source cluster even as new topics are added.

Checking consumer position

Sometimes it's useful to see the position of your consumers. We have a tool that will show the position of all consumers in a consumer group as well as how far behind the end of the log they are. To run this tool on a consumer group named my-group consuming a topic named my-topic would look like this:

```
> bin/kafka-run-class.sh kafka.tools.ConsumerOffsetChecker --zookeeper localhost:2181 --group test
group      Topic      Topic      Elapsed Offset      Logfile      Lag      Owner
my-group   my-topic    0          0                  0          0      test_jkxps-mm-1394154511599-60744496-0
my-group   my-topic    1          0                  0          0      test_jkxps-mm-1394154522217-la0ew913-0
```

Expanding your cluster

Adding servers to a Kafka cluster is easy, just assign them a unique broker id and start up Kafka on your new servers. However these new servers will not automatically be assigned any data partitions, so unless partitions are moved to them they won't be doing any work until new topics are created. So usually when you add machines to your cluster you will want to migrate some existing data to these machines.

The process of migrating data is manually initiated but fully automated. Under the covers what happens is that Kafka will add the new server as a follower of the partition it is migrating and allow it to fully replicate the existing data in that partition. When the new server has fully replicated the contents of this partition and joined the in-sync replica one of the existing replicas will delete their partition's data.

The partition reassignment tool can be used to move partitions across brokers. An ideal partition distribution would ensure even data load and partition sizes across all brokers. The partition reassignment tool does not have the capability to automatically study the data distribution in a Kafka cluster and move partitions around to attain an even load distribution. As such, the admin has to figure out which topics or partitions should be moved around.

The partition reassignment tool can run in 3 mutually exclusive modes -

- `--generate`: In this mode, given a list of topics and a list of brokers, the tool generates a candidate reassignment to move all partitions of the specified topics to the new brokers. This option merely provides a convenient way to generate a partition reassignment plan given a list of topics and target brokers.
- `--execute`: In this mode, the tool kicks off the reassignment of partitions based on the user provided reassignment plan. (using the `--reassignment-json-file` option). This can either be a custom reassignment plan hand crafted by the admin or provided by using the `--generate` option
- `--verify`: In this mode, the tool verifies the status of the reassignment for all partitions listed during the last `--execute`. The status can be either of successfully completed, failed or in progress

Automatically migrating data to new machines

The partition reassignment tool can be used to move some topics off of the current set of brokers to the newly added brokers. This is typically useful while expanding an existing cluster since it is easier to move entire topics to the new set of brokers, than moving one partition at a time. When used to do this, the user should provide a list of topics that should be moved to the new set of brokers and a target list of new brokers. The tool then evenly distributes all partitions for the given list of topics across the new set of brokers. During this move, the replication factor of the topic is kept constant. Effectively the replicas for all partitions for the input list of topics are moved from the old set of brokers to the newly added brokers.

For instance, the following example will move all partitions for topics foo1,foo2 to the new set of brokers 5,6. At the end of this move, all partitions for topics foo1 and foo2 will only exist on brokers 5,6

Since, the tool accepts the input list of topics as a json file, you first need to identify the topics you want to move and create the json file as follows-

```
> cat topics-to-move.json
{"topics": [{"topic": "foo1"},
             {"topic": "foo2"}],
 "version": 1}
```

Once the json file is ready, use the partition reassignment tool to generate a candidate assignment:

```
> bin/kafka-reassign-partitions.sh --zookeeper localhost:2181 --topics-to-move-json-file topics-to-move.json --broker-list "5,6" --generate
Current partition replica assignment
```

```
{
  "version": 1,
  "partitions": [
    {
      "topic": "foo1",
      "partition": 0,
      "replicas": [1,2],
      "topic": "foo1",
      "partition": 1,
      "replicas": [3,4],
      "topic": "foo2",
      "partition": 0,
      "replicas": [1,2],
      "topic": "foo2",
      "partition": 1,
      "replicas": [3,4],
      "topic": "foo1",
      "partition": 1,
      "replicas": [2,3],
      "topic": "foo2",
      "partition": 1,
      "replicas": [2,3]}
  ]
}
```

Proposed partition reassignment configuration

```
{
  "version": 1,
  "partitions": [
    {
      "topic": "foo1",
      "partition": 0,
      "replicas": [5,6],
      "topic": "foo1",
      "partition": 1,
      "replicas": [5,6],
      "topic": "foo2",
      "partition": 0,
      "replicas": [5,6],
      "topic": "foo2",
      "partition": 1,
      "replicas": [5,6],
      "topic": "foo1",
      "partition": 1,
      "replicas": [5,6],
      "topic": "foo2",
      "partition": 1,
      "replicas": [5,6]}
  ]
}
```

The tool generates a candidate assignment that will move all partitions from topics foo1,foo2 to brokers 5,6.

Note, however, that at this point, the partition movement has not started, it merely tells you the current assignment and the proposed new assignment. The current assignment should be saved in case you want to roll back to it. The new assignment should be saved in a json file (e.g. expand-cluster-reassignment.json) to be input to the tool with the `--execute` option as follows-

```
> bin/kafka-reassign-partitions.sh --zookeeper localhost:2181 --reassignment-json-file expand-cluster-reassignment.json --execute
Current partition replica assignment

{"version":1,
 "partitions":[{"topic":"foo1","partition":0,"replicas":[1,2]],
               {"topic":"foo1","partition":0,"replicas":[3,4]],
               {"topic":"foo1","partition":2,"replicas":[1,2]],
               {"topic":"foo2","partition":0,"replicas":[3,4]],
               {"topic":"foo2","partition":1,"replicas":[2,3]],
               {"topic":"foo2","partition":1,"replicas":[2,3]]]
}

Save this to use as the --reassignment-json-file option during rollback
Successfully started reassignment of partitions
{"version":1,
 "partitions":[{"topic":"foo1","partition":2,"replicas":[5,6]],
               {"topic":"foo1","partition":0,"replicas":[5,6]],
               {"topic":"foo1","partition":2,"replicas":[5,6]],
               {"topic":"foo2","partition":0,"replicas":[5,6]],
               {"topic":"foo2","partition":0,"replicas":[5,6]],
               {"topic":"foo2","partition":1,"replicas":[5,6]],
               {"topic":"foo2","partition":1,"replicas":[5,6]]]
}

```

Finally, the `--verify` option can be used with the tool to check the status of the partition reassignment. Note that the same `expand-cluster-reassignment.json` (used with the `--execute` option) should be used with the `--verify` option

```
> bin/kafka-reassign-partitions.sh --zookeeper localhost:2181 --reassignment-json-file expand-cluster-reassignment.json --verify
Status of partition reassignment:
Reassignment of partition [foo1,0] completed successfully
Reassignment of partition [foo1,1] is in progress
Reassignment of partition [foo1,2] is in progress
Reassignment of partition [foo2,0] completed successfully
Reassignment of partition [foo2,1] completed successfully
Reassignment of partition [foo2,2] completed successfully

```

Custom partition reassignment

The partition reassignment tool can also be used to selectively move replicas of a partition to a specific set of brokers. When used in this manner, it is assumed that the user knows the reassignment plan and does not require the tool to generate a candidate reassignment, effectively skipping the `--generate` step and moving straight to the `--execute` step

For instance, the following example moves partition 0 of topic `foo1` to brokers 5,6 and partition 1 of topic `foo2` to brokers 2,3

The first step is to hand craft the custom reassignment plan in a json file:

```
> cat custom-reassignment.json
{"version":1,"partitions":[{"topic":"foo1","partition":0,"replicas":[5,6]},{"topic":"foo2","partition":1,"replicas":[2,3]]}]

```

Then, use the json file with the `--execute` option to start the reassignment process:

```
> bin/kafka-reassign-partitions.sh --zookeeper localhost:2181 --reassignment-json-file custom-reassignment.json --execute
Current partition replica assignment

{"version":1,
 "partitions":[{"topic":"foo1","partition":0,"replicas":[1,2]],
               {"topic":"foo2","partition":1,"replicas":[3,4]]]
}

Save this to use as the --reassignment-json-file option during rollback
Successfully started reassignment of partitions
{"version":1,
 "partitions":[{"topic":"foo1","partition":0,"replicas":[5,6]],
               {"topic":"foo2","partition":1,"replicas":[2,3]]]
}

```

The `--verify` option can be used with the tool to check the status of the partition reassignment. Note that the same `expand-cluster-reassignment.json` (used with the `--execute` option) should be used with the `--verify` option

```
bin/kafka-reassign-partitions.sh --zookeeper localhost:2181 --reassignment-json-file custom-reassignment.json --verify
Status of partition reassignment:
Reassignment of partition [foo1,0] completed successfully
Reassignment of partition [foo2,1] completed successfully

```

Decommissioning brokers

The partition reassignment tool does not have the ability to automatically generate a reassignment plan for decommissioning brokers yet. As such, the admin has to come up with a reassignment plan to move the replica for all partitions hosted on the broker to be decommissioned, to the rest of the brokers. This can be relatively tedious as the reassignment needs to ensure that all the replicas are not moved from the decommissioned broker to only one other broker. To make this process effortless, we plan to add tooling support for decommissioning brokers in the future.

Increasing replication factor

Increasing the replication factor of an existing partition is easy. Just specify the extra replicas in the custom reassignment json file and use it with the `--execute` option to increase the replication factor of the specified partitions.

For instance, the following example increases the replication factor of partition 0 of topic `foo` from 1 to 3. Before increasing the replication factor, the partition's only replica existed on broker 5. As part of increasing the replication factor, we will add more replicas on brokers 6 and 7.

The first step is to hand craft the custom reassignment plan in a json file:

```
> cat increase-replication-factor.json
{"version":1,
 "partitions":[{"topic":"foo","partition":0,"replicas":[5,6,7]]}]

```

Then, use the json file with the `--execute` option to start the reassignment process:

```
> bin/kafka-reassign-partitions.sh --zookeeper localhost:2181 --reassignment-json-file increase-replication-factor.json --execute
Current partition replica assignment

{"version":1,
 "partitions":[{"topic":"foo","partition":0,"replicas":[5]]]
}

Save this to use as the --reassignment-json-file option during rollback
Successfully started reassignment of partitions
{"version":1,
 "partitions":[{"topic":"foo","partition":0,"replicas":[5,6,7]]]
}

```

The `--verify` option can be used with the tool to check the status of the partition reassignment. Note that the same `increase-replication-factor.json` (used with the `--execute` option) should be used with the `--verify` option

```
bin/kafka-reassign-partitions.sh --zookeeper localhost:2181 --reassignment-json-file increase-replication-factor.json --verify
Status of partition reassignment:
Reassignment of partition [foo,0] completed successfully

```

You can also verify the increase in replication factor with the `kafka-topics` tool:

```
> bin/kafka-topics.sh --zookeeper localhost:2181 --topic foo --describe
Topic:foo PartitionCount:1 ReplicationFactor:3 Configs:
Topic: foo Partition: 0 Leader: 5 Replicas: 5,6,7 Isr: 5,6,7

```

Setting quotas

It is possible to set default quotas that apply to all client ids by setting these configs on the brokers. By default, each client id receives an unlimited quota. The following sets the default quota per producer and consumer client id to 10MB/sec.

```
quota.producer.default=10485760
quota.consumer.default=10485760

```

It is also possible to set custom quotas for each client.

```
> bin/kafka-configs.sh --zookeeper localhost:2181 --alter --add-config 'producer_byte_rate=1024,consumer_byte_rate=2048' --entity-name clientA --entity-type client
Updated config for clientid: "clientA".

```

Here's how to describe the quota for a given client.

```
> ./kafka-configs.sh --zookeeper localhost:2181 --describe --entity-name clientA --entity-type clients
Configs for clientid:clientA are producer_byte_rate=1024,consumer_byte_rate=2048

```

6.2 Datacenters

Some deployments will need to manage a data pipeline that spans multiple datacenters. Our recommended approach to this is to deploy a local Kafka cluster in each datacenter with application instances in each datacenter interacting only with their local cluster and mirroring between clusters (see the documentation on the `mirrormaker2` tool for how to do this).

This deployment pattern allows datacenters to act as independent entities and allows us to manage and tune inter-datacenter replication centrally. This allows each facility to stand alone and operate even if the inter-datacenter links are unavailable: when this occurs the mirroring falls behind until the link is restored at which time it catches up.

For applications that need a global view of all data you can use mirroring to provide clusters which have aggregate data mirrored from the local clusters in all datacenters. These aggregate clusters are used for reads by applications that require the full data set.

This is not the only possible deployment pattern. It is possible to read from or write to a remote Kafka cluster over the WAN, though obviously this will add whatever latency is required to get the cluster.

Kafka naturally batches data in both the producer and consumer so it can achieve high-throughput even over a high-latency connection. To allow this though it may be necessary to increase the TCP socket buffer sizes for the producer, consumer, and broker using the `socket.send.buffer.bytes` and `socket.receive.buffer.bytes` configurations. The appropriate way to set this is documented [here](#).

It is generally not advisable to run a single Kafka cluster that spans multiple datacenters over a high-latency link. This will incur very high replication latency both for Kafka writes and Zookeeper writes, and neither Kafka nor Zookeeper will remain available in all locations if the network between locations is unavailable.

6.3 Kafka Configuration

Important Client Configurations

The most important producer configurations control

- compression
- sync vs async production
- batch size (for sync producers)

The most important consumer configuration is the fetch size.

All configurations are documented in the `config` section.

As Production Server Config

Here is our server production server configuration:

```
# Replication configurations
num.replicas.fetchers=4
replica.fetch.max.bytes=1048576
replica.fetch.wait.max.ms=500
replica.high.watermark.checkpoint.interval.ms=5000
replica.socket.timeout.ms=30000
replica.socket.receive.buffer.bytes=65536
replica.log.time.max.ms=10000

controller.socket.timeout.ms=30000
controller.message.queue.size=10

# Log configuration
num.partitions=8
message.max.bytes=100000
auto.create.topics.enable=true
log.index.interval.bytes=4096
log.index.size.max.bytes=10485760
log.retention.hours=168
log.flush.interval.ms=1000
log.flush.interval.messages=20000
log.flush.scheduler.interval.ms=2000
log.roll.hours=168
log.retention.check.interval.ms=300000
log.segment.bytes=1073741824

# ZK configuration
zookeeper.connection.timeout.ms=6000
zookeeper.sync.time.ms=2000

# Socket server configuration
num.io.threads=8
num.network.threads=8
socket.request.max.bytes=104857600
socket.receive.buffer.bytes=1048576
socket.send.buffer.bytes=1048576
queued.max.requests=16
fetch.purgatory.purge.interval.requests=100
producer.purgatory.purge.interval.requests=100
```

Our client configuration varies a fair amount between different use cases.

Java Version

From a security perspective, we recommend you use the latest released version of JDK 1.8 as older freely available versions have disclosed security vulnerabilities. LinkedIn is currently running JDK 1.8 u51 (looking to upgrade to a newer version) with the G1 collector. If you decide to use the G1 collector (the current default) and you are still on JDK 1.7, make sure you are on u51 or newer. LinkedIn tried out u21 in testing, but they had a number of problems with the GC implementation in that version. LinkedIn's tuning looks like this:

```
-Xms6g -Xmx6g -XX:MetaspaceSize=96m -XX:+UseG1GC
-XX:MaxGCFreeHeapSize=20 -XX:InitiatingHeapOccupancyPercent=35 -XX:G1HeapRegionSize=16M
-XX:G1OldMetaspaceSize=50 -XX:MaxMetaspaceFreeRatio=80
```

For reference, here are the stats on one of LinkedIn's busiest clusters (at peak): 40 brokers - 50k partitions (replication factor 2) - 400k messages/sec in - 300 MB/sec inbound, 1 GB/sec outbound. The tuning looks fairly aggressive, but all of the brokers in that cluster have a 90% GC pause time of about 21ms, and they're doing less than 1 young GC per second.

6.4 Hardware and OS

We are using dual quad-core Intel Xeon machines with 24GB of memory.

You need sufficient memory to buffer active readers and writers. You can do a back-of-the-envelope estimate of memory needs by assuming you want to be able to buffer for 30 seconds and compute your memory need as `write_throughput*30`.

The disk throughput is important. We have 8x7200 rpm SATA drives. In general disk throughput is the performance bottleneck, and more disks is more better. Depending on how you configure flush behavior you may or may not benefit from more expensive disks (if you force flush often then higher RPM SAS drives may be better).

OS

Kafka should run well on any unix system and has been tested on Linux and Solaris.

We have seen a few issues running on Windows and Windows is not currently a well supported platform though we would be happy to change that.

You likely don't need to do much OS-level tuning though there are a few things that will help performance.

Two configurations that may be important:

- We upped the number of file descriptors since we have lots of topics and lots of connections.
- We upped the max socket buffer size to enable high-performance data transfer between data centers [described here](#).

Disk and Filesystem

We recommend using multiple drives to get good throughput and not sharing the same drives used for Kafka data with application logs or other OS filesystem activity to ensure good latency. You can either RAID these drives together into a single volume or format and mount each drive as its own directory. Since Kafka has replication the redundancy provided by RAID can also be provided at the application level. This choice has several tradeoffs.

If you configure multiple data directories partitions will be assigned round-robin to data directories. Each partition will be entirely in one of the data directories. If data is not well balanced among partitions this can lead to load imbalance between disks.

RAID can potentially do better at balancing load between disks (although it doesn't always seem to) because it balances load at a lower level. The primary downside of RAID is that it is usually a big performance hit for write throughput and reduces the available disk space.

Another potential benefit of RAID is the ability to tolerate disk failures. However our experience has been that rebuilding the RAID array is so I/O intensive that it effectively disables the server, so this does not provide much real availability improvement.

Applications, OS, and File Management

Kafka always immediately writes all data to the filesystem and supports the ability to configure the flush policy that controls when data is forced out of the OS cache and onto disk using the `flush` and `flush.messages`. This flush policy can be controlled to force data to disk after a period of time or after a certain number of messages has been written. There are several choices in this configuration.

Kafka must eventually call `fsync` to know that data was flushed. When recovering from a crash for any log segment not known to be `fsync'd` Kafka will check the integrity of each message by checking its CRC and also rebuild the accompanying offset index file as part of the recovery process executed on startup.

Note that durability in Kafka does not require syncing data to disk, as a failed node will always recover from its replicas.

We recommend using the default flush settings which disable application `fsync` entirely. This means relying on the background flush done by the OS and Kafka's own background flush. This provides the best of all worlds for most uses: no knobs to tune, great throughput and latency, and full recovery guarantees. We generally feel that the guarantees provided by replication are stronger than `sync` to local disk, however the paranoid still may

prefer having both and application level fsync policies are still supported.

The drawback of using application level flush settings are that this is less efficient in it's disk usage pattern (it gives the OS less leeway to re-order writes) and it can introduce latency as fsync in most Linux filesystems blocks writes to the file whereas the background flushing does much more granular page-level locking.

In general you don't need to do any low-level tuning of the filesystem, but in the next few sections we will go over some of this in case it is useful.

Understanding Linux OS flush behavior

In Linux, data written to the filesystem is maintained in `pages` until it must be written out to disk (due to an application-level fsync or the OS's own flush policy). The flushing of data is done by a set of background threads called `pdflush` (or in post 2.6.32 kernels "flusher threads").

`pdflush` has a configurable policy that controls how much dirty data can be maintained in cache and for how long before it must be written back to disk. This policy is described [here](#). When `pdflush` cannot keep up with the rate of data being written it will eventually cause the writing process to block incurring latency in the writes to slow down the accumulation of data.

You can see the current state of OS memory usage by doing

```
> cat /proc/meminfo
```

The meaning of these values are described in the link above.

Using `pagecache` has several advantages over an in-process cache for storing data that will be written out to disk:

- The I/O scheduler will batch together consecutive small writes into bigger physical writes which improves throughput.
- The I/O scheduler will attempt to re-sequence writes to minimize movement of the disk head which improves throughput.
- It automatically uses all the free memory on the machine

Ext4 Notes

`Ext4` may or may not be the best filesystem for Kafka. Filesystems like `XFS` supposedly handle locking during fsync better. We have only tried `Ext4`, though.

It is not necessary to tune these settings, however those wanting to optimize performance have a few knobs that will help:

- `data=writeback`: `Ext4` defaults to `data=ordered` which puts a strong order on some writes. Kafka does not require this ordering as it does very paranoid data recovery on all unflushed log. This setting removes the ordering constraint and seems to significantly reduce latency.
- `Disable journaling`: Journaling is a tradeoff: it makes reboots faster after server crashes but it introduces a great deal of additional locking which adds variance to write performance. Those who don't care about reboot time and want to reduce a major source of write latency spikes can turn off journaling entirely.
- `commit=rum_secs`: This tunes the frequency with which `ext4` commits to its metadata journal. Setting this to a lower value reduces the loss of unflushed data during a crash. Setting this to a higher value will improve throughput.
- `noatime`: This setting controls additional ordering guarantees when using `data=writeback` mode. This should be safe with Kafka as we do not depend on write ordering and improves throughput and latency.
- `delalloc`: Delayed allocation means that the filesystem avoid allocating any blocks until the physical write occurs. This allows `ext4` to allocate a large extent instead of smaller pages and helps ensure the data is written sequentially. This feature is great for throughput. It does seem to involve some locking in the filesystem which adds a bit of latency variance.

6.5 Monitoring

Kafka uses Yammer Metrics for metrics reporting in both the server and the client. This can be configured to report stats using pluggable stats reporters to hook up to your monitoring system.

The easiest way to see the available metrics to fire up `console` and point it at a running kafka client or server; this will all browsing all metrics with `JMX`.

We pay particular we do graphing and alerting on the following metrics:

Description	YBson name	Normal value:
Message in rate	kafka.server.type=BrokerTopicMetrics,name=MessagesInPerSec	
Byte in rate	kafka.server.type=BrokerTopicMetrics,name=BytesInPerSec	
Request rate	kafka.network.type=RequestMetrics,name=RequestsPerSec,request=(Produce)FetchConsumer FetchFollower)	
Byte out rate	kafka.server.type=BrokerTopicMetrics,name=BytesOutPerSec	
Log flush rate and time	kafka.log.type=LogFlushStats,name=LogFlushRateAndTimeMs	
# of under replicated partitions (ISR - [all replicas])	kafka.server.type=ReplicaManager,name=UnderReplicatedPartitions	0
Is controller active on broker	kafka.controller.type=KafkaController,name=ActiveControllerCount	only one broker in the cluster should have 1
Leader election rate	kafka.controller.type=ControllerStats,name=LeaderElectionRateAndTimeMs	non-zero when there are broker failures
Unclean leader election rate	kafka.controller.type=ControllerStats,name=UncleanLeaderElectionsPerSec	0
Partition counts	kafka.server.type=ReplicaManager,name=PartitionCount	mostly even across brokers
Leader replica counts	kafka.server.type=ReplicaManager,name=LeaderCount	mostly even across brokers
ISR shrink rate	kafka.server.type=ReplicaManager,name=IsShrinkPerSec	if a broker goes down, ISR for some of the partitions will shrink. When that broker is up again, ISR will be expanded once the replicas are fully caught up. Other than that, the expected value for both ISR shrink rate and expansion rate is 0.
ISR expansion rate	kafka.server.type=ReplicaManager,name=IsExpandsPerSec	See above
Max lag in messages b/w follower and leader replicas	kafka.server.type=ReplicaFetcherManager,name=MaxLag,clientId=Replica	lag should be proportional to the maximum batch size of a produce request.
Lag in messages per follower replica	kafka.server.type=FetcherLagMetrics,name=ConsumerLag,clientId=[...],(w+1)topic=[...],(w+1)partition=[0-9]+)	lag should be proportional to the maximum batch size of a produce request.
Requests waiting in the producer purgatory	kafka.server.type=ProduceRequestPurgatory,name=PurgatorySize	non-zero if ack=1 is used
Requests waiting in the fetch purgatory	kafka.server.type=FetchRequestPurgatory,name=PurgatorySize	size depends on fetch.wait.max.ms in the consumer
Request total time	kafka.network.type=RequestMetrics,name=TotalTimeMs,request=(Produce)FetchConsumer FetchFollower)	broken into queue, local, remote and response send time
Time the request waiting in the request queue	kafka.network.type=RequestMetrics,name=QueueTimeMs,request=(Produce)FetchConsumer FetchFollower)	
Time the request being processed at the leader	kafka.network.type=RequestMetrics,name=LocalTimeMs,request=(Produce)FetchConsumer FetchFollower)	
Time the request waits for the follower	kafka.network.type=RequestMetrics,name=RemoteTimeMs,request=(Produce)FetchConsumer FetchFollower)	non-zero for produce requests when ack=1

Time to send the response	kafka.network.type=RequestMetrics,name=ResponseSendTimeMs,request=(Produce FetchConsumer FetchFollower)	
Number of messages the consumer lags behind the producer by	kafka.consumer.type=ConsumerFetcherManager,name=MaxLag,clientid=[-].{w +}	
The average fraction of time the network processors are idle	kafka.network.type=SocketServer,name=NetworkProcessorAvgIdlePercent	between 0 and 1, ideally > 0.3
The average fraction of time the request handler threads are idle	kafka.server.type=KafkaRequestHandlerPool,name=RequestHandlerAvgIdlePercent	between 0 and 1, ideally > 0.3
Quota metrics per client id	kafka.server.type=(Produce Fetch),clientid=[-].{w +}	Two attributes: throttle-time indicates the amount of time in ms the client id was throttled. Ideally = 0. byte-rate indicates the data produce/consume rate of the client in bytes/sec.

[New producer monitoring](#)

The following metrics are available on new producer instances.

Metric/Attribute name	Description	MBean name
waiting-threads	The number of user threads blocked waiting for buffer memory to enqueue their records	kafka.producer.type=producer-metrics,clientid=[-].{w +}
buffer-total-bytes	The maximum amount of buffer memory the client can use (whether or not it is currently used).	kafka.producer.type=producer-metrics,clientid=[-].{w +}
buffer-available-bytes	The total amount of buffer memory that is not being used (either unallocated or in the free list).	kafka.producer.type=producer-metrics,clientid=[-].{w +}
bufferpool-wait-time	The fraction of time an appender waits for space allocation.	kafka.producer.type=producer-metrics,clientid=[-].{w +}
batch-size-avg	The average number of bytes sent per partition per request.	kafka.producer.type=producer-metrics,clientid=[-].{w +}
batch-size-max	The max number of bytes sent per partition per request.	kafka.producer.type=producer-metrics,clientid=[-].{w +}
compression-rate-avg	The average compression rate of record batches.	kafka.producer.type=producer-metrics,clientid=[-].{w +}
record-queue-time-avg	The average time in ms record batches spent in the record accumulator.	kafka.producer.type=producer-metrics,clientid=[-].{w +}
record-queue-time-max	The maximum time in ms record batches spent in the record accumulator	kafka.producer.type=producer-metrics,clientid=[-].{w +}
request-latency-avg	The average request latency in ms	kafka.producer.type=producer-metrics,clientid=[-].{w +}
request-latency-max	The maximum request latency in ms	kafka.producer.type=producer-metrics,clientid=[-].{w +}
record-send-rate	The average number of records sent per second.	kafka.producer.type=producer-metrics,clientid=[-].{w +}
records-per-request-avg	The average number of records per request.	kafka.producer.type=producer-metrics,clientid=[-].{w +}
record-retry-rate	The average per second number of retried record sends	kafka.producer.type=producer-metrics,clientid=[-].{w +}
record-error-rate	The average per second number of record sends that resulted in errors	kafka.producer.type=producer-metrics,clientid=[-].{w +}
record-size-max	The maximum record size	kafka.producer.type=producer-metrics,clientid=[-].{w +}
record-size-avg	The average record size	kafka.producer.type=producer-metrics,clientid=[-].{w +}
requests-in-flight	The current number of in-flight requests awaiting a response.	kafka.producer.type=producer-metrics,clientid=[-].{w +}
metadata-age	The age in seconds of the current producer metadata being used.	kafka.producer.type=producer-metrics,clientid=[-].{w +}
connection-close-rate	Connections closed per second in the window.	kafka.producer.type=producer-metrics,clientid=[-].{w +}
connection-creation-rate	New connections established per second in the window.	kafka.producer.type=producer-metrics,clientid=[-].{w +}
network-io-rate	The average number of network operations (reads or writes) on all connections per second.	kafka.producer.type=producer-metrics,clientid=[-].{w +}
outgoing-byte-rate	The average number of outgoing bytes sent per second to all servers.	kafka.producer.type=producer-metrics,clientid=[-].{w +}
request-rate	The average number of requests sent per second.	kafka.producer.type=producer-metrics,clientid=[-].{w +}
request-size-avg	The average size of all requests in the window.	kafka.producer.type=producer-metrics,clientid=[-].{w +}
request-size-max	The maximum size of any request sent in the window.	kafka.producer.type=producer-metrics,clientid=[-].{w +}
incoming-byte-rate	Bytes/second read off all sockets	kafka.producer.type=producer-metrics,clientid=[-].{w +}
response-rate	Responses received sent per second.	kafka.producer.type=producer-metrics,clientid=[-].{w +}
select-rate	Number of times the I/O layer checked for new I/O to perform per second	kafka.producer.type=producer-metrics,clientid=[-].{w +}
io-wait-time-ns-avg	The average length of time the I/O thread spent waiting for a socket ready for reads or writes in nanoseconds.	kafka.producer.type=producer-metrics,clientid=[-].{w +}
io-wait-ratio	The fraction of time the I/O thread spent waiting.	kafka.producer.type=producer-metrics,clientid=[-].{w +}
io-time-ns-avg	The average length of time for I/O per select call in nanoseconds.	kafka.producer.type=producer-metrics,clientid=[-].{w +}
io-ratio	The fraction of time the I/O thread spent doing I/O	kafka.producer.type=producer-metrics,clientid=[-].{w +}
connection-count	The current number of active connections.	kafka.producer.type=producer-metrics,clientid=[-].{w +}
outgoing-byte-rate	The average number of outgoing bytes sent per second for a node.	kafka.producer.type=producer-node-metrics,clientid=[-].{w +},nodeid=[0-9 +]
request-rate	The average number of requests sent per second for a node.	kafka.producer.type=producer-node-metrics,clientid=[-].{w +},nodeid=[0-9 +]
request-size-avg	The average size of all requests in the window for a node.	kafka.producer.type=producer-node-metrics,clientid=[-].{w +},nodeid=[0-9 +]
request-size-max	The maximum size of any request sent in the window for a node.	kafka.producer.type=producer-node-metrics,clientid=[-].{w +},nodeid=[0-9 +]
incoming-byte-rate	The average number of responses received per second for a node.	kafka.producer.type=producer-node-metrics,clientid=[-].{w +},nodeid=[0-9 +]
request-latency-avg	The average request latency in ms for a node.	kafka.producer.type=producer-node-metrics,clientid=[-].{w +},nodeid=[0-9 +]
request-latency-max	The maximum request latency in ms for a node.	kafka.producer.type=producer-node-metrics,clientid=[-].{w +},nodeid=[0-9 +]
response-rate	Responses received sent per second for a node.	kafka.producer.type=producer-node-metrics,clientid=[-].{w +},nodeid=[0-9 +]
record-send-rate	The average number of records sent per second for a topic.	kafka.producer.type=producer-topic-metrics,clientid=[-].{w +},topicid=[-].{w +}
byte-rate	The average number of bytes sent per second for a topic.	kafka.producer.type=producer-topic-metrics,clientid=[-].{w +},topicid=[-].{w +}

compression-rate	The average compression rate of record batches for a topic.	kafka.producertype=producer-topic-metrics,client-id=[-\.w +],topic=[-\.w +]
record-retry-rate	The average per second number of retried record sends for a topic.	kafka.producertype=producer-topic-metrics,client-id=[-\.w +],topic=[-\.w +]
record-error-rate	The average per second number of record sends that resulted in errors for a topic.	kafka.producertype=producer-topic-metrics,client-id=[-\.w +],topic=[-\.w +]
produce-throttle-time-max	The maximum time in ms a request was throttled by a broker.	kafka.producertype=producer-topic-metrics,client-id=[-\.w +]
produce-throttle-time-avg	The average time in ms a request was throttled by a broker.	kafka.producertype=producer-topic-metrics,client-id=[-\.w +]

We recommend monitor GC time and other stats and various server stats such as CPU utilization, I/O service time, etc. On the client side, we recommend monitor the message/byte rate (global and per topic), request rate/size/time, and on the consumer side, max lag in messages among all partitions and min fetch request rate. For a consumer to keep up, max lag needs to be less than a threshold and min fetch rate needs to be larger than 0.

Audit:

The final alerting we do is on the correctness of the data delivery. We audit that every message that is sent is consumed by all consumers and measure the lag for this to occur. For important topics we alert if a certain completeness is not achieved in a certain time period. The details of this are discussed in KAFKA-260.

6.7.7 ZooKeeper

Stable version

At LinkedIn, we are running ZooKeeper 3.2. Version 3.3 has known serious issues regarding ephemeral node deletion and session expirations. After running into those issues in production, we upgraded to 3.3.4 and have been running that smoothly for over a year now.

Operationalizing ZooKeeper

Operationally, we do the following for a healthy ZooKeeper installation:

- Redundancy in the physical/hardware/network layout: try not to put them all in the same rack, decent (but don't go nuts) hardware, try to keep redundant power and network paths, etc.
- I/O segregation: if you do a lot of write type traffic you'll almost definitely want the transaction logs on a different disk group than application logs and snapshots (the write to the ZooKeeper service has a synchronous write to disk, which can be slow).
- Application segregation: Unless you really understand the application patterns of other apps that you want to install on the same box, it can be a good idea to run ZooKeeper in isolation (though this can be a balancing act with the capabilities of the hardware).
- Use care with virtualization: it can work, depending on your cluster layout and read/write patterns and SLAs, but the I/O overheads introduced by the virtualization layer can add up and throw off ZooKeeper, as it can be very time sensitive.
- ZooKeeper configuration and monitoring: It's Java, make sure you give it 'enough' heap space (We usually run them with 3-5G, but that's mostly due to the data set size we have here). Unfortunately we don't have a good formula for it. As far as monitoring, both JMX and the 4 letter words (4lw) commands are very useful, they do overlap in some cases (and in those cases we prefer the 4 letter commands, they seem more predictable, or at the very least, they work better with the LI monitoring infrastructure)
- Don't overbuild the cluster: large clusters, especially in a write heavy usage pattern, means a lot of intracluster communication (quorums on the writes and subsequent cluster member updates), but don't underbuild it (and risk swamping the cluster).
- Try to run on a 3-5 node cluster: ZooKeeper writes use quorums and inherently that means having an odd number of machines in a cluster. Remember that a 5 node cluster will cause writes to slow down compared to a 3 node cluster, but will allow more fault tolerance.

Overall, we try to keep the ZooKeeper system as small as will handle the load (plus standard growth capacity planning) and as simple as possible. We try not to do anything fancy with the configuration or application layout as compared to the official release as well as keep it as self contained as possible. For these reasons, we tend to skip the OS packaged versions, since it has a tendency to try to put things in the OS standard hierarchy, which can be 'messy', for want of a better way to word it.

7. Security

7.1 Security Overview

In release 0.9.0.0, the Kafka community added a number of features that, used either separately or together, increases security in a Kafka cluster. These features are considered to be of beta quality. The following security measures are currently supported:

1. Authentication of connections to brokers from clients (producers and consumers), other brokers and tools, using either SSL or SASL (Kerberos)
2. Authentication of connections from brokers to ZooKeeper
3. Encryption of data transferred between brokers and clients, between brokers, or between brokers and tools using SSL. (Note that there is a performance degradation when SSL is enabled, the magnitude of which depends on the CPU type and the JVM implementation.)
4. Authorization of read / write operations by clients
5. Authorization is pluggable and integration with external authorization services is supported

It's worth noting that security is optional - non-secured clusters are supported, as well as a mix of authenticated, unauthenticated, encrypted and non-encrypted clients. The guides below explain how to configure and use the security features in both clients and brokers.

7.2 Encryption and Authentication using SSL

Apache Kafka allows clients to connect over SSL. By default SSL is disabled but can be turned on as needed.

1. Generate SSL key and certificate for each of the brokers

The first step of deploying HTTPS is to generate the key and the certificate for each machine in the cluster. You can use Java's keytool utility to accomplish this task. We will generate the key into a temporary keystore initially so that we can export and sign it later with CA.

```
keytool -keystore server.keystore.jks -alias localhost -validity (validity) -genkey
```

You need to specify two parameters in the above command:

1. keystore: the keystore file that stores the certificate. The keystore file contains the private key of the certificate; therefore, it needs to be kept safely.
2. validity: the valid time of the certificate in days.

Ensure that common name (CN) matches exactly with the fully qualified domain name (FQDN) of the server. The client compares the CN with the DNS domain name to ensure that it is indeed connecting to the desired server, not the malicious one.

2. Creating your own CA

After the first step, each machine in the cluster has a public-private key pair, and a certificate to identify the machine. The certificate, however, is unsigned, which means that an attacker can create such a certificate to pretend to be any machine.

Therefore, it is important to prevent forged certificates by signing them for each machine in the cluster. A certificate authority (CA) is responsible for signing certificates. CA works like a government that issues passports—the government stamps (signs) each passport so that the passport becomes difficult to forge. Other governments verify the stamps to ensure the passport is authentic. Similarly, the CA signs the certificates, and the cryptography guarantees that a signed certificate is computationally difficult to forge. Thus, as long as the CA is a genuine and trusted authority, the clients have high assurance that they are connecting to the authentic machines.

```
openssl req -new -x509 -keyout ca-key -out ca-cert -days 365
```

The generated CA is simply a public-private key pair and certificate, and it is intended to sign other certificates.

The next step is to add the generated CA to the ""clients' truststore"" so that the clients can trust this CA:

```
keytool -keystore server.truststore.jks -alias CARoot -import -file ca-cert
```

Note If you configure the Kafka brokers to require client authentication by setting `ssl.client.auth` to be "requested" or "required" on the `HostAndBrokerConfig` then you must provide a truststore for the Kafka brokers as well and it should have all the CA certificates that clients keys were signed by.

```
keytool -keystore client.truststore.jks -alias CARoot -import -file ca-cert
```

In contrast to the keystore in step 1 that stores each machine's own identity, the truststore of a client stores all the certificates that the client should trust. Importing a certificate into one's truststore also means trusting all certificates that are signed by that certificate. As the analogy above, trusting the government (CA) also means trusting all passports (certificates) that it has issued. This attribute is called

the chain of trust, and it is particularly useful when deploying SSL on a large Kafka cluster. You can sign all certificates in the cluster with a single CA, and have all machines share the same truststore that trusts the CA. That way all machines can authenticate all other machines.

3. [Signing the certificate](#)

The next step is to sign all certificates generated by step 1 with the CA generated in step 2. First, you need to export the certificate from the keystore:

```
keytool -keystore server.keystore.jks -alias localhost -certreq -file cert-file
```

Then sign it with the CA:

```
openssl x509 -req -CA ca-cert -CAkey ca-key -in cert-file -out cert-signed -days (validity) -Ccreateserial -passin pass:(ca-password)
```

Finally, you need to import both the certificate of the CA and the signed certificate into the keystore:

```
keytool -keystore server.keystore.jks -alias CARoot -import -file ca-cert
keytool -keystore server.keystore.jks -alias localhost -import -file cert-signed
```

The definitions of the parameters are the following:

- 1. keystore: the location of the keystore
- 2. ca-cert: the certificate of the CA
- 3. ca-key: the private key of the CA
- 4. ca-password: the passphrase of the CA
- 5. cert-file: the exported, unsigned certificate of the server
- 6. cert-signed: the signed certificate of the server

Here is an example of a bash script with all above steps. Note that one of the commands assumes a password of 'test1234', so either use that password or edit the command before running it.

```
#!/bin/bash
#Step 1
keytool -keystore server.keystore.jks -alias localhost -validity 365 -genkey
#Step 2
openssl req -new -x509 -keyout ca-key -out ca-cert -days 365
keytool -keystore server.truststore.jks -alias CARoot -import -file ca-cert
keytool -keystore client.truststore.jks -alias CARoot -import -file ca-cert
#Step 3
keytool -keystore server.keystore.jks -alias localhost -certreq -file cert-file
openssl x509 -req -CA ca-cert -CAkey ca-key -in cert-file -out cert-signed -days 365 -Ccreateserial -passin pass:test1234
keytool -keystore server.keystore.jks -alias CARoot -import -file ca-cert
keytool -keystore server.keystore.jks -alias localhost -import -file cert-signed
```

4. [Configuring Kafka brokers](#)

Kafka Brokers support listening for connections on multiple ports. We need to configure the following property in server.properties, which must have one or more comma-separated values:

listeners

If SSL is not enabled for inter-broker communication (see below for how to enable it), both PLAINTEXT and SSL ports will be necessary

```
listeners=PLAINTEXT://host.name:port,SSL://host.name:port
```

Following SSL configs are needed on the broker side

```
ssl.keystore.location=/var/private/ssl/kafka.server.keystore.jks
ssl.keystore.password=test1234
ssl.key.password=test1234
ssl.truststore.location=/var/private/ssl/kafka.server.truststore.jks
ssl.truststore.password=test1234
```

Optional settings that are worth considering:

- 1. ssl.client.auth=none ("required" => client authentication is required, "requested" => client authentication is requested and client without certs can still connect. The usage of "requested" is discouraged as it provides a false sense of security and misconfigured clients will still connect successfully.)
- 2. ssl.cipher.suites (Optional). A cipher suite is a named combination of authentication, encryption, MAC and key exchange algorithm used to negotiate the security settings for a network connection using TLS or SSL network protocol. (Default is an empty list)
- 3. ssl.enabled.protocols=TLSv1.2,TLSv1.1,TLSv1 (Not out the SSL protocols that you are going to accept from clients. Do note that SSL is deprecated in favor of TLS and using SSL in production is not recommended)
- 4. ssl.keystore.type=JKS
- 5. ssl.truststore.type=JKS

If you want to enable SSL for inter-broker communication, add the following to the broker properties file (it defaults to PLAINTEXT)

```
security.inter.broker.protocol=SSL
```

Due to import regulations in some countries, the Oracle implementation limits the strength of cryptographic algorithms available by default. If stronger algorithms are needed (for example, AES with 256 bit keys), the `JDK Unlimited Strength Jurisdiction Policy Files` must be obtained and installed in the `JDK/JRE`. See the [JDK/JavaSE/Tools/Documentation](#) for more information.

Once you start the broker you should be able to see in the server.log

```
with addresses: PLAINTEXT -> EndPoint(192.168.64.1,9092,PLAINTEXT),SSL -> EndPoint(192.168.64.1,9093,SSL)
```

To check quickly if the server keystore and truststore are setup properly you can run the following command

```
openssl s_client -debug -connect localhost:9093 -tlsl
```

(Note: TLSv1 should be listed under ssl.enabled.protocols)

In the output of this command you should see server's certificate:

```
-----BEGIN CERTIFICATE-----
(Variable sized random bytes)
-----END CERTIFICATE-----
subject=C=US/ST=CA/L=Santa Clara/O=org/CN=8riharsha Chintalapani
issuer=C=US/ST=CA/L=Santa Clara/O=org/CN=org/CN=kafka/emailAddress=test@test.com
```

If the certificate does not show up or if there are any other error messages than your keystore is not setup properly.

5. [Configuring Kafka clients](#)

SSL is supported only for the new Kafka Producer and Consumer, the older API is not supported. The configs for SSL will be same for both producer and consumer.

If client authentication is not required in the broker, then the following is a minimal configuration example:

```
security.protocol=SSL
ssl.truststore.location=/var/private/ssl/kafka.client.truststore.jks
ssl.truststore.password=test1234
```

If client authentication is required, then a keystore must be created like in step 1 and the following must also be configured:

```
ssl.keystore.location=/var/private/ssl/kafka.client.keystore.jks
ssl.keystore.password=test1234
ssl.key.password=test1234
```

Other configuration settings that may also be needed depending on our requirements and the broker configuration:

- 1. ssl.provider (Optional). The name of the security provider used for SSL connections. Default value is the default security provider of the JVM.
- 2. ssl.cipher.suites (Optional). A cipher suite is a named combination of authentication, encryption, MAC and key exchange algorithm used to negotiate the security settings for a network connection using TLS or SSL network protocol.
- 3. ssl.enabled.protocols=TLSv1.2,TLSv1.1,TLSv1. It should list at least one of the protocols configured on the broker side
- 4. ssl.truststore.type=JKS
- 5. ssl.keystore.type=JKS

Examples using console-producer and console-consumer:

```
kafka-console-producer.sh --broker-list localhost:9093 --topic test --producer.config client-ssl.properties
kafka-console-consumer.sh --bootstrap-server localhost:9093 --topic test --zoo-consumer --consumer.config client-ssl.properties
```

7. [Authentication using SNI/SSL](#)

1. [Prerequisites](#)

1. [Kerberos](#)

If your organization is already using a Kerberos server (for example, by using Active Directory), there is no need to install a new server just for Kafka. Otherwise you will need to install one, your Linux vendor likely has packages for Kerberos and a short guide on how to install and configure it ([Ubuntu](#), [Redhat](#)). Note that if you are using Oracle Java, you will need to download JCE policy files for your Java version and copy them to `$JAVA_HOME/jre/lib/security`.

2. [Create Kerberos Principals](#)

If you are using the organization's Kerberos or Active Directory server, ask your Kerberos administrator for a principal for each Kafka broker in your cluster and for every operating system user that will access Kafka with Kerberos authentication (via clients and tools).
If you have installed your own Kerberos, you will need to create these principals yourself using the following commands:

```
sudo /usr/sbin/kadmin.local -q 'addprinc -randkey kafka/(hostname)$(REALM)'  
sudo /usr/sbin/kadmin.local -q 'xctadd -s /etc/security/keytabs/kafkabroker.keytab kafka/(hostname)$(REALM)'
```

- 3. Make sure all hosts can resolve the hostnames using `hostname`. It is a Kerberos requirement that all your hosts can be resolved with their FQDNs.

2. [Configuring the Kafka brokers](#)

- 1. Add a suitably modified JAAS file similar to the one below to each Kafka broker's config directory, let's call it `kafka_server_jaas.conf` for this example (note that each broker should have its own keytab):

```
KafkaServer {  
  com.sun.security.auth.module.Krb5LoginModule required  
  useKeyTab=true  
  storeKey=true  
  keyTab="/etc/security/keytabs/kafka_server.keytab"  
  principal="kafka/kafkal.hostname.com@EXAMPLE.COM";  
}  
  
// Zookeeper client authentication  
Client {  
  com.sun.security.auth.module.Krb5LoginModule required  
  useKeyTab=true  
  storeKey=true  
  keyTab="/etc/security/keytabs/kafka_client.keytab"  
  principal="kafka/kafkal.hostname.com@EXAMPLE.COM";  
}
```

- 2. Pass the JAAS and optionally the `krb5` file locations as `JVM` parameters to each Kafka broker (see [here](#) for more details):

```
-Djava.security.krb5.conf=/etc/kafka/krb5.conf  
-Djava.security.auth.login.config=/etc/kafka/kafka_server_jaas.conf
```

- 3. Make sure the keytabs configured in the JAAS file are readable by the operating system user who is starting kafka broker.
- 4. Configure a SASL port in `server.properties`, by adding at least one of `SASL_PLAINTEXT` or `SASL_SSL` to the `listeners` parameter, which contains one or more comma-separated values:

```
listeners=SASL_PLAINTEXT://host.name:port
```

If `SASL_SSL` is used, then [SSL must also be configured](#). If you are only configuring a SASL port (or if you want the Kafka brokers to authenticate each other using SASL) then make sure you set the same SASL protocol for inter-broker communication:

```
security.inter.broker.protocol=SASL_PLAINTEXT (or SASL_SSL)
```

We must also configure the service name in `server.properties`, which should match the principal name of the Kafka brokers. In the above example, principal is `"kafka/kafkal.hostname.com@EXAMPLE.COM"`, so:

```
sasl.kerberos.service.name=kafka
```

Important notes:

- 1. `KafkaServer` is a section name in JAAS file used by each `KafkaServer/Broker`. This section tells the broker which principal to use and the location of the keytab where this principal is stored. It allows the broker to login using the keytab specified in this section.
- 2. Client section is used to authenticate a SASL connection with zookeeper. It also allows the brokers to set SASL ACL on zookeeper nodes which locks these nodes down so that only the brokers can modify it. It is necessary to have the same principal name across all brokers. If you want to use a section name other than `Client`, set the system property `zookeeper.sasl.client` to the appropriate name (e.g., `-Dzookeeper.sasl.client=KClient`).
- 3. Zookeeper uses 'zookeeper' as the service name by default. If you want to change this, set the system property `zookeeper.sasl.client.service.name` to the appropriate name (e.g., `-Dzookeeper.sasl.client.service.name=zk`).

3. [Configuring the Kafka clients](#)

SASL authentication is only supported for the new kafka producer and consumer, the older API is not supported. To configure SASL authentication on the clients:

- 1. Clients (producers, consumers, connect workers, etc) will authenticate to the cluster with their own principal (usually with the same name as the user running the client), so obtain or create these principals as needed. Then create a JAAS file for each principal. The `KafkaClient` section describes how the clients like producer and consumer can connect to the Kafka Broker. The following is an example configuration for a client using a keytab (recommended for long-running processes):

```
KafkaClient {  
  com.sun.security.auth.module.Krb5LoginModule required  
  useKeyTab=true  
  storeKey=true  
  keyTab="/etc/security/keytabs/kafka_client.keytab"  
  principal="kafka-client-1@EXAMPLE.COM";  
}  
  
For command-line utilities like kafka-console-consumer or kafka-console-producer, kinit can be used along with "useTicketCache=true" as in:
```

```
KafkaClient {  
  com.sun.security.auth.module.Krb5LoginModule required  
  useTicketCache=true;  
}
```

- 2. Pass the JAAS and optionally `krb5` file locations as `JVM` parameters to each client `JVM` (see [here](#) for more details):

```
-Djava.security.krb5.conf=/etc/kafka/krb5.conf  
-Djava.security.auth.login.config=/etc/kafka/kafka_client_jaas.conf
```

- 3. Make sure the keytabs configured in the `kafka_client_jaas.conf` are readable by the operating system user who is starting kafka client.
- 4. Configure the following properties in producer properties or consumer properties:

```
security.protocol=SASL_PLAINTEXT (or SASL_SSL)  
sasl.kerberos.service.name=kafka
```

4. [Interoperating Security Features in a Running Cluster](#)

You can secure a running cluster via one or more of the supported protocols discussed previously. This is done in phases:

- Incrementally bounce the cluster nodes to open additional secured port(s).
- Restart clients using the secured rather than PLAINTEXT port (assuming you are securing the client-broker connection).
- Incrementally bounce the cluster again to enable broker-to-broker security (if this is required)
- A final incremental bounce to close the PLAINTEXT port.

The specific steps for configuring SSL and SASL are described in sections 7.2 and 7.3. Follow these steps to enable security for your desired protocol(s).

The security implementation lets you configure different protocols for both broker-client and broker-broker communication. These must be enabled in separate bounces. A PLAINTEXT port must be left open throughout so brokers and/or clients can continue to communicate.

When performing an incremental bounce stop the brokers cleanly via a `SIGTERM`. It's also good practice to wait for restarted replicas to return to the ISR list before moving onto the next node.

As an example, say we wish to encrypt both broker-client and broker-broker communication with SSL. In the first incremental bounce, a SSL port is opened on each node:

```
listeners=PLAINTEXT://broker1:9091,SSL://broker1:9092
```

We then restart the clients, changing their config to point at the newly opened, secured port:

```
bootstrap.servers = [broker1:9092,...]  
security.protocol = SSL  
...etc
```

In the second incremental server bounce we instruct Kafka to use SSL as the broker-broker protocol (which will use the same SSL port):

```
listeners=PLAINTEXT://broker1:9091,SSL://broker1:9092  
security.inter.broker.protocol=SSL
```

In the final bounce we secure the cluster by closing the PLAINTEXT port:

```
listeners=SSL://broker1:9092  
security.inter.broker.protocol=SSL
```

Alternatively we might choose to open multiple ports so that different protocols can be used for broker-broker and broker-client communication. Say we wished to use SSL encryption throughout (i.e. for broker-broker and broker-client communication) but we'd like to add SASL authentication to the broker-client connection also. We would achieve this by opening two additional ports during the first bounce:

```
listeners=PLAINTEXT://broker1:9091,SSL://broker1:9092,SASL_SSL://broker1:9093
```

We would then restart the clients, changing their config to point at the newly opened, SASL & SSL secured port:

```
bootstrap.servers = [broker1:9093,...]
security.protocol = SASL_SSL
...etc
```

The second server bounce would switch the cluster to use encrypted broker-broker communication via the SSL port we previously opened on port 9092:

```
listeners=PLAINTEXT://broker1:9091,SSL://broker1:9092,SASL_SSL://broker1:9093
security.inter.broker.protocol=SSL
```

The final bounce secures the cluster by closing the PLAINTEXT port.

```
listeners=SSL://broker1:9092,SASL_SSL://broker1:9093
security.inter.broker.protocol=SSL
```

ZooKeeper can be secured independently of the Kafka cluster. The steps for doing this are covered in section 7.5.2.

7.6.AuthorizationandACLs

Kafka ships with a pluggable Authorizer and an out-of-box authorizer implementation that uses zookeeper to store all the acls. Kafka acls are defined in the general format of "Principal P is [Allowed/Denied] Operation O From Host H On Resource R". You can read more about the acl structure on KIP-11. In order to add, remove or list acls you can use the Kafka authorizer CLI. By default, if a Resource R has no associated acls, no one other than super users is allowed to access R. If you want to change that behavior, you can include the following in broker.properties.

```
allow.everyone.if.no.acl.found=true
```

One can also add super users in broker.properties like the following (note that the delimiter is semicolon since SSL user names may contain comma).

```
super.users=User:Bob;User:Alice
```

By default, the SSL user name will be of the form

"CN=writeruser,OJ=Unknown,C=Unknown,L=Unknown,ST=Unknown,C=Unknown". One can change that by setting a customized PrincipalBuilder in broker.properties like the following.

```
principal.builder.class=CustomizedPrincipalBuilderClass
```

By default, the SASL user name will be the primary part of the Kerberos principal. One can change that by setting sasl.kerberos.principal.to.local.rules to a customized rule in broker.properties. The format of sasl.kerberos.principal.to.local.rules is a list where each rule works in the same way as the auth_to_local in [Kerberos configuration file \(krb5.conf\)](#). Each rule starts with RULE: and contains an expression in the format (istring)(regex)(pattern)(replacement)(g). See the kerberos documentation for more details. An example of adding a rule to properly translate user@MYDOMAIN.COM to user while also keeping the default rule in place is:

```
sasl.kerberos.principal.to.local.rules=RULE:[!@!$%&'](.*)@MYDOMAIN.COM)/8.*///,DEFAULT
```

Command-LineInterface

Kafka Authorization management CLI can be found under bin directory with all the other CLIs. The CLI script is called **kafka-acls.sh**. Following lists all the options that the script supports:

Option	Description	Default	Option type
-add	Indicates to the script that user is trying to add an acl.		Action
-remove	Indicates to the script that user is trying to remove an acl.		Action
-list	Indicates to the script that user is trying to list acls.		Action
-authorizer	Fully qualified class name of the authorizer.	kafka.security.auth.SimpleAclAuthorizer	Configuration
-authorizer-properties	Key/value pairs that will be passed to authorizer for initialization. For the default authorizer the example values are: zookeeper.connect=localhost:2181		Configuration
-cluster	Specifies cluster as resource.		Resource
-topic [topic-name]	Specifies the topic as resource.		Resource
-group [group-name]	Specifies the consumer group as resource.		Resource
-allow-principal	Principal is in PrincipalType name format that will be added to ACL with Allow permission. You can specify multiple -allow-principal in a single command.		Principal
-deny-principal	Principal is in PrincipalType name format that will be added to ACL with Deny permission. You can specify multiple -deny-principal in a single command.		Principal
-allow-host	IP address from which principals listed in -allow-principal will have access.	if -allow-principal is specified defaults to " * " which translates to "all hosts"	Host
-deny-host	IP address from which principals listed in -deny-principal will be denied access.	if -deny-principal is specified defaults to " * " which translates to "all hosts"	Host
-operation	Operation that will be allowed or denied. Valid values are: Read, Write, Create, Delete, Alter, Describe, ClusterAction, All	All	Operation
-producer	Convenience option to add/remove acls for producer role. This will generate acls that allows WRITE, DESCRIBE on topic and CREATE on cluster.		Convenience
-consumer	Convenience option to add/remove acls for consumer role. This will generate acls that allows READ, DESCRIBE on topic and READ on consumer group.		Convenience

Examples

• AddingAcls

Suppose you want to add an acl "Principals User:Bob and User:Alice are allowed to perform Operation Read and Write on Topic Test-Topic from IP 198.51.100.0 and IP 198.51.100.1". You can do that by executing the CLI with following options:

```
bin/kafka-acls.sh --authorizer-properties zookeeper.connect=localhost:2181 --add --allow-principal User:Bob --allow-principal User:Alice --allow-host 198.51.100.
```

By default all principals that don't have an explicit acl that allows access for an operation to a resource are denied. In rare cases where an allow acl is defined that allows access to all but some principal we will have to use the -deny-principal and -deny-host option. For example, if we want to allow all users to Read from Test-topic but only deny User:BadBob from IP 198.51.100.3 we can do so using following commands:

```
bin/kafka-acls.sh --authorizer-properties zookeeper.connect=localhost:2181 --add --allow-principal User:* --allow-host * --deny-principal User:BadBob --deny-host
```

Note that '-allow-host' and '-deny-host' only support IP addresses (hostnames are not supported). Above examples add acls to a topic by specifying -topic [topic-name] as the resource option. Similarly user can add acls to cluster by specifying -cluster and to a consumer group by specifying -group [group-name].

• RemovingAcls

Removing acls is pretty much the same. The only difference is instead of -add option users will have to specify -remove option. To remove the acls added by the first example above we can execute the CLI with following options:

```
bin/kafka-acls.sh --authorizer-properties zookeeper.connect=localhost:2181 --remove --allow-principal User:Bob --allow-principal User:Alice --allow-host 198.51.
```

• ListingAcls

We can list acls for any resource by specifying the -list option with the resource. To list all acls for Test-topic we can execute the CLI with following options:

```
bin/kafka-acls.sh --authorizer-properties zookeeper.connect=localhost:2181 --list --topic Test-topic
```

• Addingorremovingapincipalasproducerorconsumer

The most common use case for acl management are adding/removing a principal as producer or consumer so we added convenience options to handle these cases. In order to add User:Bob as a producer of Test-topic we can execute the following command:

```
bin/kafka-acls.sh --authorizer-properties zookeeper.connect=localhost:2181 --add --allow-principal User:Bob --producer --topic Test-topic
```

Similarly to add Alice as a consumer of Test-topic with consumer group Group-1 we just have to pass

```
--consumer option:
bin/kafka-acls.sh --authorizer-properties zookeeper.connect=localhost:2181 --add --allow-principal User:Bob --consumer --topic test-topic --group Group-1

Note that for consumer option we must also specify the consumer group. In order to remove a principal from producer or consumer role we just need to pass --remove option.
```

7.5.ZooKeeper Authentication

7.5.12.How to enable

To enable ZooKeeper authentication on brokers, there are two necessary steps:

- 1. Create a JAAS login file and set the appropriate system property to point to it as described above
- 2. Set the configuration property zookeeper.security.acl to true in each broker to true

The metadata stored in ZooKeeper is such that only brokers will be able to modify the corresponding znodes, but znodes are world readable. The rationale behind this decision is that the data stored in ZooKeeper is not sensitive, but inappropriate manipulation of znodes can cause cluster disruption. We also recommend limiting the access to ZooKeeper via network segmentation (only brokers and some admin tools need access to ZooKeeper if the new consumer and new producer are used).

7.5.13.Migrating clusters

If you are running a version of Kafka that does not support security or simply with security disabled, and you want to make the cluster secure, then you need to execute the following steps to enable ZooKeeper authentication with minimal disruption to your operations:

- 1. Perform a rolling restart setting the JAAS login file, which enables brokers to authenticate. At the end of the rolling restart, brokers are able to manipulate znodes with strict ACLs, but they will not create znodes with those ACLs
- 2. Perform a second rolling restart of brokers, this time setting the configuration parameter zookeeper.security.acl to true, which enables the use of secure ACLs when creating znodes
- 3. Execute the ZKSecurityMigrator tool. To execute the tool, there is this script - ./bin/zookeeper-security-migration.sh with zookeeper.security.acl set to secure. This tool traverses the corresponding sub-trees changing the ACLs of the znodes

It is also possible to turn off authentication in a secure cluster. To do it, follow these steps:

- 1. Perform a rolling restart of brokers setting the JAAS login file, which enables brokers to authenticate, but setting zookeeper.security.acl to false. At the end of the rolling restart, brokers stop creating znodes with secure ACLs, but are still able to authenticate and manipulate all znodes
- 2. Execute the ZKSecurityMigrator tool. To execute the tool, run this script - ./bin/zookeeper-security-migration.sh with zookeeper.security.acl set to insecure. This tool traverses the corresponding sub-trees changing the ACLs of the znodes
- 3. Perform a second rolling restart of brokers, this time omitting the system property that sets the JAAS login file

Here is an example of how to run the migration tool:

```
./bin/zookeeper-security-migration --zookeeper.acl=secure --zookeeper.connection=localhost:2181
```

Run this to see the full list of parameters:

```
./bin/zookeeper-security-migration --help
```

7.5.14.Migrating the ZooKeeper ensemble

It is also necessary to enable authentication on the ZooKeeper ensemble. To do it, we need to perform a rolling restart of the server and set a few properties. Please refer to the ZooKeeper documentation for more detail:

- 1. [Apache ZooKeeper documentation](#)
- 2. [Apache ZooKeeper wiki](#)

8.Kafka Connect

8.1.Overview

Kafka Connect is a tool for scalability and reliably streaming data between Apache Kafka and other systems. It makes it simple to quickly define connectors that move large collections of data into and out of Kafka. Kafka Connect can ingest entire databases or collect metrics from all your application servers into Kafka topics, making the data available for stream processing with low latency. An export job can deliver data from Kafka topics into secondary storage and query systems or into batch systems for offline analysis. Kafka Connect features include:

- **Common framework for Kafka connectors:** Kafka Connect standardizes integration of other data systems with Kafka, simplifying connector development, deployment, and management
- **Distributed and standalone modes:** scale up to a large, centrally managed service supporting an entire organization or scale down to development, testing, and small production deployments
- **REST interface:** submit and manage connectors to your Kafka Connect cluster via an easy to use REST API
- **Automatic offset management:** with just a little information from connectors, Kafka Connect can manage the offset commit process automatically so connector developers do not need to worry about this error prone part of connector development
- **Distributed and scalability defaults:** Kafka Connect builds on the existing
- **Streaming batch integration:** leveraging Kafka's existing capabilities, Kafka Connect is an ideal solution for bridging streaming and batch data systems

8.2.User Guide

The quickstart provides a brief example of how to run a standalone version of Kafka Connect. This section describes how to configure, run, and manage Kafka Connect in more detail.

Running Kafka Connect

Kafka Connect currently supports two modes of execution: standalone (single process) and distributed. In standalone mode all work is performed in a single process. This configuration is simpler to setup and get started with and may be useful in situations where only one worker makes sense (e.g. collecting log files), but it does not benefit from some of the features of Kafka Connect such as fault tolerance. You can start a standalone process with the following command:

```
> bin/connect-standalone.sh config/connect-standalone.properties connector1.properties [connector2.properties ...]
```

The first parameter is the configuration for the worker. This includes settings such as the Kafka connection parameters, serialization format, and how frequently to commit offsets. The provided example should work well with a local cluster running with the default configuration provided by config/server.properties. It will require tweaking to use with a different configuration or production deployment. The remaining parameters are connector configuration files. You may include as many as you want, but all will execute within the same process (on different threads). Distributed mode handles automatic balancing of work, allows you to scale up (or down) dynamically, and offers fault tolerance both in the active tasks and for configuration and offset commit data. Execution is very similar to standalone mode:

```
> bin/connect-distributed.sh config/connect-distributed.properties
```

The difference is in the class which is started and the configuration parameters which change how the Kafka Connect process decides where to store configurations, how to assign work, and where to store offsets. In particular, the following configuration parameters are critical to set before starting your cluster:

- **group.id** (default connect-cluster) - unique name for the cluster, used in forming the Connect cluster group; note that this **must not conflict** with consumer group IDs
- **config.storage.topic** (default connect-configs) - topic to use for storing connector and task configurations; note that this should be a single partition, highly replicated topic
- **offset.storage.topic** (default connect-offsets) - topic to use for this topic should have many partitions and be replicated

Note that in distributed mode the connector configurations are not passed on the command line. Instead, use the REST API described below to create, modify, and destroy connectors.

Configuring Connectors

Connector configurations are simple key value mappings. For standalone mode these are defined in a properties file and passed to the Connect process on the command line. In distributed mode, they will be included in the JSON payload for the request that creates (or modifies) the connector. Most configurations are connector dependent, so they can't be outlined here. However, there are a few common options:

- **name** - Unique name for the connector. Attempting to register again with the same name will fail.
- **connector.class** - The Java class for the connector
- **tasks.max** - The maximum number of tasks that should be created for this connector. The connector may create fewer tasks if it cannot achieve this level of parallelism.

Sink connectors also have one additional option to control their input:

- `topics` - A list of topics to use as input for this connector

For any other options, you should consult the documentation for the connector.

REST API

Since Kafka Connect is intended to be run as a service, it also supports a REST API for managing connectors. By default this service runs on port 8083. The following are the currently supported endpoints:

- `GET /connectors` - return a list of active connectors
- `POST /connectors` - create a new connector; the request body should be a JSON object containing a `string name` field and a `object config` field with the connector configuration parameters
- `GET /connectors/{name}` - get information about a specific connector
- `GET /connectors/{name}/config` - get the configuration parameters for a specific connector
- `PUT /connectors/{name}/config` - update the configuration parameters for a specific connector
- `GET /connectors/{name}/tasks` - get a list of tasks currently running for a connector
- `DELETE /connectors/{name}` - delete a connector, halting all tasks and deleting its configuration

8.3 Connector Development Guide

This guide describes how developers can write new connectors for Kafka Connect to move data between Kafka and other systems. It briefly reviews a few key concepts and then describes how to create a simple connector.

Core Concepts and APIs

Connectors and Tasks

To copy data between Kafka and another system, users create a `Connector` for the system they want to pull data from or push data to. Connectors come in two flavors: `SourceConnectors` import data from another system (e.g. `JDBCSourceConnector` would import a relational database into Kafka) and `SinkConnectors` export data (e.g. `HDFSSinkConnector` would export the contents of a Kafka topic to an HDFS file). Connectors do not perform any data copying themselves; their configuration describes the data to be copied, and the `Connector` is responsible for breaking that job into a set of `Tasks` that can be distributed to workers. These `Tasks` also come in two corresponding flavors: `SourceTask` and `SinkTask`. With an assignment in hand, each `Task` must copy its subset of the data to or from Kafka. In Kafka Connect, it should always be possible to frame these assignments as a set of input and output streams consisting of records with consistent schemas. Sometimes this mapping is obvious: each file in a set of log files can be considered a stream with each parsed line forming a record using the same schema and offsets stored as byte offsets in the file. In other cases it may require more effort to map to this model: a `JDBC` connector can map each table to a stream, but the offset is less clear. One possible mapping uses a timestamp column to generate queries incrementally returning new data, and the last queried timestamp can be used as the offset.

Streams and Records

Each stream should be a sequence of key-value records. Both the keys and values can have complex structure -- many primitive types are provided, but arrays, objects, and nested data structures can be represented as well. The runtime data format does not assume any particular serialization format; this conversion is handled internally by the framework. In addition to the key and value, records (both those generated by sources and those delivered to sinks) have associated stream IDs and offsets. These are used by the framework to periodically commit the offsets of data that have been processed so that in the event of failures, processing can resume from the last committed offsets, avoiding unnecessary reprocessing and duplication of events.

Dynamic Connectors

Not all jobs are static, so `Connector` implementations are also responsible for monitoring the external system for any changes that might require reconfiguration. For example, in the `JDBCSourceConnector` example, the `Connector` might assign a set of tables to each `Task`. When a new table is created, it must discover this so it can assign the new table to one of the `Tasks` by updating its configuration. When it notices a change that requires reconfiguration (or a change in the number of `Tasks`), it notifies the framework and the framework updates any corresponding `Tasks`.

Developing a Simple Connector

Developing a connector only requires implementing two interfaces, the `Connector` and `Task`. A simple example is included with the source code for Kafka in the `file` package. This connector is meant for use in standalone mode and has implementations of a `SourceConnector/SourceTask` to read each line of a file and emit it as a record and a `SinkConnector/SinkTask` that writes each record to a file. The rest of this section will walk through some code to demonstrate the key steps in creating a connector, but developers should also refer to the full example source code as many details are omitted for brevity.

Connector Example

We'll cover the `SourceConnector` as a simple example. `SinkConnector` implementations are very similar. Start by creating the class that inherits from `SourceConnector` and add a couple of fields that will store parsed configuration information (the filename to read from and the topic to send data to):

```
public class FileStreamSourceConnector extends SourceConnector {
    private String filename;
    private String topic;
```

The easiest method to fill in is `getTaskClass()`, which defines the class that should be instantiated in worker processes to actually read the data:

```
@Override
public Class getTaskClass() {
    return FileStreamSourceTask.class;
}
```

We will define the `FileStreamSourceTask` class below. Next, we add some standard lifecycle methods, `start()` and `stop()`:

```
@Override
public void start(Map<String, String> props) {
    // The complete version includes error handling as well.
    filename = props.get(FILE_CONFIG);
    topic = props.get(TOPIC_CONFIG);
}

@Override
public void stop() {
    // Nothing to do since no background monitoring is required.
}
```

Finally, the real core of the implementation is in `getTaskConfigs()`. In this case we're only handling a single file, so even though we may be permitted to generate more tasks as per the `maxTasks` argument, we return a list with only one entry:

```
@Override
public List<Map<String, String>> getTaskConfigs(int maxTasks) {
    ArrayList<Map<String, String>> configs = new ArrayList<>();
    // Only one input stream makes sense.
    Map<String, String> config = new Map<>();
    if (filename != null)
        config.put(FILE_CONFIG, filename);
    config.put(TOPIC_CONFIG, topic);
    configs.add(config);
    return configs;
}
```

Even with multiple tasks, this method implementation is usually pretty simple. It just has to determine the number of input tasks, which may require contacting the remote service it is pulling data from, and then divvy them up. Because some patterns for splitting work among tasks are so common, some utilities are provided in `ConnectorUtils` to simplify these cases. Note that this simple example does not include dynamic input. See the discussion in the next section for how to trigger updates to task configs.

Task Example: SourceTask

Next we'll describe the implementation of the corresponding `SourceTask`. The implementation is short, but too long to cover completely in this guide. We'll use pseudo code to describe most of the implementation, but you can refer to the source code for the full example. Just as with the connector, we need to create a class inheriting from the appropriate base `Task` class. It also has some standard lifecycle methods:

```
public class FileStreamSourceTask extends SourceTask<Object, Object> {
    String filename;
    InputStream stream;
    String topic;

    public void start(Map<String, String> props) {
        filename = props.get(FileStreamSourceConnector.FILE_CONFIG);
        stream = openInputStream(filename);
        topic = props.get(FileStreamSourceConnector.TOPIC_CONFIG);
    }

    @Override
    public synchronized void stop() {
        stream.close()
    }
}
```

These are slightly simplified versions, but show that these methods should be relatively simple and the only work they should perform is allocating or freeing resources. There are two points to note about this implementation. First, the `start()` method does not yet handle resuming from a previous offset, which will be addressed in a later section. Second, the `stop()` method is synchronous. This will be necessary because `SourceTasks` are given a dedicated thread which they can block indefinitely, so they need to be stopped with a call from a different thread in the `Worker`. Next, we implement the main functionality of the task, the `poll()` method which gets events from the input system and returns a `List<SourceRecord>`:

```
@Override
public List<SourceRecord> poll() throws InterruptedException {
    try {
        ArrayList<SourceRecord> records = new ArrayList<>();
        while (streamValid(stream) && records.isEmpty()) {
            SourceOffset line = readOffsetLine(stream);
            if (line != null) {
                Map sourcePartition = Collections.singletonMap("filename", filename);
                Map sourceOffset = Collections.singletonMap("position", streamOffset);
                records.add(new SourceRecord(sourcePartition, sourceOffset, topic, Schema.STRING_SCHEMA, line));
            } else {
                Thread.sleep(1);
            }
        }
        return records;
    } catch (IOException e) {
        // Underlying stream was killed, probably as a result of calling stop. Allow to return
        // null, and driving thread will handle any shutdown if necessary.
    }
    return null;
}
```

Again, we've omitted some details, but we can see the important steps: the `poll()` method is going to be called repeatedly, and for each call it will loop trying to read records from the file. For each line it reads, it also tracks the file offset. It uses this information to create an output `SourceRecord` with four pieces of information: the source partition (there is only one, the single file being read), source offset (byte offset in the file), output topic name, and output value (the line, and we include a schema indicating this value will always be a string). Other variants of the `SourceRecord` constructor can also include a specific output partition and a key. Note that this implementation uses the normal `Java Input Stream` interface and may sleep if data is not available. This is acceptable because Kafka Connect provides each task with a dedicated thread. While task implementations have to conform to the basic `poll()` interface, they have a lot of flexibility in how they are implemented. In this case, an NIO-based implementation would be more efficient, but this simple approach works, is quick to implement, and is compatible with older versions of Java.

SinkTasks

The previous section described how to implement a simple `SourceTask`. Unlike `SourceConnector` and `SinkConnector`, `SourceTask` and `SinkTask` have very different interfaces because `SourceTasks` uses a pull interface and `SinkTasks` uses a push interface. Both share the common lifecycle methods, but the `SinkTask` interface is quite different:

```
public abstract class SinkTask implements Task {
    public void initialize(SinkTaskContext context) { ... }

    public abstract void put(Collection<SinkRecord> records);
    public abstract void flush(Map<TopicPartition, Long> offsets);
}
```

The `SinkTask` documentation contains full details, but this interface is nearly as simple as the `SourceTask`. The `put()` method should contain most of the implementation, accepting sets of `SinkRecords`, performing any required translation, and storing them in the destination system. This method does not need to ensure the data has been fully written to the destination system before returning. In fact, in many cases internal buffering will be useful so an entire batch of records can be sent at once, reducing the overhead of inserting events into the downstream data store. The `SinkRecords` contain essentially the same information as `SourceRecords`: Kafka topic, partition, offset and the event key and value. The `flush()` method is used during the offset commit process, which allows tasks to recover from failures and resume from a safe point such that no events will be missed. The method should push any outstanding data to the destination system and then block until the write has been acknowledged. The `offsets` parameter can often be ignored, but is useful in some cases where implementations want to store offset information in the destination store to provide exactly-once delivery. For example, an HDFS connector could do this and use atomic move operations to make sure the `flush()` operation atomically commits the data and offsets to a final location in HDFS.

Resuming from Previous Offsets

The `SourceTask` implementation included a stream ID (the input filename) and offset (position in the file) with each record. The framework uses this to commit offsets periodically so that in the case of a failure, the task can recover and minimize the number of events that are reprocessed and possibly duplicated (or to resume from the most recent offset if Kafka Connect was stopped gracefully, e.g. in standalone mode or due to a job reconfiguration). This commit process is completely automated by the framework, but only the connector knows how to seek back to the right position in the input stream to resume from that location. To correctly resume upon startup, the task can use the `SourceContext` passed into its `initialize()` method to access the offset data. In `initialize()`, we would add a bit more code to read the offset (if it exists) and seek to that position:

```
Stream = new FileInputStream(filename);
Map<String, Object> offset = context.offsetStorageReader().offset(Collections.singletonMap(FILENAME_FIELD, filename));
if (offset != null) {
    Long lastRecordedOffset = (Long) offset.get("position");
    if (lastRecordedOffset != null)
        seekToOffset(stream, lastRecordedOffset);
}
```

Of course, you might need to read many keys for each of the input streams. The `OffsetStorageReader` interface also allows you to issue bulk reads to efficiently load all offsets, then apply them by seeking each input stream to the appropriate position.

Dynamic Input/Output Streams

Kafka Connect is intended to define bulk data copying jobs, such as copying an entire database rather than creating many jobs to copy each table individually. One consequence of this design is that the set of input or output streams for a connector can vary over time. Source connectors need to monitor the source system for changes, e.g. table additions/deletions in a database. When they pick up changes, they should notify the framework via the `ConnectorContext` object that reconfiguration is necessary. For example, in a `SourceConnector`:

```
if (inputsChanged())
    this.context.requestTaskReconfiguration();
```

The framework will promptly request new configuration information and update the tasks, allowing them to gracefully commit their progress before reconfiguring them. Note that in the `SourceConnector`, this monitoring is currently left up to the connector implementation. If an extra thread is required to perform this monitoring, the connector must allocate it itself. Ideally this code for monitoring changes would be isolated to the `Connector` and tasks would not need to worry about them. However, changes can also affect tasks, most commonly when one of their input streams is destroyed in the input system, e.g. if a table is dropped from a database. If the `Task` encounters the issue before the `Connector`, which will be common if the `Connector` needs to poll for changes, the `Task` will need to handle the subsequent error. Thankfully, this can usually be handled simply by catching and handling the appropriate exception. `SinkConnectors` usually only have to handle the addition of streams, which may translate to new entries in their outputs (e.g. a new database table). The framework manages any changes to the Kafka input, such as when the set of input topics changes because of a regex subscription. `SinkTasks` should expect new input streams, which may require creating new resources in the downstream system, such as a new table in a database. The trickiest situation to handle in these cases may be conflicts between multiple `SinkTasks` seeing a new input stream for the first time and simultaneously trying to create the new resource. `SinkConnectors`, on the other hand, will generally require no special code for handling a dynamic set of streams.

Working with Schemas

The `FileStream` connectors are good examples because they are simple, but they also have trivially structured data – each line is just a string. Almost all practical connectors will need schemas with more complex data formats. To create more complex data, you'll need to work with the Kafka Connect data API. Most structured records will need to interact with two classes in addition to primitive types: `Schema` and `Struct`. The API documentation provides a complete reference, but here is a simple example creating a `Schema` and `Struct`:

```
Schema schema = SchemaBuilder.struct().name("NAME")
    .field("name", Schema.STRING_SCHEMA)
    .field("age", Schema.INT_SCHEMA)
    .field("isAdmin", new SchemaBuilder.boolean().defaultValue(false).build())
    .build();

Struct struct = new Struct(schema)
    .put("name", "Barbara Lislov")
    .put("age", 75)
    .build();
```

If you are implementing a source connector, you'll need to decide when and how to create schemas. Where possible, you should avoid recomputing them as much as possible. For example, if your connector is guaranteed to have a fixed schema, create it statically and reuse a single instance. However, many connectors will have dynamic schemas. One simple example of this is a database connector. Considering even just a single

table, the schema will not be predefined for the entire connector (as it varies from table to table). But it also may not be fixed for a single table over the lifetime of the connector since the user may execute an `ALTER TABLE` command. The connector must be able to detect these changes and react appropriately. Sink connectors are usually simpler because they are consuming data and therefore do not need to create schemas. However, they should take just as much care to validate that the schemas they receive have the expected format. When the schema does not match -- usually indicating the upstream producer is generating invalid data that cannot be correctly translated to the destination system -- sink connectors should throw an exception to indicate this error to the system.

