

**Table 1: Description of the 34 Spark configuration parameters.**

Configuration Parameters—Description	Range	Default
<b>spark.reducer.maxSizeInFlight</b> —Maximum size of map outputs to fetch simultaneously from each reduce task, in MB.	2–128	48
<b>spark.shuffle.file.buffer</b> —Size of the in-memory buffer for each shuffle file output stream, in KB.	2–128	32
<b>spark.shuffle.sort.bypassMergeThreshold</b> —Avoid merge-sorting data if there is no map-side aggregation.	100–1000	200
<b>spark.speculation.interval</b> —How often Spark will check for tasks to speculate, in millisecond.	10–100	100
<b>spark.speculation.multiplier</b> —How many times slower a task is than the median to be considered for speculation.	1–5	1.5
<b>spark.speculation.quantile</b> —Percentage of tasks which must be complete before speculation is enabled.	0–1	0.75
<b>spark.broadcast.blockSize</b> —Size of each piece of a block for TorrentBroadcastFactory, in MB.	2–128	4
<b>spark.io.compression.codec</b> —The codec used to compress internal data such as RDD partitions, and so on.	snappy, lz4, lz4	snappy
<b>spark.io.compression.lz4.blockSize</b> —Block size used in LZ4 compression, in KB.	2–128	32
<b>spark.io.compression.snappy.blockSize</b> —Block size used in snappy, in KB.	2–128	32
<b>spark.kryo.referenceTracking</b> —Whether to track references to the same object when serializing data with Kryo	true,false	true
<b>spark.kryoserializer.buffer.max</b> —Maximum allowable size of Kryo serialization buffer, in MB.	8–128	64
<b>spark.kryoserializer.buffer</b> —Initial size of Kryo's serialization buffer, in KB.	2–128	64
<b>spark.driver.cores</b> —Number of cores to use for the driver process.	1–30	1
<b>spark.executor.cores</b> —The number of cores to use on each executor.	4–30	core #
<b>spark.executor.instances</b> —The number of executors for static allocation.	6–10	2
<b>spark.driver.memory</b> —Amount of memory to use for the driver process, in MB.	1024–36864	1024
<b>spark.executor.memory</b> —Amount of memory to use per executor process, in MB.	7168–36864	1024
<b>spark.storage.memoryMapThreshold</b> —Size of a block above which Spark maps when reading a block from disk, in MB.	50–500	2
<b>spark.network.timeout</b> —Default timeout for all network interactions, in second.	20–500	120
<b>spark.locality.wait</b> —How long to launch a data-local task before giving up, in second.	1–10	3
<b>spark.scheduler.revive.interval</b> —The interval length for the scheduler to revive the worker resource, in second.	2–50	1
<b>spark.task.maxFailures</b> —Number of task failures before giving up on the job.	1–8	4
<b>spark.shuffle.compress</b> —Whether to compress map output files.	true,false	true
<b>spark.memory.fraction</b> —Fraction of (heap space - 300 MB) used for execution and storage.	0.5–1	0.75
<b>spark.shuffle.spill.compress</b> —Whether to compress data spilled during shuffles.	true,false	true
<b>spark.speculation</b> —If set to "true", performs speculative execution of tasks.	true,false	false
<b>spark.broadcast.compress</b> —Whether to compress broadcast variables before sending them. Generally a good idea.	true,false	true
<b>spark.rdd.compress</b> —Whether to compress serialized RDD partitions.	true,false	false
<b>spark.serializer</b> —Class to use for serializing objects that are sent over the network or need to be cached in serialized form.	java,kryo	java
<b>spark.memory.storageFraction</b> —Amount of storage memory immune to eviction, expressed as a fraction of the size of the region set aside by <i>spark.memory.fraction</i> .	0.5–1	0.5
<b>spark.default.parallelism</b> —The largest number of partitions in a parent RDD for distributed shuffle operations.	8–50	#
<b>spark.memory.offHeap.enabled</b> —If true, Spark will attempt to use off-heap memory for certain operations.	true,false	false
<b>spark.memory.offHeap.size</b> —The absolute amount of memory which can be used for off-heap allocation, in MB.	10–1000	0