

Let's dumb-proof data pipelines

Ganesh Chand
Staff Solutions Consultant, Databricks

Why make your data pipelines dumb-proof?

User error can cause production outages

- Missed SLAs
- Downstream impact
- Business impacts
- Data Decay

Common causes of ETL accidents due to human error

- Configuration error
- Resource misuse
- Operator error
- Defects

Typical Situations

- Environment changes
- Emergency fixes
- Regular maintenance



Secure your loads and double check!

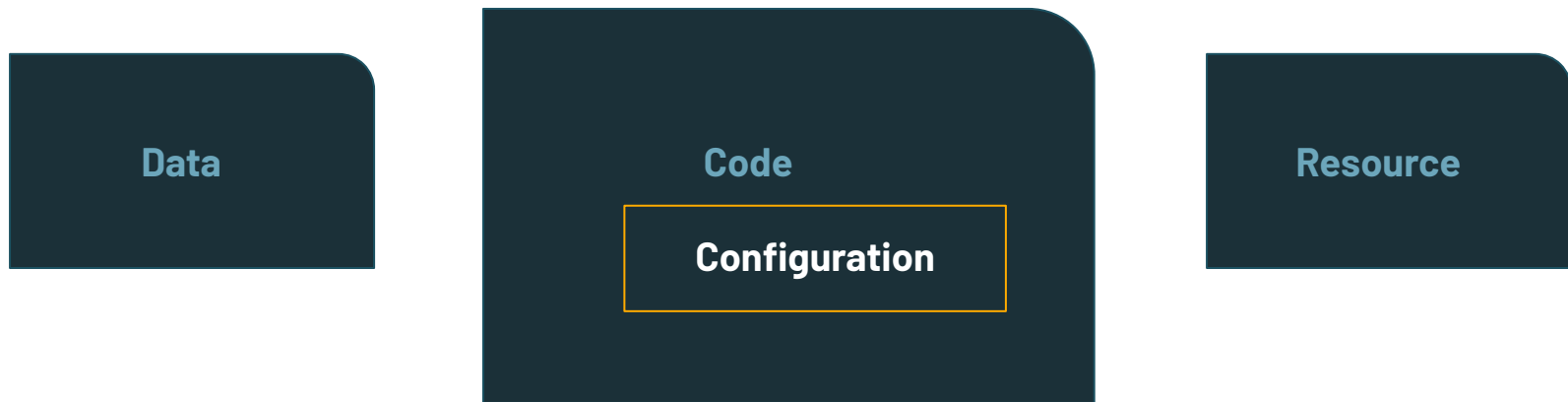
How to make your data pipelines dumb-proof?

Data

Code

Resource

How to make your data pipelines dumb-proof?



How to make your data pipelines dumb-proof?

0

Hard coded Data Pipeline

1

Parameterized Data Pipeline

2

Parameterized Data Pipeline with **externalized** configuration parameters

3

Parameterized Data Pipeline with **boilerplate free** configuration loading code and **validation**

beginner

intermediate

advanced

expert

Hard coded Data Pipelines

```
val input = spark.readStream
  .format("rate")
  .option("rowsPerSecond", 1)
  .load()
```

```
spark.conf.set("spark.sql.shuffle.partitions", 2)
spark.conf.set("spark.databricks.delta.properties.defaults.appendOnly", "true")
```

```
val output = input
// add audit columns
.withColumn("etl_insert_timestamp", current_timestamp())
.withColumn("event_date", to_date('timestamp'))
```

```
output.writeStream.format("delta")
  .option("mergeSchema", "true")
  .option("checkpointLocation", "/tmp/deltalake/dumbproof/bronze/rate/0/_checkpoint")
  .partitionBy("event_date")
  .outputMode("append")
  .trigger(ProcessingTime("30 seconds"))
  .start("/tmp/deltalake/dumbproof/bronze/rate/0")
```



Environment changes
User error

Fixing Hard coded Data Pipelines

- **Parameterize** values that change across environments
- **Validate** user inputs
- Refactor code

Parameters & Input Validation

Checkpoint Path : Input Rate (rows per second) : Make table Apped only : | v

Output Table Path : Spark SQL Shuffle Partitions : Trigger Internal :

```
2 // provide reasonal defaults where possible
3 dbutils.widgets.text("rowsPerSecond", "100", "Input Rate (rows per second)") // optional
4 dbutils.widgets.text("outputPath", "", "Output Table Path") // required
5 dbutils.widgets.text("checkpointLocation", "", "Checkpoint Path") // optional
6 dbutils.widgets.text("triggerInterval", "30 seconds", "Trigger Internal") // optional
7 dbutils.widgets.text("shufflePartitions", "10", "Spark SQL Shuffle Partitions") // optional
8 dbutils.widgets.dropdown("isTableAppendOnly", "false", Array("true", "false"), "Make table Apped only") // optional
```

Command took 0.47 seconds -- by ganesh@databricks.com at 5/2/2021, 11:28:01 PM on Shared Autoscaling

Cmd 3

```
1 val rowsPerSecond: String = dbutils.widgets.get("rowsPerSecond")
2 val outputPath: String = dbutils.widgets.get("outputPath")
3 val checkpointLocation: String = {
4   val checkpointLocationInput = dbutils.widgets.get("checkpointLocation")
5   if(checkpointLocationInput.isEmpty) s"$outputPath/_checkpoint"
6   else checkpointLocationInput
7 }
8 val triggerIntervalInput: String = dbutils.widgets.get("triggerInterval")
9 val triggerInterval: Option[Duration] = Try(Duration(triggerIntervalInput)).toOption
10 val shufflePartitions = dbutils.widgets.get("shufflePartitions")
11 val isTableAppendOnly = dbutils.widgets.get("isTableAppendOnly")
12
13 // validation here for instant feedback on change in user input
14 require(triggerInterval.isDefined, s"unable to parse triggerInterval '$triggerIntervalInput' to Duration")
15 require(rowsPerSecond.toInt > 0, s"input rate value must be greater than 0. Got $rowsPerSecond")
16 require(!outputPath.isEmpty, "output path is a required parameter")
```

Define Parameters

Process User Inputs

Validate & Fail on
invalid input

Parameterized Data Pipeline

```
val input = spark.readStream
  .format("rate")
  .option("rowsPerSecond", rowsPerSecond)
  .load()
```

```
spark.conf.set("spark.sql.shuffle.partitions", shufflePartitions)
spark.conf.set("spark.databricks.delta.properties.defaults.appendOnly", isTableAppendOnly)
```

```
val output = input
  // add audit columns
  .withColumn("etl_insert_timestamp", current_timestamp())
  .withColumn("event_date", to_date('timestamp'))
```

```
val writer = output.writeStream.format("delta")
  .queryName("ingest_rateSource_into_delta")
  .option("mergeSchema", "true")
  .option("checkpointLocation", checkpointLocation)
  .partitionBy("event_date")
  .outputMode("append")
  .trigger(ProcessingTime(triggerInterval.get))
  .start(outputPath)
```

```
2 // provide reasonable defaults where possible
3 dbutils.widgets.text("rowsPerSecond", "100", "Input Rate (rows per second)") // optional
4 dbutils.widgets.text("outputPath", "", "Output Table Path") // required
5 dbutils.widgets.text("checkpointLocation", "", "Checkpoint Path") // optional
6 dbutils.widgets.text("triggerInterval", "30 seconds", "Trigger Interval") // optional
7 dbutils.widgets.text("shufflePartitions", "10", "Spark SQL Shuffle Partitions") // optional
8 dbutils.widgets.dropdown("isTableAppendOnly", "false", Array("true", "false"), "Make table Apped only") // optional

Command took 0.47 seconds -- by ganeshdatabricks.com at 5/2/2021, 11:28:01 PM on Shared Autoscaling

Cmd 3

1 val rowsPerSecond: String = dbutils.widgets.get("rowsPerSecond")
2 val outputPath: String = dbutils.widgets.get("outputPath")
3 val checkpointLocation: String = {
4   val checkpointLocationInput = dbutils.widgets.get("checkpointLocation")
5   if (checkpointLocationInput.isEmpty) s"$outputPath/_checkpoint"
6   else checkpointLocationInput
7 }
8 val triggerIntervalInput: String = dbutils.widgets.get("triggerInterval")
9 val triggerInterval: Option[Duration] = Try(Duration(triggerIntervalInput)).toOption
10 val shufflePartitions = dbutils.widgets.get("shufflePartitions")
11 val isTableAppendOnly = dbutils.widgets.get("isTableAppendOnly")
12
13 // validation here for instant feedback on change in user input
14 require(triggerInterval.isDefined, s"unable to parse triggerInterval '$triggerIntervalInput' to Duration")
15 require(rowsPerSecond.toInt > 0, s"input rate value must be greater than 0. Got $rowsPerSecond")
16 require(!outputPath.isEmpty, "output path is a required parameter")
```

✓ Environment changes
✓ User error

- ❑ Input parameters > 10
- ❑ Type casting parameter value
- ❑ Tracking parameter changes

Externalizing Configuration

- Choose configuration Format
- Configuration Library
- Refactor Code

Configuration Format

Many Options

- ☐ ENV VARIABLE
- ☐ INI
- ☐ XML
- ☐ JSON
- ☐ PROPERTIES
- ☐ HOCCON
- ☐ YAML
- ☐ TOML

Why not JSON (JavaScript Object Notation)

JSON isn't a Good Configuration Language

- Verbose
- Lack of comments
- Lack of multiline support
- Lack of substitutions and includes

Configuration Format

Many Options

- ☐ ENV VARIABLE
- ☐ INI
- ☐ XML
- ☐ JSON
- ☐ PROPERTIES
- ☒ **HOCON**
- ☐ YAML
- ☐ TOML

Why HOCON (Human-Optimized Config Object Notation)

HOCON is a purpose-built Configuration format

- Concise
- Supports comments
- Supports multiline String
- Supports substitutions and includes
- a JSON superset
- Python/Java/Scala libraries

Official HOCON Spec: <https://github.com/lightbend/config/blob/master/HOCON.md>

Configuration in **JSON** Format

```
{
  "dev": {
    "outputPath": "/tmp/deltalake/dev/sample/1/rate",
    "checkpointPath": "/tmp/deltalake/dev/sample/1/rate/_checkpoint",
    "rowsPerSecond": 1,
    "triggerInterval": "30s",
    "tuning": [
      "spark.sql.shuffle.partitions=10"
    ]
  },
  "prod": {
    "outputPath": "/tmp/deltalake/prod/sample/1/rate",
    "checkpointPath": "/tmp/deltalake/prod/sample/1/rate/_checkpoint",
    "rowsPerSecond": 1000,
    "triggerInterval": "1m",
    "tuning": [
      "spark.sql.shuffle.partitions=10"
    ]
  }
}
```

Configuration in **HOCON** format

```
// iteration 1: removed unwanted curly braces, colons and double quotes

dev {
  outputPath: /tmp/deltalake/dev/sample/1/rate
  checkpointPath: /tmp/deltalake/dev/sample/1/rate/_checkpoint
  rowsPerSecond: 1,
  triggerInterval: 30s,
  tuning: ["spark.sql.shuffle.partitions=10"]
},

prod {
  outputPath: /tmp/deltalake/prod/sample/1/rate
  checkpointPath: /tmp/deltalake/prod/sample/1/rate/_checkpoint
  rowsPerSecond: 1000,
  triggerInterval: 1m,
  tuning: ["spark.sql.shuffle.partitions=10"]
}
```

Optimized Configuration in **HOCON** format

```
// iteration 1: removed unwanted curly braces, colons and double quotes
// iteration 2: extracted common values to common config object and checkpoint path is set to relative to output
path

common {
  rootPath = /tmp/deltalake
  tuning = ["spark.sql.shuffle.partitions=10"]
}

dev {
  outputPath: ${common.rootPath}/dev/sample/1/rate,
  checkpointPath: ${dev.outputPath}_checkpoint,
  rowsPerSecond: 1,
  triggerInterval: 30s,
  tuning: ${common.tuning}
},

prod {
  outputPath: ${common.rootPath}/prod/sample/1/rate,
  checkpointPath: ${prod.outputPath}/_checkpoint,
  rowsPerSecond: 1000,
  triggerInterval: 1m,
  tuning: ${common.tuning}
}
```


Readable and maintainable Configuration

```
{
  "dev": {
    "outputPath": "/tmp/deltalake/dev/sample/1/rate",
    "checkpointPath": "/tmp/deltalake/dev/sample/1/rate/_checkpoint",
    "rowsPerSecond": 1,
    "triggerInterval": "30s",
    "tuning": [
      "spark.sql.shuffle.partitions=10"
    ]
  },
  "prod": {
    "outputPath": "/tmp/deltalake/prod/sample/1/rate",
    "checkpointPath": "/tmp/deltalake/prod/sample/1/rate/_checkpoint",
    "rowsPerSecond": 1000,
    "triggerInterval": "1m",
    "tuning": [
      "spark.sql.shuffle.partitions=10"
    ]
  }
}
```



```
common {
  rootPath = /tmp/deltalake
  tuning = ["spark.sql.shuffle.partitions=10"]
}

dev {
  outputPath: ${common.rootPath}/dev/sample/1/rate,
  checkpointPath: ${dev.outputPath}_checkpoint,
  rowsPerSecond: 1,
  triggerInterval: 30s,
  tuning: ${common.tuning}
},

prod {
  outputPath: ${common.rootPath}/prod/sample/1/rate,
  checkpointPath: ${prod.outputPath}_checkpoint,
  rowsPerSecond: 1000,
  triggerInterval: 1m,
  tuning: ${common.tuning}
}
```

Configuration Library

Typesafe - Configuration library for JVM languages by Lightbend

- Implemented in plain Java with **zero** dependencies
- Supports files in three formats: Java properties, JSON, and a human-friendly JSON superset (HOCON)
- Included in Databricks Runtime

```
import com.typesafe.config._

val conf: Config =
  ConfigFactory.parseString(dbutils.fs.head("dbfs:/path/to/data-pipeline.conf")).resolve

val rowsPerSecond: Long = conf.getLong("rowsPerSecond")
val outputPath: String = conf.getString("outputPath")
val checkpointLocation: String = conf.getString("checkpointPath")
val triggerInterval: ScalaDuration = ScalaDuration(conf.getString("triggerInterval"))
val sparkTuningConfs: Seq[String] = conf.getStringList("tuning").asScala
```

Refactor Code - Loading and Parsing Configuration

Environment : `dev`

Pipeline Config file : `/dumbproof/deltalake/c`

```
1 dbutils.widgets.removeAll()
2 dbutils.widgets.text("env", "dev", "Environment")
3 dbutils.widgets.text("configFilePath", "", "Pipeline Config file")
4
5 val env = dbutils.widgets.get("env")
6 val configFilePath = dbutils.widgets.get("configFilePath")
7
8
9 import com.typesafe.config._
10 val conf = ConfigFactory.parseString(dbutils.fs.head(configFilePath)).resolve.getConfig(env)
11
12
13 val rowsPerSecond: Long = conf.getLong("rowsPerSecond")
14 val outputPath: String = conf.getString("outputPath")
15 val checkpointLocation: String = conf.getString("checkpointPath")
16 val triggerInterval: ScalaDuration = ScalaDuration(conf.getString("triggerInterval"))
17 val sparkTuningConfs: Seq[String] = conf.getStringList("tuning").asScala
18 require(rowsPerSecond.toInt > 0, s"input rate value must be greater than 0. Got $rowsPerSecond")
```

Parameterized Data Pipeline

```
val input = spark.readStream
  .format("rate")
  .option("rowsPerSecond", rowsPerSecond)
  .load()

sparkTuningConfs
  .foreach(config => spark.sql($"SET $config"))

val output = input
  // add audit columns
  .withColumn("etl_insert_timestamp", current_timestamp)
  .withColumn("event_date", to_date('timestamp))

val writer = output.writeStream.format("delta")
  .queryName("ingest_rateSource_into_delta")
  .option("mergeSchema", "true")
  .option("checkpointLocation", checkpointLocation)
  .partitionBy("event_date")
  .outputMode("append")
  .trigger(ProcessingTime(triggerInterval))
  .start(outputPath)
```

Environment: dev

Pipeline Config file: /dumbproof/deltalake/c

```
1 dbutils.widgets.removeAll()
2 dbutils.widgets.text("env", "dev", "Environment")
3 dbutils.widgets.text("configFilePath", "", "Pipeline Config file")
4
5 val env = dbutils.widgets.get("env")
6 val configFilePath = dbutils.widgets.get("configFilePath")
7
8
9 import com.typesafe.config._
10 val conf = ConfigFactory.parseString(dbutils.fs.head(configFilePath)).resolve.getConfig(env)
11
12
13 val rowsPerSecond: Long = conf.getLong("rowsPerSecond")
14 val outputPath: String = conf.getString("outputPath")
15 val checkpointLocation: String = conf.getString("checkpointPath")
16 val triggerInterval: ScalaDuration = ScalaDuration(conf.getString("triggerInterval"))
17 val sparkTuningConfs: Seq[String] = conf.getStringList("tuning").asScala
18 require(rowsPerSecond.toInt > 0, s"input rate value must be greater than 0. Got $rowsPerSecond")
```

- ✓ Environment changes
- ✓ User error
- ✓ Input parameters > 10
- ✓ Type casting parameter value
- ✓ Tracking parameter changes

- ❑ Reduce boilerplate code
- ❑ Enhanced validation

Boilerplate free configuration code

PureConfig - Scala library for loading configuration files

- Reads HOCON, .properties, JSON to native Scala classes in a boilerplate-free way
- Supports sealed traits, case classes, collections, optional values out-of-the-box
- Config Writers to generate configuration files

Official Docs: <https://pureconfig.github.io/>

Sample Code

```
9  import scala.concurrent.duration.{Duration => ScalaDuration}
10 import org.apache.spark.sql.streaming.Trigger
11
12 sealed trait StreamingTrigger
13 case object TriggerOnce extends StreamingTrigger
14 case class TriggerWithInterval(interval: ScalaDuration) extends StreamingTrigger
15
16 case class PipelineConf(
17     outputPath: String,
18     checkpointPath: String,
19     rowsPerSecond: Long,
20     streamingTrigger: StreamingTrigger,
21     tuning: Seq[String]
22 ) {
23     require(rowsPerSecond > 0, s"input rate value must be greater than 0. Got $rowsPerSecond")
24     require(!outputPath.isEmpty, "table output path cannot be empty")
25     require(
26         tuning.forall { sparkConfString =>
27             sparkConfString.split("=").size == 2
28         }, "" "Invalid spark configuration found. A valid spark configuration must be of type "<key>=<value>""
29     )
30     val trigger: Trigger = streamingTrigger match {
31         case TriggerOnce => Trigger.Once()
32         case TriggerWithInterval(duration) => Trigger.ProcessingTime(duration)
33     }
34 }
35
36 import pureconfig._
37 import pureconfig.generic.auto._
38
39 val configStr: String =
40     """
41     | output-path = /table/output/path/,
42     | checkpoint-path = /table/checkpoint/path,
43     | rows-per-second = 1,
44     | streaming-trigger { type = "trigger-once" }
45     | tuning = ["spark.sql.shuffle.partitions=10"]
46     """
47     .stripMargin
48
49 val pipelineConf: PipelineConf = ConfigSource.string(configStr).load[PipelineConf].right.get
50
51 output
52     .writeStream
53     .trigger(pipelineConf.trigger) // pipelineConf.trigger: org.apache.spark.sql.streaming.Trigger = OneTimeTrigger
54     .format("delta")
55     .option("checkpointLocation", pipelineConf.checkpointPath)
56     .start(pipelineConf.outputPath)
```

- ✓ Environment changes
- ✓ User error
- ✓ Input parameters > 10
- ✓ Type casting parameter value
- ✓ Tracking parameter changes
- ✓ Reduce boilerplate code
- ✓ Enhanced validation

Summary

- User errors are inevitable. Make it hard, if not impossible, to make mistakes
- Catch errors and fail as early as you can
 - Compile Time > Unit Test > Integration Test > Runtime (before Spark Job starts)
- Externalize configuration to track changes
- Use libraries to reduce boilerplate

Source Code and examples available at

<https://github.com/ganeshchand/lets-dumb-proof-data-pipeline>

Resources

.