



inovex

**The productivity booster for your Data
Projects with Scala 3 & scala-cli**

こんにちは - About me

Kolja Maier

- Data & ML Engineering
@ inovex 🇩🇪
- ~10 years of Scala
experience



inovex

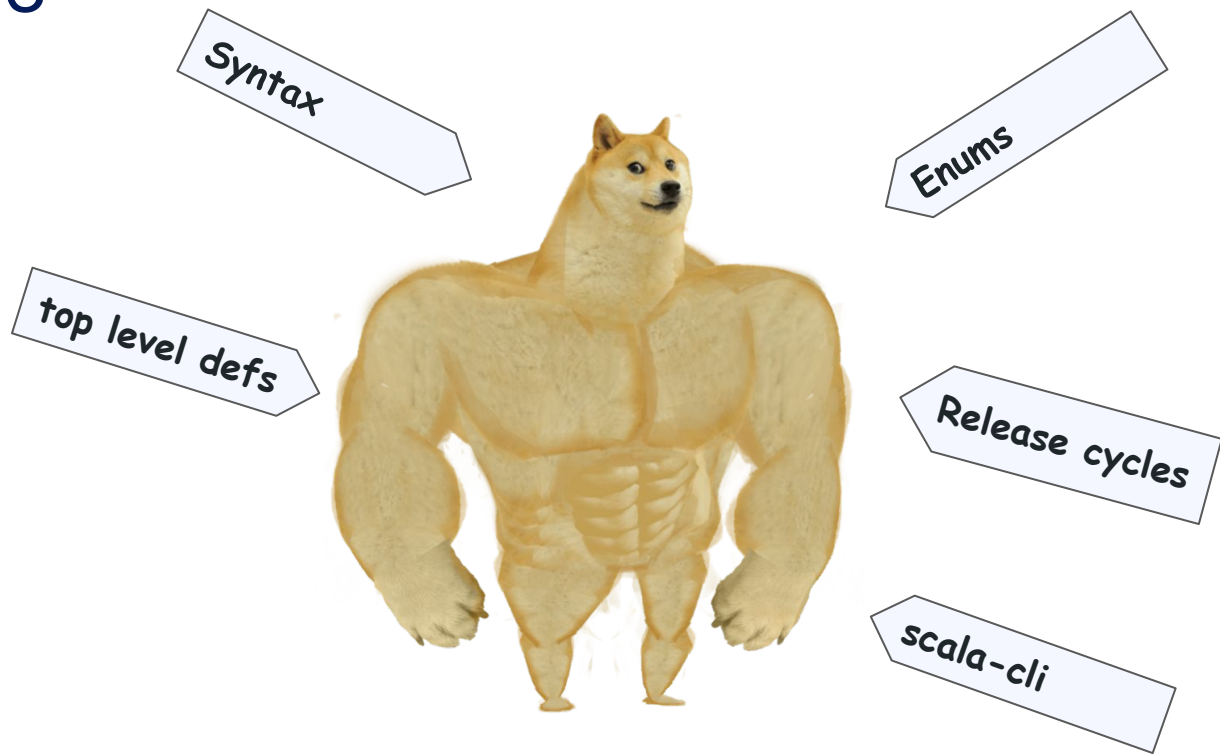
Big Data

Battle proof architecture



Productivity

Scala 3



scala-cli examples

Spark



```
1 scala-cli run --spark SparkExample.scala
```

```
people
  .toDS
  .group
    ⚙️ groupBy(cols: Column*): RelationalGroup... (cols: Column*): RelationalGrou ×
    ⚙️ groupBy(col1: String, cols: String*): R... pedDataset
    ⚙️ groupByKey[K](func: MapFunction[Person,...
    ⚙️ groupByKey[K: Encoder](func: Person => ...

Groups the Dataset using the specified
columns, so we can run aggregation on
them. See RelationalGroupedDataset for all
the available aggregate functions.

// Compute the average for all num
ds.groupBy($"department").avg()

// Compute the max age and average
ds.groupBy($"department", $"gender
  "salary" -> "avg",
  "age" -> "max"
))
```


Spark

```
1 people.toDS
2   .groupByKey(_.name)
3   .mapGroups { (name, peopleWithSameName) =>
```

	_1	_2	_3	_4
Marukami	49	75	63.0	
Mishima	25	25	25.0	
Suzuki	16	41	28.5	

Spark

```
1 people.toDS
2   .groupByKey(_.name)
3   .mapGroups { (name, peopleWithSameName) =>
4     val ages = peopleWithSameName.map(_._age).toList
5     (
6       name,
7       ages.min,
8       ages.max,
9       ages.sum.toDouble / ages.length
10    )
11  }
```

	_1	_2	_3	_4
Marukami	49	75	63.0	
Mishima	25	25	25.0	
Suzuki	16	41	28.5	

Spark

```
1 people.toDS
2   .groupByKey(_.name)
3   .mapGroups { (name, peopleWithSameName) =>
4     val ages = peopleWithSameName.map(_._age).toList
5     (
6       name,
7       ages.min,
8       ages.max,
9       ages.sum.toDouble / ages.length
10    )
11  }
12  .map(statistics =>
13    (statistics._1, statistics._2, statistics._3, statistics._4)
14  )
15  .show()
```

	_1	_2	_3	_4
Marukami	49	75	63.0	
Mishima	25	25	25.0	
Suzuki	16	41	28.5	

Spark



```
1 people.toDS.createOrReplaceTempView("people")
2
3 val sqlDF =
4   spark.sql(
5     """
6     SELECT name, COUNT(*) AS cnt
7     FROM people
8     GROUP BY name
9     """
10  )
11 sqlDF.show()
```

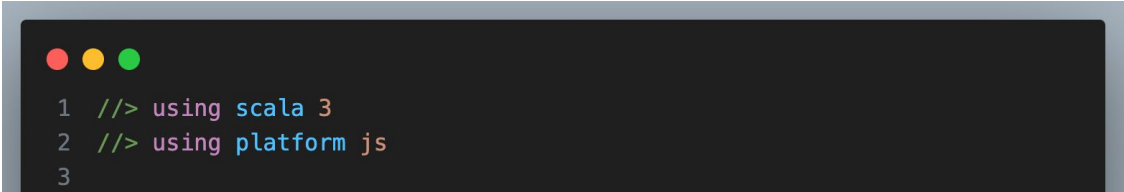
name	cnt
Suzuki	2
Mishima	1
Marukami	3

Spark



```
1 //> using scala 3
2 //> using dep org.apache.spark:spark-sql-api_2.13:3.5.1
3 //> using dep io.github.vincenzobaz::spark-scala3-encoders:0.2.6
4
5 import org.apache.spark.sql.SparkSession
6 import scala3encoders.given
```

BigQuery UDFs



```
1 //> using scala 3
2 //> using platform js
3
```

BigQuery UDFs



```
1 //> using scala 3
2 //> using platform js
3
4 import scala.scalajs.js
5 import scala.scalajs.js.annotation.*
6
```

BigQuery UDFs

```
1 //> using scala 3
2 //> using platform js
3
4 import scala.scalajs.js
5 import scala.scalajs.js.annotation.*
6
7 enum Flavor:
8   case 🥬, 🐔, 🐄, 🐷, 🍌, 🐘
```


BigQuery UDFs

```
1  //> using scala 3
2  //> using platform js
3
4  import scala.scalajs.js
5  import scala.scalajs.js.annotation.*
6
7  enum Flavor:
8    case 🥬, 🐔, 🐷, 🍣, 🐖
9
10 @JSImportTopLevel("dailyFlavor")
11 def dailyFlavor(b: String): String =
12   Flavor.valueOf(b) match
13     case Flavor.🥬 => "Veggie day"
14     case Flavor.🍣 => "Seafood lover!"
15     case _ => "...and so on"
```

BigQuery UDFs

```
CREATE TEMP FUNCTION dailyFlavor(f STRING)  
RETURNS STRING
```

BigQuery UDFs

```
CREATE TEMP FUNCTION dailyFlavor(f STRING)  
RETURNS STRING  
LANGUAGE js  
OPTIONS (  
  library=[ 'gs://$GCS_BUCKET/main.js' ])
```

BigQuery UDFs

```
CREATE TEMP FUNCTION dailyFlavor(f STRING)
RETURNS STRING
LANGUAGE js
OPTIONS (
  library=[ 'gs://$GCS_BUCKET/main.js' ])
AS r"""
  return dailyFlavor(f)
""";
```

BigQuery UDFs

```
CREATE TEMP FUNCTION dailyFlavor(f STRING)
RETURNS STRING
LANGUAGE js
OPTIONS (
  library=['gs://$GCS_BUCKET/main.js'])
AS r"""
  return dailyFlavor(f)
""";
SELECT dailyFlavor('🥬') AS result;
```

```
+-----+
|  result  |
+-----+
| Veggie day |
+-----+
```

Kafka - Data modelling

```
1 case class Ramen(brand: String, flavor: Flavor, price: Double)
```

```
1 enum Flavor:  
2   case 🌿, 🐔, 🐘, 🍜, 🐷
```

Kafka - Data modelling

```
1 case class Ramen(brand: String, flavor: Flavor, price: Double)
2 object Ramen:
3   given JsonCodec[Ramen] = DeriveJsonCodec.gen[Ramen]
```


```
1 enum Flavor:
2   case 🍃, 🐔, 🍷, 🍜, 🍝, 🍖
3 object Flavor:
4   given JsonCodec[Flavor] =
5     JsonCodec[String].transform(Flavor.valueOf, _.toString)
```

Kafka - Data modelling

```
1 case class Ramen(brand: String, flavor: Flavor, price: Double)
2 object Ramen:
3   given JsonCodec[Ramen] = DeriveJsonCodec.gen[Ramen]
4   val valueSerializer = Serde.string.inmapM[Any, Ramen](s =>
5     ZIO
6       .fromEither(s.fromJson[Ramen])
7       .mapError(e => new RuntimeException(e))
8   )(r => ZIO.succeed(r.toJson))
```

```
1 enum Flavor:
2   case 🍜, 🐔, 🍱, 🍲, 🍤, 🍖
3 object Flavor:
4   given JsonCodec[Flavor] =
5     JsonCodec[String].transform(Flavor.valueOf, _.toString)
```


Kafka - Producer



```
1 ZStream
2   .fromZIO(Random.nextIntBetween(0, Int.MaxValue))
3   .forever
```

Kafka - Producer

```
1 ZStream
2   .fromZIO(Random.nextIntBetween(0, Int.MaxValue))
3   .forever
4   .mapZIO { id =>
5     Producer.produce[Any, Int, Ramen](
6       topic = "ramen",
7       key = id,
```

Kafka - Producer

```
1 ZStream
2   .fromZIO(Random.nextIntBetween(0, Int.MaxValue))
3   .forever
4   .mapZIO { id =>
5     Producer.produce[Any, Int, Ramen](
6       topic = "ramen",
7       key = id,
8       value = Ramen(
9         SRandom.shuffle(List("Nissin", "Ichiran", "Samyang")).head,
10        Flavor.fromOrdinal(id % 5),
11        SRandom.between(0.0, 1.0)
12      ),
```

Kafka - Producer

```
1 ZStream
2   .fromZIO(Random.nextIntBetween(0, Int.MaxValue))
3   .forever
4   .mapZIO { id =>
5     Producer.produce[Any, Int, Ramen](
6       topic = "ramen",
7       key = id,
8       value = Ramen(
9         SRandom.shuffle(List("Nissin", "Ichiran", "Samyang")).head,
10        Flavor.fromOrdinal(id % 5),
11        SRandom.between(0.0, 1.0)
12      ),
13       keySerializer = Serde.int,
14       Ramen.valueSerializer
15     )
16   }
```

```
Received Ramen(Ichiran, 🐔, 0.022659367901554872)
Sent value ramen-0@32776
Received Ramen(Ichiran, 🌿, 0.2837352181068149)
Sent value ramen-0@32777
^CSent value ramen-0@32778
Sent value ramen-0@32779
Sent value ramen-0@32780
Received Ramen(Samyang, 🐷, 0.15919173807975895)
Sent value ramen-0@32781
Sent value ramen-0@32782
Received Ramen(Ichiran, 🐷, 0.3246139886443965)
Sent value ramen-0@32783
Received Ramen(Nissin, 🐔, 0.14962751137279784)
Sent value ramen-0@32784
Received Ramen(Samyang, 🍜, 0.09595954071228774)
```

	Offset	Partition
-	0	0

KeyValueHeaders

{

"brand": "Samyang",
"flavor": "🐔",
"price": 0.8224032060

}

Summary

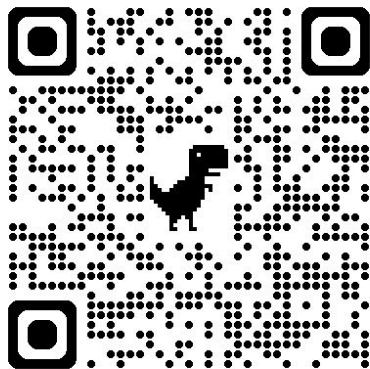
Summary

- It is crucial to have a setup for validating data tooling with fast evolving business requirements
 - Scala 3 powerful features and scala-cli represent a strong combination!
 - Disclaimer: Use what fits best
- Further productivity pointers
 - scala-cli toolkit
 - just

Thank you!

Kolja Maier

[Code to this talk](#)



inovex

