

## GSoC Final Report

# Spark Backend Support for Gora (GORA-386) Final Report

Furkan KAMACI<sup>1</sup>

<sup>1</sup> *Istanbul Technical University*

**Abstract**—This reports explains the overall progress for GSoC 2015.

**Index Terms**— Apache Gora, GSoC, Spark

## I. INTRODUCTION

GORA-386 aims to develop a Spark Backend for Apache Gora. Apache Gora open source framework provides an in-memory data model and persistence for big data. Gora supports persisting to column stores, key value stores, document stores and RDBMSs, and analyzing the data with extensive Apache Hadoop MapReduce support. Even Spark is so powerful compared to Map/Reduce which Gora currently supports; there was no Spark backend for Gora and my GSoC project is aimed that.

During GSoC period, I've blogged at my personal website (<http://furkankamaci.com/>) and created a fork from Apache Gora's master branch and worked on it: <https://github.com/kamaci/gora>

This report explains what is done during my GSoC period. Here is the tag of this report (you can check for whole list from here: <https://cwiki.apache.org/confluence/display/GORA/Spark+Backend+Support+for+Gora+%28GORA-386%29+Reports>) :

<b>Project Name</b>	Spark Backend Support for Gora (GORA-386)
<b>Project URL</b>	<a href="https://issues.apache.org/jira/browse/GORA-386">https://issues.apache.org/jira/browse/GORA-386</a>
<b>Report #</b>	9
<b>Report Compiled By</b>	Furkan KAMACI
<b>Report Date</b>	01 Sep 2015

## II. PROPOSED SCHEDULE&TIMELINE

Suggested schedule and timeline is as follows:

- 1) Analyzing The Problem (1 Week)
  - a) Problem will be analyzed with more detail.
- 2) Spark - Map/Reduce Translation (4 weeks)
  - a) Behaviors specific to Hadoop's implementation rather than the idea of MapReduce in the abstract will be recreated according to Spark.
- 3) Gora – Translated Map/Reduce Integration (4 weeks)
  - a) Gora has its own input/output format and Spark uses RDD. After necessary implementation is done for Spark – Map/Reduce translation, Gora Integration will be done. (2 weeks).
- 4) Generic Abstraction Layer Backend (1 week)
  - a) A generic layer backend that let other projects' integration as like Spark.
- 5) Test (1 week)
  - a) Implementation tests will be written and run.
- 6) Documentation (1 week)
  - a) Documentation will be prepared.

### III. COMMUNITY BONDING

At community bonding period, Apache Gora documentation is read and Apache Gora source code is checked to be more familiar with project. Related projects are analyzed including Apache Flink and Apache Crunch to implement a Spark backend into Apache Gora.

An issue from Jira is picked up (<https://issues.apache.org/jira/browse/GORA-262>) and fixed at Apache Gora. Security is always an issue when writing into Solr and resolved issue GORA-262 is about adding support for HTTPClient authentication in gora-solr module.

On the other hand, Apache Gora documentation is organized and introduced with Apache Gora dev team. I've also blogged at my personal website about my acceptance of GSoC (<http://furkankamaci.com/gsoc-2015-acceptance-for-apache-gora/>)

### IV. PRE MIDTERM PERIOD

I've separated this period into three main parts as followed:

#### Part 1

Due to implementing this project needs an infrastructure about Apache Spark, I've started with analyzing Spark's first papers. I've analyzed "Spark: Cluster Computing with Working" ([http://www.cs.berkeley.edu/~matei/papers/2010/hotcloud\\_spark.pdf](http://www.cs.berkeley.edu/~matei/papers/2010/hotcloud_spark.pdf)) and "Resilient Distributed Datasets: A Fault-Tolerant Abstraction for In-Memory Cluster Computing" ([https://www.cs.berkeley.edu/~matei/papers/2012/nsdi\\_spark.pdf](https://www.cs.berkeley.edu/~matei/papers/2012/nsdi_spark.pdf)). I've published two posts about Spark and Cluster Computing (<http://furkankamaci.com/spark-and-cluster-computing/>) and Resilient Distributed Datasets (<http://furkankamaci.com/resilient-distributed-datasets-rdds/>) at my personal blog.

I've followed Apache Spark documentation and developed examples to analyze RDDs.

#### Part 2

Apache Gora's GoraInputFormat class and Spark's newHadoopRDD method are analyzed. An example application implemented to read data from Hbase. Here is the most important line of code for Apache Gora and Spark integration:

```
JavaPairRDD<ImmutableBytesWritable, Result> hbaseRDD = sc.newAPIHadoopRDD(conf, TableInputFormat.class, ImmutableBytesWritable.class, Result.class);
```

Apache Gora supports reading/writing data from/to Hadoop files. Spark has a method for generating an RDD compatible with Hadoop files. So, an architecture is designed which creates a bridge between GoraInputFormat and RDD due to both of them support Hadoop files.

#### Part 3

A base class is created for Apache Gora and Spark integration named as: GoraSparkEngine. It has initialize methods that takes Spark context, data store, optional Hadoop configuration and returns an RDD. One of its initialize method's signature and JavaDoc is as follows:

```
/**
 * Initializes a {@link JavaPairRDD} from given Spark context and data store.
 * If given data store is {@link Configurable} and has not a configuration
 * than a Hadoop configuration is created otherwise existed configuration is
 * used.
 *
 * @param sparkContext Spark context
 * @param dataStore Data store
 * @return initialized rdd
 */
public JavaPairRDD<K, V> initialize(JavaSparkContext sparkContext,
    DataStore<K, V> dataStore) {}
```

Apache Gora has a tutorial module and there is an example for map reduce jobs which is named as LogAnalytics. It is documented as follows:

*LogAnalytics is the tutorial class to illustrate Gora MapReduce API. The analytics mapreduce job reads the web access data stored earlier by the LogManager, and calculates the aggregate daily pageviews. The output of the job is stored in a Gora compatible data store.*

After implementing a base for GoraSpark engine, I've developed a new example aligned to LogAnalytics named as: LogAnalyticsSpark. I've developed map and reduce parts (except for writing results into database) which does the same thing as LogAnalytics and also something more i.e. printing number of lines in tables.

Here is the part of how to run GoraSparkEngine:

```
SparkConf sparkConf = new SparkConf().setAppName("Gora Integration Application").setMaster("local");

Class[] c = new Class[1];
c[0] = Pageview.class;
sparkConf.registerKryoClasses(c);

JavaSparkContext sc = new JavaSparkContext(sparkConf);
Configuration hadoopConf = new Configuration();
DataStore<Long, Pageview> dataStore = DataStoreFactory.getDataStore(inStoreClass, Long.class, Pageview.class,
hadoopConf);
JavaPairRDD<Long, Pageview> goraRDD = goraSparkEngine.initialize(sc, dataStore);
```

When we get an RDD, we can do the operations over it as like making operations on any other RDDs which is not created over Apache Gora. Here is the map function used at calculation:

```
/**
 * map function used in calculation
 */
private static Function<Pageview, Tuple2<Tuple2<String, Long>, Long>> mapFunc = new Function<Pageview,
Tuple2<Tuple2<String, Long>, Long>>() {
    @Override
    public Tuple2<Tuple2<String, Long>, Long> call(Pageview pageview)
        throws Exception {
        String url = pageview.getUrl().toString();
        Long day = getDay(pageview.getTimestamp());
        Tuple2<String, Long> keyTuple = new Tuple2<>(url, day);
        return new Tuple2<>(keyTuple, 1L);
    }
};
```

Here is the way that map function is called on an RDD which is created via GoraSparkEngine with GoraInputFormat:

```
JavaRDD<Tuple2<Tuple2<String, Long>, Long>> mappedGoraRdd = goraRDD.values().map(mapFunc);
```

## V. AFTER MIDTERM PERIOD

After reading data via Gora with Spark, I've aimed to write RDDs into different data stores via Apache Gora and did it. Here is how to persist RDD data via Gora:

```
JavaPairRDD<String, MetricDatum> reducedGoraRdd =
JavaPairRDD.fromJavaRDD(mappedGoraRdd).reduceByKey(redFunc).mapToPair(metricFunc);
```

Here is the reduce and metric functions used at calculation:

```
/**
 * reduce function used in calculation
 */
private static Function2<Long, Long, Long> redFunc = new Function2<Long, Long, Long>() {
    @Override
    public Long call(Long aLong, Long aLong2) throws Exception {
        return aLong + aLong2;
    }
};

/**
 * metric function used after map phase
 */
private static PairFunction<Tuple2<Tuple2<String, Long>, Long>, String, MetricDatum> metricFunc = new
PairFunction<Tuple2<Tuple2<String, Long>, Long>, String, MetricDatum>() {
    @Override
    public Tuple2<String, MetricDatum> call(
        Tuple2<Tuple2<String, Long>, Long> tuple2LongTuple2) throws Exception {
```

```
String dimension = tuple2LongTuple2._1()._1();
Long timestamp = tuple2LongTuple2._1()._2();
MetricDatum metricDatum = new MetricDatum();
metricDatum.setMetricDimension(dimension);
metricDatum.setTimestamp(timestamp);
String key = metricDatum.getMetricDimension().toString();
key += "_" + Long.toString(timestamp);
metricDatum.setMetric(tuple2LongTuple2._2());
return new Tuple2<>(key, metricDatum);
}
};
```

I've improved design architecture and GoraSparkEngine class and also improved LogAnalyticsSpark example (which reads data from Hbase and writes result into Solr) and written test classes for my implementation.

I've published a post about Apache Gora: "*In-Memory Data Model and Persistence for Big Data*" (<http://furkankamaci.com/in-memory-data-model-and-persistence-for-big-data/>) and another one for my GSoC project: "*Spark Backend for Apache Gora*" (<http://furkankamaci.com/spark-backend-for-apache-gora/>).

## VI. CONCLUSION

Gora is a powerful project, roughly, which can work like Hibernate of NoSQL world and one can run Map/Reduce jobs on it. Even Spark is so powerful compared to Map/Reduce which Gora currently supports; there was no Spark backend for Gora and I've implemented it with my GSoC project.

One can read data from a data store (i.e. Hbase), run Spark codes on it (map/reduce) and write results into same or another data store (i.e. Solr). GoraSparkEngine gives a Spark backend capability to Apache Gora and I think that it will make Gora much more powerful.

## VII. REFERENCES

<https://issues.apache.org/jira/browse/GORA-386>  
<http://furkankamaci.com/>  
<https://github.com/kamaci/gora>  
<https://issues.apache.org/jira/browse/GORA-262>  
<http://furkankamaci.com/gsoc-2015-acceptance-for-apache-gora/>  
[http://www.cs.berkeley.edu/~matei/papers/2010/hotcloud\\_spark.pdf](http://www.cs.berkeley.edu/~matei/papers/2010/hotcloud_spark.pdf)  
[https://www.cs.berkeley.edu/~matei/papers/2012/nsdi\\_spark.pdf](https://www.cs.berkeley.edu/~matei/papers/2012/nsdi_spark.pdf)  
<http://furkankamaci.com/spark-and-cluster-computing/>  
<http://furkankamaci.com/resilient-distributed-datasets-rdds/>  
<https://cwiki.apache.org/confluence/display/GORA/Spark+Backend+Support+for+Gora+%28GORA-386%29+Reports>  
<http://gora.apache.org/current/tutorial.html#introduction>  
[http://en.wikipedia.org/wiki/Apache\\_Gora](http://en.wikipedia.org/wiki/Apache_Gora)  
<http://furkankamaci.com/in-memory-data-model-and-persistence-for-big-data/>  
<http://furkankamaci.com/spark-backend-for-apache-gora/>