

¿Qué ofrece Autentia Real Business Solutions S.L?

Somos su empresa de **Soporte a Desarrollo Informático**.
 Ese apoyo que siempre quiso tener...

1. Desarrollo de componentes y proyectos a medida



2. Auditoría de código y recomendaciones de mejora

3. Arranque de proyectos basados en nuevas tecnologías

1. Definición de frameworks corporativos.
2. Transferencia de conocimiento de nuevas arquitecturas.
3. Soporte al arranque de proyectos.
4. Auditoría preventiva periódica de calidad.
5. Revisión previa a la certificación de proyectos.
6. Extensión de capacidad de equipos de calidad.
7. Identificación de problemas en producción.



4. Cursos de formación (impartidos por desarrolladores en activo)

Spring MVC, JSF-PrimeFaces /RichFaces,
 HTML5, CSS3, JavaScript-jQuery

Control de autenticación y
 acceso (Spring Security)
 UDDI

JPA-Hibernate, MyBatis
 Motor de búsqueda empresarial (Solr)
 ETL (Talend)

Gestor portales (Liferay)
 Gestor de contenidos (Alfresco)
 Aplicaciones híbridas

Web Services
 Rest Services
 Social SSO
 SSO (Cas)

Dirección de Proyectos Informáticos.
 Metodologías ágiles
 Patrones de diseño
 TDD

Tareas programadas (Quartz)
 Gestor documental (Alfresco)
 Inversión de control (Spring)

BPM (jBPM o Bonita)
 Generación de informes (JasperReport)
 ESB (Open ESB)



» Estás en: [Inicio](#) » [Tutoriales](#) » [Trident, un compañero de viaje para tratar con Storm](#)



Juan Alonso Ramos

Consultor tecnológico de desarrollo de proyectos informáticos.

Ingeniero en Informática, especialidad en Ingeniería del Software

Puedes encontrarme en Autentia: Ofrecemos de servicios soporte a desarrollo, factoría y formación

Somos expertos en Java/J2EE



[Ver todos los tutoriales del autor](#)



Catálogo de servicios Autentia



Fecha de publicación del tutorial: 2014-10-01

Tutorial visitado 9 veces [Descargar en PDF](#)

Trident, un compañero de viaje para tratar con Storm

0. Índice de contenidos.

- 1. Introducción.
- 2. Entorno.
- 3. Introducción a Trident.
- 3.1 Function.
- 3.2 Filter.
- 3.3 Aggregator.
- 3.4 Métodos para el procesamiento de streams.
- 3.5 State.
- 4. Implementar un trending topics de Twitter con Trident.
- 5. Conclusiones.

1. Introducción.

En un tutorial anterior vimos una introducción a Apache Storm como un sistema para procesamiento de streams en tiempo real. En este tutorial vamos a ver el API Trident, una abstracción por encima de Storm para facilitar la creación de las topologías. Se trata de una aproximación a lo que es Pig o Cascading para Hadoop, salvando las distancias.

Con Trident podemos configurar una topología que procese una entrada de dato asemejándola a una query SQL, es decir a partir de una fuente de datos podemos manipularla seleccionando y procesando los datos que necesitemos y finalmente persistimos los resultados en alguna unidad de almacenamiento ya sea el HDFS o una base de datos NoSQL.

Puedes descargar el código del tutorial desde [aquí](#).

2. Entorno.

El tutorial se ha realizado con el siguiente entorno:

- MacBook Pro 15' (2.4 GHz Intel Core i5, 8GB DDR3 SDRAM).
- Oracle Java SDK 1.7.0_60
- Apache Storm 0.9.2-incubating
- Twitter4j 4.0.1

3. Introducción a Trident

Trident es una API para manejar streams de datos. Sobre estos datos podemos realizar funciones de extracción, formateo, agrupamiento, sumas, cuentas, etc. para finalmente persistir el resultado en algún sistema de almacenamiento. A continuación explicamos las funciones más importantes:

3.1 Function

En programación una función recibe una serie de datos en la entrada, realiza algún cálculo con esos datos y devuelve un resultado. Una función en Trident es exactamente lo mismo. Recibirá un conjunto de campos en la entrada y emitirá en la salida una o más tuplas. Las funciones en Trident serían el equivalente a los Bolts de Storm, digamos que se colocan en mitad del flujo de datos pudiéndoles llegar parte de esos datos. Una vez procesados dentro de la función los emiten pasando a formar parte del flujo original.



Síguenos a través de:



Últimas Noticias

» [Curso JBoss de Red Hat](#)

» Si eres el responsable o líder técnico, considérate desafortunado. No puedes culpar a nadie por ser gris

» [Portales, gestores de contenidos documentales y desarrollos a medida](#)

» [Comentando el libro Start-up Nation, La historia del milagro económico de Israel, de Dan Senor & Salu Singer](#)

» [Screencasts de programación narrados en Español](#)

[Histórico de noticias](#)

Últimos Tutoriales

» [Cómo se trabaja con un Croma y para qué sirve](#)

» [Creación de un módulo AMP de Alfresco con arquetipo Maven](#)

» [Introducción a Apache Storm](#)

» [Canon AX10: Una cámara de video para profesionales y aficionados.](#)

» [Cómo crear tu CV en](#)

Un ejemplo de función es la siguiente:

```
1 import storm.trident.operation.BaseFunction;
2 import storm.trident.operation.TridentCollector;
3 import storm.trident.tuple.TridentTuple;
4 import backtype.storm.tuple.Values;
5
6 public class ParamSplitter extends BaseFunction {
7     public void execute(TridentTuple tuple, TridentCollector collector) {
8         final String[] values = tuple.getString(0).split(",");
9
10        for (String value : values) {
11            collector.emit(new Values(value));
12        }
13    }
14 }
```

Imagina que nos llega una cadena de texto con diferentes valores separados por comas, cada uno representa un valor diferente. Por ejemplo si fuera un dato que nos devuelve un sensor de temperatura podría mandarnos una cadena con la fecha, hora, ciudad, temperatura, etc. Cada dato podría ser tratado de forma independiente por lo que antes de nada lo partimos y lo metemos de nuevo en el flujo para que otras funciones lo puedan tratar convenientemente. Recordar que cuando diseñamos este tipo de sistemas debemos tener en mente en todo momento la escalabilidad por lo que Storm ya está pensado para que cada función se pueda ejecutar en máquinas independientes.

3.2 Filter

Un filtro es una función que establece una condición para que el dato pueda o no continuar la cadena si cumple la regla definida en el filtro.

```
1 import storm.trident.operation.BaseFilter;
2 import storm.trident.tuple.TridentTuple;
3
4 public class PositiveFilter extends BaseFilter {
5     public boolean isKeep(TridentTuple tuple) {
6         return tuple.getInteger(0) > 0;
7     }
8 }
```

3.3 Aggregator

Una operación muy común cuando procesamos datos de un mismo tipo es la función de agregación por lo que Trident dispone de un método para crear nuestros propios agregadores o bien nos proporciona los más comunes como la suma o cuenta de valores:

Existen 3 tipos de agregadores: CombinerAggregator, ReducerAggregator y Aggregator.

Un CombinerAggregator debe devolver una tupla con un único campo de salida. Durante la ejecución del flujo se llamará al agregador con cada tupla de entrada y se irán combinando valores hasta el final del stream. Es útil también para combinar resultados en un único nodo antes de transferirlo por la red optimizando los recursos.

```
1 package storm.trident.operation.builtin;
2
3 import clojure.lang.Numbers;
4 import storm.trident.operation.CombinerAggregator;
5 import storm.trident.tuple.TridentTuple;
6
7 public class Sum implements CombinerAggregator<Number> {
8
9     @Override
10    public Number init(TridentTuple tuple) {
11        return (Number) tuple.getValue(0);
12    }
13
14    @Override
15    public Number combine(Number val1, Number val2) {
16        return Numbers.add(val1, val2);
17    }
18
19    @Override
20    public Number zero() {
21        return 0;
22    }
23 }
```

Storm llamará al método `init()` con cada tupla y posteriormente al `combine` hasta que la partición sea procesada. Los valores que se le pasan al `combine()` son parcialmente agregados.

El `ReducerAggregator` es bastante similar al `CombinerAggregator` con la diferencia que Storm llamará al `reduce()` hasta que la partición sea procesada completamente.

Por último tenemos el interfaz `Aggregate` para poder implementar nuestros propios agregadores:

```
1 import storm.trident.tuple.TridentTuple;
2
3 public interface Aggregator<T> extends Operation {
4     T init(Object batchId, TridentCollector collector);
5     void aggregate(T val, TridentTuple tuple, TridentCollector collector);
6     void complete(T val, TridentCollector collector);
7 }
```

Lo interesante de un `Aggregator` es que le podemos pasar a la función cualquier tipo por lo que es muy flexible pudiendo pasarle por ejemplo mapas, muy útiles para agrupar por distintas claves y poder hacer una cuenta sobre ellas.

3.4 Métodos para el procesamiento de streams

- `each`

Este método se utiliza para indicar la función o filtro que vamos a utilizar para procesar una tupla del stream de datos.

```
1 | stream.each(new Fields("str"), new ParamSplitter(), new Fields("date, hour, city, ?
```

formato europeo

Últimos Tutoriales del Autor

- » [Introducción a Apache Storm](#)
- » [Primeros pasos con Neo4j](#)
- » [Testing de Hadoop con MRUnit](#)
- » [Introducción a Spring Data Hadoop](#)
- » [Implementar una función UDF de Apache Pig](#)

Categorías del Tutorial

 Big Data

Indicamos que todas las tuplas del stream de entrada se pasen por la función ParamSplitter que como vimos antes se encargaba de dividir los campos separados por comas. La función emitirá tuplas en la salida pero es mediante el método each() donde configuramos el nombre de estas tuplas, similar a como cuándo configuramos los spouts y bols en una topología de Storm.

- **partitionAggregate**

Ejecuta una función de agregación y su resultado reemplazará las tuplas de entrada. Se ejecuta con todas las tuplas de todos los nodos.

```
1 stream.partitionAggregate(new Fields("values"), new Sum(), new Fields("sum"))
2
```

- **aggregate**

Ejecuta una función de agregación en un único nodo de forma aislada.

```
1 stream.partitionAggregate(new Fields("values"), new Count(), new Fields("count"))
2
```

- **project**

El método project() sirve para mantener en el flujo únicamente los campos especificados en la operación. Por ejemplo si el flujo tuviera los campos ["a","b","c","d"], si únicamente queremos quedarnos con los campos "a" y "b" realizaríamos la proyección:

```
1 stream.project(new Fields("a", "b"))
2
```

- **parallelismHint**

Configura la operación sobre una función para que sea ejecutada con el grado de paralelismo que le indiquemos.

```
1 stream.each(new Fields("str"), new ParamSplitter(), new Fields("date, hour, city,
2
```

- **partitionBy**

Encamina las tuplas para que sean procesadas las del mismo tipo en el mismo nodo de destino.

```
1 stream.partitionBy(new Fields("date"))
2
```

- **shuffle**

Utilizar el algoritmo round robin para redistribuir equitativamente las tuplas.

- **groupBy**

Se utiliza para agrupar las tuplas por un tipo.

```
1 stream.groupBy(new Fields("hashtag"))
2
```

3.5 State

Para almacenar el estado de las operaciones Trident soporta múltiples fuentes, desde almacenarlo en memoria, persistir en HDFS o almacenar los resultados en una base de datos NoSQL como Cassandra, Memcached, Redis, etc.

La administración del estado es tolerante a fallos mediante el procesamiento de las tuplas por lotes más pequeños asignando identificadores únicos que a la hora de persistir son consultados por el estado para mantener la consistencia de los datos.

```
1 TridentTopology topology = new TridentTopology();
2 TridentState wordCounts = topology.newStream("spout", spout)
3     .each(new Fields("sentence"), new Split(), new Fields("word"))
4     .groupBy(new Fields("word"))
5     .persistentAggregate(new MemoryMapState.Factory(), new Count(), new Fields("count"))
```

La lógica necesaria para gestionar las transacciones se realiza en la clase MemoryMapState proporcionada por el API de Trident.

La topología Storm creada a través de Trident se define mediante la clase TridentTopology donde se le indica el flujo de entrada de datos o spout. Este spout realiza la ingesta de datos metiendo streams de sentencias que serán procesadas por la función Split() que se encargará de separarlas por palabras. Una vez que se van separando se realiza una agrupación por cada palabra y una cuenta por número de apariciones que se almacenará en memoria.

4. Implementar un trending topics de Twitter con Trident.

Para ilustrar con un ejemplo lo que hemos visto sobre el API Trident vamos a construir una topología que se encargue de consumir los tweets recibidos de Twitter, extraiga los hashtags que tenga el tweet, si es que tiene alguno, y realice una cuenta para finalmente sacar una lista de los trending topics.

La topología será bastante parecida a la del tutorial de Storm. Para empezar vamos a crear un spout para recoger los tweets. En este caso nuestro spout implementa el interfaz IBatchSpout. A Twitter le vamos a pedir tweets que contengan información sobre equipos de fútbol para ver sobre qué se habla más en relación a los equipos de fútbol de España. Que me perdonen los aficionados del resto de equipos de fútbol pero por simplificar únicamente he puesto los 3 primeros en la clasificación de la temporada 13-14. Podéis probar metiendo diferentes topics si lo preferís ;)

```
1 import java.util.Map;
2 import java.util.concurrent.LinkedBlockingQueue;
3
4 import storm.trident.operation.TridentCollector;
5 import storm.trident.spout.IBatchSpout;
6 import twitter4j.FilterQuery;
7 import twitter4j.StallWarning;
8 import twitter4j.Status;
9 import twitter4j.StatusDeletionNotice;
10 import twitter4j.StatusListener;
11 import twitter4j.TwitterStream;
12 import twitter4j.TwitterStreamFactory;
```

```

13 import backtype.storm.Config;
14 import backtype.storm.task.TopologyContext;
15 import backtype.storm.tuple.Fields;
16 import backtype.storm.tuple.Values;
17 import backtype.storm.utils.Utils;
18
19 @SuppressWarnings({"serial", "rawtypes"})
20 public class TwitterConsumerBatchSpout implements IBatchSpout {
21
22     private LinkedBlockingQueue<Status> queue;
23     private TwitterStream twitterStream;
24
25     @Override
26     public void open(Map conf, TopologyContext context) {
27         this.twitterStream = new TwitterStreamFactory().getInstance();
28         this.queue = new LinkedBlockingQueue<Status>();
29
30         final StatusListener listener = new StatusListener() {
31
32             @Override
33             public void onStatus(Status status) {
34                 queue.offer(status);
35             }
36
37             @Override
38             public void onDeletionNotice(StatusDeletionNotice sdn) {
39             }
40
41             @Override
42             public void onTrackLimitationNotice(int i) {
43             }
44
45             @Override
46             public void onScrubGeo(long l, long ll) {
47             }
48
49             @Override
50             public void onException(Exception e) {
51             }
52
53             @Override
54             public void onStallWarning(StallWarning warning) {
55             }
56         };
57
58         twitterStream.addListener(listener);
59
60         final FilterQuery query = new FilterQuery();
61         query.track(new String[]{"atleti", "fcbarcelona", "realmadrid"});
62         query.language(new String[]{"es"});
63
64         twitterStream.filter(query);
65     }
66
67     @Override
68     public void emitBatch(long batchId, TridentCollector collector) {
69         final Status status = queue.poll();
70         if (status == null) {
71             Utils.sleep(50);
72         } else {
73             collector.emit(new Values(status));
74         }
75     }
76
77     @Override
78     public void ack(long batchId) {
79     }
80
81     @Override
82     public void close() {
83         twitterStream.shutdown();
84     }
85
86     @Override
87     public Map getComponentConfiguration() {
88         return new Config();
89     }
90
91     @Override
92     public Fields getOutputFields() {
93         return new Fields("tweet");
94     }
95 }

```

Lo siguiente será crear la topología. Creamos un método main donde configuramos la topología, por simplicidad la arrancaremos en modo local. Construimos la topología trident pasándole el spout creado anteriormente.

```

1
2 import java.io.IOException;
3
4 import storm.trident.TridentTopology;
5 import storm.trident.operation.builtin.Count;
6 import storm.trident.operation.builtin.Debug;
7 import storm.trident.spout.IBatchSpout;
8 import storm.trident.testing.MemoryMapState;
9 import backtype.storm.Config;
10 import backtype.storm.LocalCluster;
11 import backtype.storm.generated.StormTopology;
12 import backtype.storm.tuple.Fields;
13
14 import com.autentia.tutoriales.functions.HashtagExtractor;
15 import com.autentia.tutoriales.spout.TwitterConsumerBatchSpout;
16
17 public class TrendingTopicsTridentTopology {
18
19     public static StormTopology createTopology(IBatchSpout spout) throws IOException {
20         final TridentTopology topology = new TridentTopology();
21

```

```

22     topology.newStream("spout", spout)
23         .each(new Fields("tweet"), new HashtagExtractor(), new Fields("hashtag"))
24         .groupBy(new Fields("hashtag"))
25         .persistentAggregate(new MemoryMapState.Factory(), new Count(), new Fiel
26         .newValuesStream()
27         .each(new Fields("hashtag", "count"), new Debug());
28
29     return topology.build();
30 }
31
32 public static void main(String[] args) {
33     final Config conf = new Config();
34     final LocalCluster local = new LocalCluster();
35     final IBatchSpout spout = new TwitterConsumerBatchSpout();
36
37     try {
38         local.submitTopology("hashtag-count-topology", conf, createTopology(spout));
39     } catch (IOException e) {
40         throw new RuntimeException(e);
41     }
42 }
43 }

```

Toda topología debe recibir un stream de datos, en nuestro caso el TwitterConsumerBatchSpout que irá metiendo en el sistema los tweets recibidos, a continuación configuramos las operaciones. En primer lugar, cada tweet se procesa en la función HashtagExtractor para extraer únicamente los hashtags que son pasados de nuevo al flujo para posteriormente ser agrupados. Posteriormente a la agrupación se realiza una cuenta sobre ellos siendo necesario almacenar el estado de la misma, en este caso se realiza en memoria.

Para terminar y poder ver los resultados obtenidos pasamos los valores "hashtag" y "count" por la función Debug que los imprime por consola.

La función HashtagExtractor es muy sencilla:

```

1     import storm.trident.operation.BaseFunction;
2     import storm.trident.operation.TridentCollector;
3     import storm.trident.tuple.TridentTuple;
4     import twitter4j.HashtagEntity;
5     import twitter4j.Status;
6     import backtype.storm.tuple.Values;
7
8     @SuppressWarnings("serial")
9     public class HashtagExtractor extends BaseFunction {
10
11         @Override
12         public void execute(TridentTuple tuple, TridentCollector collector) {
13             final Status status = (Status) tuple.get(0);
14
15             for (HashtagEntity hashtag : status.getHashtagEntities()) {
16                 collector.emit(new Values(hashtag.getText()));
17             }
18         }
19     }

```

Recoge la tupla cero que contiene el tweet, extrae sus hashtags y los emite a la topología.

Si ejecutamos la clase TrendingTopicsTridentTopology veremos cómo van apareciendo hashtags por consola con un contador que irá incrementándose sucesivamente.

```

1     DEBUG: [FelizLunes, 5]
2     DEBUG: [futbol, 8]
3     DEBUG: [FCBarcelona, 15]
4     DEBUG: [Baloncesto, 2]
5     DEBUG: [HalaMadrid, 11]
6     DEBUG: [HM, 2]
7     DEBUG: [ReyesDeEuropa, 2]
8     DEBUG: [RealMadrid, 28]
9     DEBUG: [Atleti, 10]
10    DEBUG: [Cholismo, 3]
11    DEBUG: [ElChiringuitoDeNeox, 2]
12    DEBUG: [LigaBBVA, 2]

```

5. Conclusiones.

Storm va madurando poco a poco y se está posicionando en el mercado como uno de los mejores productos para procesamiento de datos en real time junto a otros como Spark.

Los usuarios de Storm estamos de enhorabuena ya que recientemente (el 29 de septiembre de 2014) Storm ha pasado a ser considerado como un Top-Level Project (TLP) dentro de Apache, muy buena noticia ya que supone la graduación de este fantástico framework de procesamiento en tiempo real.

Puedes descargar el código del tutorial desde [aquí](#).

Espero que te haya sido de ayuda.

Un saludo.

Juan

A continuación puedes evaluarlo:

[Regístrate para evaluarlo](#)

★★★★★

Por favor, vota +1 o compártelo si te pareció interesante

Anímate y coméntanos lo que pienses sobre este **TUTORIAL**:

» **Regístrate** y accede a esta y otras ventajas «



Esta obra está licenciada bajo licencia Creative Commons de Reconocimiento-No comercial-Sin obras derivadas 2.5

PUSH THIS | Page Pushers | Community | Help?

0 people brought clicks to this page

no clicks

       

powered by [karmacracy](#)