

Storm Adaptive Schedulers – User Manual

This document contains the instructions for compiling, deploying and using the general-purpose adaptive schedulers for Storm (<http://storm-project.net/>), developed by MIDLAB research group (<http://www.dis.uniroma1.it/~midlab/>). For detailed info about the software, please refer to

Leonardo Aniello, Roberto Baldoni, Leonardo Querzoni,
"Adaptive Online Scheduling in Storm",
7th ACM International Conference on Distributed Event-Based Systems (DEBS 2013),
available online at <http://www.dis.uniroma1.it/~midlab/articoli/ABQ13storm.pdf>.

Requirements

Current implementation only works in UNIX-like environments.

Where to download

The zip file containing the source code is available online at
<http://www.dis.uniroma1.it/~midlab/software/storm-adaptive-schedulers.zip>.

How to compile

The `src` directory in the zip contains all the sources required for the compilation. The external libraries to include are

- Storm 0.8.1 libraries
(<https://dl.dropbox.com/u/133901206/storm-0.8.1.zip>)
 - `storm-0.8.1.jar`
 - `lib/*.jar`
- Commons DBCP 1.4
(<http://archive.apache.org/dist/commons/dbcp/binaries/commons-dbc-1.4-bin.zip>)
 - `commons-dbc-1.4.jar`
- Commons Pool 1.5.6
(<http://archive.apache.org/dist/commons/pool/binaries/commons-pool-1.5.6-bin.zip>)
 - `commons-pool-1.5.6.jar`

It is recommended to compile using Java 1.6.x.

How to deploy

Generally, a Storm deployment includes a *master node* (where nimbus service runs) and several *slave nodes* (where supervisor services run). In each node (either master or slave), these jars have to be copied into the `lib` directory of Storm installation

- the jar resulting from the compilation of Storm adaptive schedulers
- `commons-dbc-1.4.jar`

- commons-pool-1.5.6.jar

The online scheduler needs a MySQL database to store and retrieve online statistics about inter-executor traffic and CPU load. The required MySQL version is 5.5.x. The SQL script `storm-scheduler.sql` included into the zip is for creating all the required DB tables. A MySQL Java connector 5.1.x is required at runtime to enable the interaction with the DB. It can be downloaded from <http://downloads.mysql.com/archives.php?p=mysql-connector-java-5.1&o=other> and has to be copied into the `lib` directory of Storm installation, for each node (either master or slave).

In order to configure the master node, two files have to be modified: `storm.yaml` (into `conf` folder of Storm installation) and `db.ini` (to be put into the root directory of Storm installation). Next two tables describe how to configure such files.

storm.yaml In this configuration file, the format to use for the single parameter-value pair is <code>parameter: "value"</code> All these parameters are mandatory		
Parameter	Description	Allowed values
<code>storm.scheduler</code>	Specifies the fully qualified name of the class to use as scheduler	For the offline scheduler: <code>midlab.storm.scheduler.OfflineScheduler</code> For the online scheduler: <code>midlab.storm.scheduler.OnlineScheduler</code>
<code>traffic.improvement</code>	Specifies the minimum percent inter-node traffic improvement required to apply a new schedule ¹	An integer between 1 and 100 (for online scheduler only)
<code>reschedule.timeout</code>	Specifies the minimum amount of time (in seconds) that has to pass before a new schedule can be applied	A positive integer (for online scheduler only)

db.ini In this configuration file, the format to use for the single parameter-value pair is <code>parameter=value</code> All these parameters are mandatory		
Parameter	Description	Allowed values
<code>jdbc.driver</code>	Specifies the fully qualified name of the class to use as DB connector	If MySQL 5.5.x is used, then the correct value is <code>com.mysql.jdbc.Driver</code>
<code>data.connection.uri</code>	Specifies the DB connection URI	If MySQL 5.5.x is used, then the correct value is (it's a single value to be written within a single line) <code>jdbc:mysql://<DB_IP_ADDRESS>/storm-scheduler?user=<USER>&password=<PWD></code> where <ul style="list-style-type: none"> • <code><DB_IP_ADDRESS></code> is the IP address of the host where MySQL DBMS is running

¹ Let T be current inter-node traffic, and T' the estimate of inter-node traffic produced by a new schedule S . If $(T - T') / T \geq \text{traffic.improvement} / 100$, then schedule S is applied.

		<ul style="list-style-type: none"> • <USER> is the username to use for the connection • <PWD> is the password of the specified username
validation.query	Specifies the query to use for checking whether a DB connection is still alive	SELECT

In order to configure a slave node, only `db.ini` (to be put into the root directory of Storm installation) has to be configured. In addition to the three parameters previously described in the table for the `db.ini` for the master node (`jdbc.driver`, `data.connection.uri`, `validation.query`), others have to be configured and are described in the next table.

db.ini In this configuration file, the format to use for the single parameter-value pair is <code>parameter=value</code> All these parameters are mandatory		
Parameter	Description	Allowed values
node-name	Specifies an identification name (unique ID) for the slave node	Any string, with the constraint that any two distinct slave nodes are assigned distinct names
capacity	Specifies the maximum percentage of CPU usage allowed for this slave node ²	An integer between 1 and 100
time.window.slot.count	Specifies the number of slots for a single time window ³	A positive integer
time.window.slot.length	Specifies the length (in seconds) of a single slot of a time window	A positive integer

How to use

Some simple lines of code have to be added to spout/bolt source in order to enable the monitoring of inter-executor traffic and CPU load. In the specific each task has to notify its thread ID to a singleton object in charge of monitoring all the tasks running in a specific Java process (Storm worker).

For each spout/bolt, an instance field has to be added

```
private TaskMonitor taskMonitor;
```

Some initialization operations have to be performed in the `open()` method of each spout

```
public void open(Map conf, TopologyContext context, SpoutOutputCollector collector) {
    ....
    // register this spout instance (task) to the java process monitor
    WorkerMonitor.getInstance().setContextInfo(context);
    // create the object required to notify relevant events (also notify thread ID)
    taskMonitor = new TaskMonitor(context.getThisTaskId());
    ...
}
```

² The total capacity of a node is measured in MHz by looking at CPU characteristics; the CPU demand (in MHz) for each executor is tracked in order to estimate the expected CPU load for each slave node when a new schedule S is applied.

³ The measurements of inter-executor traffic and CPU load are carried out on a sliding time window: each time window is divided into a fixed number of equal-length slots, and the window slides by a single slot per time.

The same has to be done in the `prepare()` method of each bolt

```
public void prepare(Map stormConf, TopologyContext context, OutputCollector collector) {
    ...
    // register this spout instance (task) to the java process monitor
    WorkerMonitor.getInstance().setContextInfo(context);
    // create the object for notifying relevant events (received tuple, which in turn
    notifies thread ID)
    taskMonitor = new TaskMonitor(context.getThisTaskId());
    ...
}
```

The notification of thread ID for the spouts is done in the `nextTuple()` method

```
public void nextTuple() {
    taskMonitor.checkThreadId();
    ...
}
```

Finally, the notifications of received tuples for the bolts is done in the `execute()` method

```
public void execute(Tuple input) {
    taskMonitor.notifyTupleReceived(input);
    ...
}
```

Other parameters can be configured for each topology. Their names vary a little with respect to what is written in the paper referenced at the beginning of this document. Available parameters are described in next table. These parameters have to be specified in the source code of the topology: they have to be put into a `java.util.Map` instance that has to be passed to `submitTopology()` method.

Topology parameters		
Let <code>conf</code> be a <code>java.util.Map</code> instance, the code to use for the single parameter-value pair is <code>conf.put(parameter, value)</code>		
For a generic topology t_i , these variables are defined		
<ul style="list-style-type: none"> • E_i: number of executors • W_i: number of workers • C_i: number of components (spouts/bolts) 		
The number of slave nodes in the cluster is N .		
Parameter	Description	Allowed values
alfa	Controls the balancing of the number of executors assigned per slot. In particular, <code>alfa</code> affects the maximum number M of executors that can be placed in a single slot. The minimum value of M for a topology t_i is $\text{ceil}(E_i / W_i)$, which means that each slot roughly contains the same number of executors. The maximum number of M corresponds to the assignment where all the slots contain one executor, except for one slot that contains all the other executors, so its value is $E_i - W_i + 1$. $M(\text{alfa}) = \text{ceil}(E_i / W_i) + \text{alfa}(E_i - W_i + 1 - \text{ceil}(E_i / W_i))$	A float in the range [0,1]
beta	Controls the number of nodes to use. It works similarly to <code>alfa</code> . If <code>beta = 0</code> , then $\min = \text{ceil}(W_i / N)$ nodes are used. If <code>beta = 1</code> , then $\max = \min(W_i, N)$ nodes are used. In general, $\min + \text{beta} * (\max - \min)$ nodes are used.	A float in the range [0,1]
gamma (only for online scheduler)	Controls the maximum load acceptable by a slot, it is computed as $\text{gamma} * (\text{totalLoad} / W_i)$, where <code>totalLoad</code> is the sum of loads of all the executors of topology t_i .	A float greater than 1

delta	<p>Controls the maximum number of slots per node.</p> <p>The minimum value is $\min = \text{ceil}(W_i / N)$.</p> <p>The maximum value is $\max = W_i - N + 1$.</p> <p>The actual value is computed as $\min + \text{delta} * (\max - \min)$.</p>	A float in the range [0,1]
epsilon (only for offline scheduler)	<p>Controls when to start considering empty slots.</p> <p>During the assignment of executors for the i^{th} component, the scheduler is forced to use empty slots if $i > \text{floor}(\text{delta} * C_i)$.</p> <p>For example, if traffic is likely to be more intense among upstream components, then delta should be set large enough such that empty slots get used when upstream components are already assigned.</p>	A float in the range [0,1]