

Storm on Multi-core

Mark Nemec

4th Year Project Report
Computer Science
School of Informatics
University of Edinburgh

2015

Abstract

In recent years there has been a growing usage of data stream processing systems. However, data stream processing has been a province of distributed systems such as Apache Storm. This is because clusters offer better scale-out properties.

On the other hand, chip makers have been able to pack more and more transistors into a single processor chip. This resulted in a steady increase of the number of cores used in a single processor chip. It is possible that scale-out processors with high number of processor cores will be used for future throughput-based applications.

In this report we present Storm-MC, an API-compatible port of Apache Storm for multi-core. This system combines an easy to use API with implementation tailored to multi-core environments. This allows us to take existing applications written with Storm in mind and run them efficiently on a single server.

Through a series of benchmarks we show that running applications on Storm-MC can provide substantial improvement in throughput (up to 3.3x) compared to running them on Apache Storm in local mode.

Acknowledgements

I would like to thank my supervisor, Dr. Stratis Viglas, for providing guidance every step of the way.

Furthermore, I would like to thank my parents who supported me during my studies at the University of Edinburgh. Without them this project would not be possible.

Finally, I would like to thank Daniela Kellerova who kept me sane while I worked on this project.

Table of Contents

| | | |
|----------|--|-----------|
| I | Introduction | I |
| 1.1 | Motivation | I |
| 1.2 | Project Contributions | 2 |
| 1.3 | Structure of the Report | 2 |
| 2 | Literature Review | 3 |
| 2.1 | Data Stream | 3 |
| 2.1.1 | Comparison to a Database Management System | 3 |
| 2.1.2 | Querying a Data Stream | 4 |
| 2.2 | Multi-core | 4 |
| 2.2.1 | Advantages Over Clusters | 5 |
| 2.2.2 | Disadvantages Over Clusters | 5 |
| 2.3 | Apache Storm | 6 |
| 2.3.1 | Dependencies | 6 |
| 2.3.2 | Usage of Apache Storm | 7 |
| 2.4 | Similar Efforts | 7 |
| 2.5 | Summary | 7 |
| 3 | Background | 9 |
| 3.1 | Storm Overview | 9 |
| 3.2 | Storm Concepts | 9 |
| 3.2.1 | Core Concepts | 9 |
| 3.2.2 | Additional Concepts | 10 |
| 3.3 | Example Topology | 11 |
| 3.4 | Storm Architecture | 12 |
| 3.4.1 | Nimbus Node | 12 |
| 3.4.2 | Worker Nodes | 14 |
| 3.4.3 | Zookeeper Nodes | 14 |
| 3.5 | Serialisation | 15 |
| 4 | Bringing Storm to Multi-core | 17 |
| 4.1 | Apache Storm on Multi-core | 17 |
| 4.1.1 | Tuple Processing Overhead | 17 |
| 4.1.2 | Thread Overhead | 18 |
| 4.2 | Storm-MC Design | 20 |
| 4.2.1 | Porting Nimbus | 20 |

| | | |
|----------|---|------------|
| 4.2.2 | Porting Worker Nodes | 21 |
| 4.2.3 | Removing State | 22 |
| 4.2.4 | Removing Serialisation | 23 |
| 4.3 | Implementation Details | 23 |
| 4.3.1 | Topology Submission | 23 |
| 4.3.2 | Tuple Processing | 24 |
| 4.3.3 | Executors | 26 |
| 4.4 | Differences between Apache Storm and Storm-MC | 26 |
| 5 | Evaluation | 29 |
| 5.1 | Evaluation Metrics | 29 |
| 5.2 | System Configuration | 29 |
| 5.2.1 | Software Setup | 29 |
| 5.2.2 | Hardware Setup | 30 |
| 5.2.3 | Storm Configuration | 31 |
| 5.3 | Results | 31 |
| 5.3.1 | WordCount Topology | 32 |
| 5.3.2 | Enron Topology | 34 |
| 5.3.3 | RollingSort Topology | 35 |
| 5.4 | Challenges | 36 |
| 6 | Conclusion | 39 |
| 6.1 | Summary of Contributions | 39 |
| 6.2 | Future Work | 39 |
| A | Listings | III |

List of Figures

| | | |
|-----|---|----|
| 2.1 | Stream Querying. | 4 |
| 3.1 | WordCount topology. | 11 |
| 3.2 | Apache Storm Architecture. | 12 |
| 4.1 | Tuple processing in Apache Storm. | 17 |
| 4.2 | Thread state distribution over time | 20 |
| 4.3 | Comparison of a worker in Storm and Storm-MC | 22 |
| 4.4 | Tuple processing in Storm-MC. | 24 |
| 5.1 | CountBolt throughput in Apache Storm and Storm-MC | 33 |
| 5.2 | Number of threads used by Apache Storm and Storm-MC | 33 |
| 5.3 | Number of threads used by Apache Storm and Storm-MC. | 34 |
| 5.4 | Enron topology. | 34 |
| 5.5 | Global email throughput over time with standard error | 36 |
| 5.6 | RollingSort topology. | 36 |

List of Tables

| | | |
|-----|---|----|
| 4.1 | Storm thread usage | 19 |
| 4.2 | Storm thread states | 19 |
| 4.3 | Feature comparison of Apache Storm and Storm-MC. | 26 |
| 5.1 | Storm-MC: Component throughput in WordCount Topology. . . | 32 |
| 5.2 | Apache Storm: Component throughput in WordCount Topology. | 32 |
| 5.3 | Storm-MC: Email Throughput in Enron Topology. | 35 |
| 5.4 | Apache Storm: Email Throughput in Enron Topology. | 35 |
| 5.5 | Storm-MC: Component throughput in RollingSort Topology. . . | 37 |
| 5.6 | Apache Storm: Component throughput in RollingSort Topology. . | 37 |

List of Listings

| | | |
|---|------------------------------------|----|
| I | WordCountTopology.java | I2 |
| 2 | RandomSentenceSpout.java | IV |
| 3 | SplitSentence.java | V |
| 4 | splitsentence.py | V |
| 5 | WordCount.java | VI |

Chapter 1

Introduction

1.1 Motivation

In recent years, there has been an explosion of cloud computing software. After Google published their paper on MapReduce [1], many new open-source frameworks for distributed computation have emerged, most notably Apache Hadoop [2] for batch processing and Apache Storm [3] for real-time data stream processing.

These frameworks split the work that needs to be carried out and distribute it across nodes of a commodity hardware cluster. Commercial companies and researchers have been able to utilise these frameworks and create distributed systems which can accomplish things that would not be otherwise possible [4]. This has mostly been allowed by the low price and good scale-out properties of commodity hardware.

At the same time, chip makers have been increasing the number of cores in processors and now we are at a point where servers with 10-core processors are considered standard. Moreover, most high-end servers have multiple processor sockets thus increasing the parallelisation possible with a single machine even further.

The number of cores is increasing but the price of highly parallel machines has gone down. In 2008, a typical Hadoop node had two dual-core processors and 4 GB of random access memory (RAM) [5]. Nowadays, a server with two eight-core processors and 256 GB of RAM can be purchased for roughly \$10,000 USD [5]. Hence a single server today might have better processing power than a small cluster from a few years ago. If this trend continues there will be processors with even more cores in the near future with higher processing power than most clusters today.

Moreover, tiled processors have emerged as competitors to traditional processors in throughput-based computations [6]. These processors use a large number of tiles connected by an on-chip network and even though the single-thread performance of each tile is lower than the performance of a conventional core, the increased parallelism yields higher throughput [7].

Seeing these trends, we believe there is a place for a real-time data stream processing system running on a single multi-core machine.

It is generally believed that writing parallel software is hard. The traditional techniques of message passing and shared memory require the programmer to manage concurrency at a fairly low level. Furthermore, Apache Storm has become the *de facto* tool used in stream processing on a cluster and according to their “Powered By” page there are tens of companies already using Storm to process their real-time data streams [8]. We think that Storm’s popularity and easy to use application programming interface (API) makes it the ideal candidate for porting to multi-core.

1.2 Project Contributions

The main idea of this project is to take the existing Apache Storm project and port it to multi-core. This is implemented in Storm-MC – a library with an API compatible with Apache Storm. This compatibility enables programmers to take an existing application written with Apache Storm in mind and run it on a single multi-core server using Storm-MC. This way, they can avoid network latency and enjoy the substantial performance improvements of a shared-memory environment.

Through a series of benchmarks we show that with its simpler design, Storm-MC offers substantial improvement in throughput for data stream processing applications over Apache Storm running in local mode.

1.3 Structure of the Report

The remainder of the report is structured as follows:

- Chapter 2 presents an overview of related literature and gives background on data stream processing and multi-core architectures.
- Chapter 3 explains the concepts used in Apache Storm as well as the architecture of a Storm cluster.
- Chapter 4 describes how Apache Storm was ported to multi-core and explains the design of Storm-MC.
- Chapter 5 presents results of benchmarking Storm-MC against Apache Storm running in local mode.
- Chapter 6 summarises this report and presents considerations for future work.

Chapter 2

Literature Review

The following chapter explains the concept of a data stream (2.1), discusses advantages and disadvantages of multi-core (2.2), gives an overview of previous work on Apache Storm (2.3), and discusses other effort of porting distributed systems to multi-core (2.4).

2.1 Data Stream

Golab and Özsu [9] define data stream as as “a real-time, continuous, ordered (implicitly by arrival time or explicitly by timestamp) sequence of items. It is impossible to control the order in which items arrive, nor is it feasible to locally store a stream in its entirety.”

2.1.1 Comparison to a Database Management System

Historically, data has been stored into database management systems (DBMS) where it was later analysed the assumption being that there would be enough disk space to contain the data. This approach fits many purposes but real-time applications have recently started feeling the need to analyse rapidly changing data on-the-fly.

This has brought on an advent of data stream processing. Several stream-processing frameworks have emerged such as Apache Storm [3], Apache Spark [10], and Yahoo S4 [11]. These frameworks, usually ran on a cluster, provide the user with abstractions which greatly simplify writing a real-time data stream processing application.

Whereas DBMSs excel at getting an exact answer to a query, data streams usually provide an approximate answer. The answer is approximate because it is usually correct only within a certain window of time, the query is simplified because it can only be ran in one pass, or because it is used with a sampling rate which does not

include all events. A typical data stream analysis using windows and sampling is shown in Figure 2.1.

2.1.2 Querying a Data Stream

The assumption behind using a time window is that users of the real-time system are most likely interested in the most recent events. Sampling, on the other hand, is used to reduce the number of events used for a query [12].

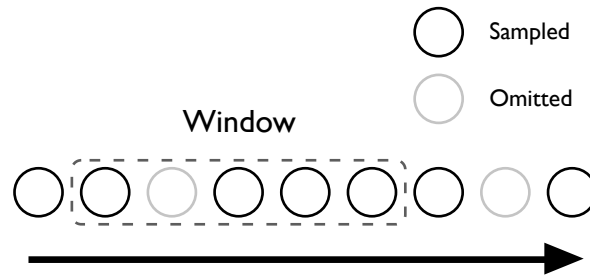


Figure 2.1: Stream Querying.

Even though the answer might only be approximate it can have great value because the query is answered at the right time. Furthermore, even though the query may run only on a subset of data it is still possible to detect trends or system failures. For example, Twitter are using Apache Storm to run real-time analysis on millions of events per second for their analytics product Answers [4].

In research, several techniques have been developed to enable real-time data stream mining. For example, the MOA environment created by Bifet, Holmes, Kirkby, *et al.* [13] enables running online machine learning methods using the WEKA machine learning workbench [14].

2.2 Multi-core

Akhter and Roberts [15] define multi-core processor as a processor that has two or more independent execution cores. This means that every thread has a hardware execution environment entirely to itself which enables threads to run in a truly parallel manner.

Running an application on multi-core can in the best case produce a speedup equivalent to the number of cores. The best case is when an application is embarrassingly parallel i.e. there is no inter-thread communication. Even though data stream processing is not an embarrassingly parallel problem, running a producer-consumer application on multiple cores can produce significant speedup as shown in [16].

2.2.1 Advantages Over Clusters

There are several reasons why someone might prefer to deploy their data stream application to a single multi-core machine over a cluster:

Communication Overhead

The latency of over-the-network communication is significantly higher than of two cores communicating on a single machine.

Lower Cost than a Data Centre

To run a distributed system on a cluster one would usually need to own a data centre. This comes with a high capital cost and increased maintenance costs than owning a single server.

More Control than with a Cloud Provider

Alternatively, one could rent out nodes on cloud computing services such as Amazon EC2 or Rackspace. While the cost of such services is acceptable the user does not have full control over their system.

2.2.2 Disadvantages Over Clusters

On the other hand, there are certain disadvantages in running a computation on a single multi-core machine rather than a cluster:

Horizontal Scaling

Commodity hardware clusters offer better horizontal scaling than a single multicore server. If one needs to add nodes to the cluster it is as easy as purchasing more commodity hardware. On a multi-core machine it is not that simple. For example, a top of the line Intel® Xeon® Processor E5-2699 has support for 36 threads. Beyond that one would need to add another socket which essentially doubles the price. Moreover, the server would need to support multiple processor sockets.

Higher Short-term Cost than with a Cloud Provider

In short-term the cost of purchasing a server may be significantly higher than renting out a cluster from a cloud computing service. Thus, it makes most sense to run an application on multi-core as a medium- to long-term investment.

More Maintenance than with a Cloud Provider

Owning a server requires more maintenance than simply renting it out from a cloud provider. A cloud provider usually does all the necessary maintenance and can provision a new machine very easily. Hence it is advantageous to use multi-core only if one can afford to maintain them as well.

2.3 Apache Storm

Apache Storm is an open source distributed real-time computation system. Storm was Originally created by Nathan Marz while working at BackType. [17] BackType was later acquired by Twitter which is when Storm became open source. Storm was incubated into Apache with version 0.9.1 and became a top-level Apache project in September 2014.

Storm was developed to run on top of a cluster where nodes execute components of a computation in parallel. Running Storm on a cluster of commodity hardware gives it good horizontal scaling properties. Running separate components in parallel allows the system to execute in real time.

2.3.1 Dependencies

Storm has five major dependencies:

Apache Zookeeper

Apache Zookeeper [18] is an open source server that allows reliable distributed coordination. Storm uses Apache Zookeeper to maintain state which is then read and written to by nodes of a Storm cluster. More detail on how Storm uses Zookeeper is given in section 3.4.3.

Apache Thrift

Apache Thrift [19] is a cross-language framework for developing services. It allows you to write a definition file for services and data types required by your application and automatically generates interface code which supports remote procedure calls and serialisation of the data types.

Kryo

Kryo [20] is a serialisation library. It is used by Storm to serialise objects when sent over the network between nodes of a cluster.

Netty

Netty [21] is an asynchronous event-driven network application framework. Storm utilises Netty to send intra-cluster messages. Thus when a node produces a result to be consumed by another node of the cluster it sends a message over the network using the TCP protocol implemented in Netty.

LMAX Disruptor

LMAX Disruptor [22] is a high-performance data structure used to exchange data between concurrent threads. It uses a lock-free implementation of a ring buffer which components of a Storm program running on the same node a cluster use to exchange messages.

2.3.2 Usage of Apache Storm

Storm works particularly well with sister Apache projects such as Apache Kafka [23] and Apache HBase [24]. Apache Kafka is a messaging broker that is often used as the missing link between producers and consumers of a cluster. Apache HBase is a big data-store that allows real-time random reads and writes modelled after Google's Bigtable project [25].

Storm is reportedly used by 81 companies listed on their website [8] and many others. Storm's popularity is one of the reasons why it was chosen for this project. Furthermore, we believe that the concepts used in Storm (explained further in section 3.2) apply to many different situations and many applications can be easily adapted to work on top of Storm.

There has been significant research in optimising computations running on Apache Storm. For example, [26] looked at how to reconfigure a Storm job by reallocating component tasks to minimise communication cost. A domain specific language for defining Storm jobs was proposed in [27]. Finally, [28] ported Storm to Haskell and looked at how to provide exactly once semantics.

2.4 Similar Efforts

Currently data stream processing is a province of distributed systems such as the ones mentioned in the previous section. Most languages support parallel execution and there are many libraries that ease the process of writing parallel programs on a single multi-core machine. However, they are not tailored to data stream processing and usually require the programmer to do the heavy lifting rather than abstract it away.

There has been effort to port Hadoop to multi-core in [5] (Hone) as well as port of Google's MapReduce in [29] (Phoenix). However, to our knowledge there has not been effort to port Storm or any other distributed real-time data stream system to multi-core.

2.5 Summary

Distributed real-time computation systems such as Apache Storm provide programmers with abstractions that make it very easy to implement a data stream processing applications on top of a cluster. However, in case of single multi-core machines there are not any obvious software choices. While several frameworks that allow the programmer to parallelise the computation exist, they are not really tailored to data stream processing.

The following chapters provide a closer analysis of Apache Storm as well as a port of Apache Storm for a single multi-core machine.

Chapter 3

Background

In this chapter we give background information necessary to understand the design of Storm-MC. We give a quick overview of Apache Storm (3.1), explain the concepts used in Storm (3.2), show an example Storm program (3.3), give details about the underlying architecture of Storm (3.4), and finally describe the serialisation used by Storm (3.5).

3.1 Storm Overview

Apache Storm was developed in a mix of Java and Clojure. As mentioned by the author of Storm in [30], writing the Storm interfaces in Java ensured large potential user-base while writing the implementation in Clojure increased productivity.

To ensure API compatibility with Storm, Storm-MC was developed using the same set of languages. This allowed for code reuse and not having to re-implement functionality already present in Storm. Hence, in the following sections we describe Storm in greater detail in hope that this will later clarify design choices made for Storm-MC.

3.2 Storm Concepts

3.2.1 Core Concepts

There are several core concepts used by Storm and hence by extension Storm-MC as well. These concepts are put together to form a simple API that allows the programmer to break down a computation into separate components and define how these components interact with each other. The three core concepts of Storm are:

Spout

A spout is a component that represents the source of a data-stream. Typically, a spout reads from a message broker such as RabbitMQ [31] or Apache Kafka but can also generate its own stream or read from somewhere like the Twitter streaming API [32].

Bolt

A bolt is a component that transforms tuples from its input data stream and emits them to its output data stream. A bolt can perform a range of functions e.g. filter out tuples based on some criteria or perform a join of two different input streams.

Topology

The programmer connects spouts and bolts in a directed graph called topology which describes how the components interact with each other. The topology is then submitted to Storm for execution.

3.2.2 Additional Concepts

Stream

A stream is defined as an unbounded sequence of tuples. Streams can be thought of as edges of a topology connecting bolts and spouts (vertices).

Tuple

A tuple wraps named fields and their values. The values of the fields can be of different types. When a component emits a tuple to a stream it sends that tuple to every bolt subscribed to the stream.

Stream Grouping

Every bolt needs to have a type of stream grouping associated with it. This grouping decides the means of distributing the tuples coming from the bolt's input streams amongst the instances of the bolt task. Following are the possible types of stream grouping:

Shuffle Randomly partition the tuples among all the bolt tasks.

Fields Hash on a subset of the tuple fields. All tuples with same values of those fields will go to same bolt task.

All Replicate the entire stream to all the bolt tasks.

Direct The producer of the tuple decides which task of the bolt will receive this tuple.

Global Send the entire stream to a single bolt task.

Local or Shuffle Prefer sending to executors in the same worker process, if that is not possible use same strategy as Shuffle.

Users are also able to specify their own custom grouping by implementing the CustomStreamGrouping interface.

All the components of a Storm topology execute in parallel. The user can specify how much parallelism he wants associated with every component and Storm spawns the necessary number of threads. This is done through a configuration file, defined in YAML, which is submitted along with the topology.

There are two additional bolts running for every topology:

Acker

The Acker bolt guarantees fault tolerance for the topology. It tracks every tuple that was produced and ensures that the tuple has been acknowledged by every bolt of the stream.

System Bolt

The System bolt is useful in two ways:

Metrics System bolt collects metrics on the local Java Virtual Machine (JVM). Other components can subscribe to these metrics and receive their values at regular intervals.

Ticks Components of a topology can subscribe to receive tick tuples in regular intervals. These tuples can be used to trigger some event of a component.

3.3 Example Topology

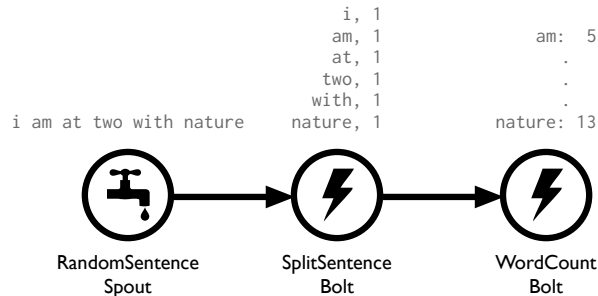


Figure 3.1: WordCount topology.

A classic example used to explain Storm topologies is the WordCount topology. In this topology, there is a spout generating random sentences, a bolt splitting the sentences on white space, and a bolt counting occurrences of every word. Figure 3.1 shows how we could represent this topology graphically.

This may seem as a simplistic example but it is useful when demonstrating how easy it is to implement a working topology using the Storm API.

Listing 1 shows how the topology is put together in Storm to form a graph of components. Storm uses the Builder design pattern [33] to build up the topology which is then submitted to Storm for execution. The last argument to the `setBolt/setSpout` method is the number of parallel tasks we want Storm to execute

for the respective component. For implementation of the spout and bolts used in this topology, refer to appendix A.

```
public class WordCountTopology {
    public static void main(String[] args) throws Exception {
        TopologyBuilder builder = new TopologyBuilder();
        builder.setSpout("spout", new RandomSentenceSpout(), 5);
        builder.setBolt("split",
            new SplitSentence(), 8).shuffleGrouping("spout");
        builder.setBolt("count",
            new WordCount(), 12).fieldsGrouping("split", new Fields("word"));
        LocalCluster cluster = new LocalCluster();
        cluster.submitTopology("word-count", conf, builder.createTopology());
    }
}
```

Listing 1: WordCountTopology.java

3.4 Storm Architecture

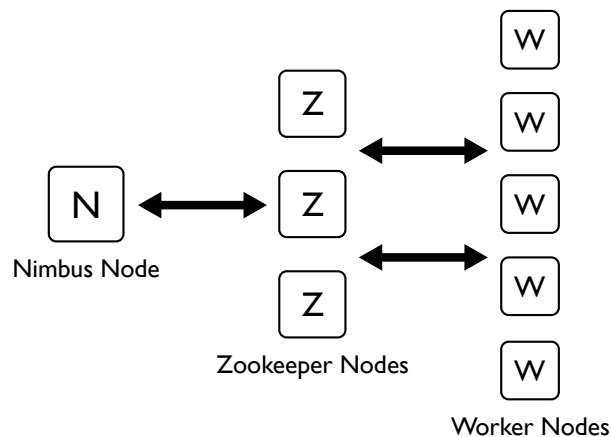


Figure 3.2: Apache Storm Architecture.

A Storm cluster adopts the Master-Worker pattern. To set up a Storm topology, the user launches daemon processes on nodes of the cluster and submits the topology to the master node, also called Nimbus. The worker nodes receive task assignments from the master and execute the tasks assigned to them. The coordination between the master node and the worker nodes is handled by nodes running Apache Zookeeper. Figure 3.2 shows a graphical representation of Storm Architecture.

3.4.1 Nimbus Node

The master node runs a server daemon called Nimbus. The main role of Nimbus is to receive topology submissions from clients. Upon receiving a topology

submission, Nimbus takes the following steps:

Validate the topology

The topology is validated using a validator to ensure that the submitted topology is valid before trying to execute it. The user can use his own validator by implementing the `ITopologyValidator` interface or use the default validator provided by Storm.

Distribute the topology source code

Nimbus ensures that the workers involved in the topology computation have the source code by sending it to all nodes of the cluster.

Schedule the topology

Nimbus runs a scheduler that distributes the work among workers of the cluster. Similarly to validation, the user can use his own scheduler by implementing the `IScheduler` interface or use the default scheduler provided by Storm. The default scheduler uses a simple Round-robin strategy [34].

Activate the topology

Nimbus transitions the topology to active state which tells the worker nodes to start executing it.

Monitor the topology

Nimbus continues to monitor the topology by reading heartbeats sent by the worker nodes to ensure that the topology is executing as expected and worker nodes have not failed.

Nimbus is an Apache Thrift [19] service (more on Thrift in section 3.5) that listens to commands submitted by clients and modifies the state of a cluster accordingly. Following are the commands supported by Nimbus:

Submit a topology

Clients can submit a topology defined in a Java Archive (JAR) file. The Nimbus service then ensures that the topology configuration and resources are distributed across the cluster and starts executing the topology as previously described.

Kill a topology

Nimbus can stop running a topology and remove it from the cluster. The cluster can still continue executing other topologies.

Activate/deactivate a topology

Topologies can be deactivated and reactivated by Nimbus. This could be useful if the spout temporarily cannot produce a stream and the user does not want the cluster to execute idly.

Rebalance a topology

Nimbus can rebalance a topology across more nodes. Thus if the number of nodes in the cluster ever changes the user can increase or decrease the number of nodes involved in the topology computation.

3.4.2 Worker Nodes

The worker nodes run a daemon called Supervisor. There are 4 layers of abstraction which control the parallelism of a worker node.

Supervisor

A supervisor is a daemon process the user runs on a worker node to make it part of the cluster. It launches worker processes and assigns them a port they can receive messages on. Furthermore, it monitors the worker processes and restarts them if they fail. A worker node runs only one supervisor process.

Worker

A worker process is assigned a port and listens to tuple messages on a socket associated with the port. A worker launches executor threads as required by the topology. Whenever it receives a tuple, it puts it on a receive queue of the target executor.

Furthermore, the worker has a transfer queue where its executors enqueue tuples ready to be sent downstream. There can be multiple worker processes running inside one supervisor.

Executor

An executor controls the parallelism within a worker process. Every executor runs in a separate thread. An executor's job is to pick up tuples from its receive queue, perform the task of a component it represents, and put the transformed tuples on the transfer queue of the worker. There can be many executors running inside one worker and an executor performs one (the usual case) or more tasks.

Task

A task represents the actual tuple processing function. However, within an executor thread all the tasks are executed sequentially. The main reason for having tasks is that the number of tasks stays the same throughout the lifetime of a topology but the number of executors can change (by rebalancing). Thus if some worker nodes in the cluster go down, the topology can continue executing with the same number of tasks as before.

3.4.3 Zookeeper Nodes

The Storm cluster contains a number of Zookeeper nodes which coordinate the communication between Nimbus and the worker nodes. Storm does this by storing the state of the cluster on the Zookeeper nodes where both Nimbus and worker nodes can access it.

The cluster state contains worker assignments, information about topologies, and heartbeats sent by the worker nodes to be read by Nimbus. Apart from the cluster state, Storm is completely stateless. Hence, if the master node or a worker node

fail the cluster continues executing and the node will get restarted if possible. The only time the cluster stops executing completely is if all the Zookeeper nodes die.

3.5 Serialisation

Since Storm topologies execute on a cluster all components need to be serialisable. This is achieved with Apache Thrift. Components are defined as Thrift objects and Thrift generates all the Java serialisation code automatically.

Furthermore, since Nimbus is a Thrift service Thrift generates all the code required for remote procedure call (RPC) support. This allows defining topologies in any of the languages supported by Thrift and easy cross-language communication with the Nimbus service.

Chapter 4

Bringing Storm to Multi-core

The following chapter explains how the Storm was ported to multi-core. We describe how Apache Storm behaves in a multi-core environment (4.1), discuss the design of Storm-MC (4.2), describe how Storm-MC was implemented (4.3), and list feature differences between Apache Storm and Storm-MC (4.4).

4.1 Apache Storm on Multi-core

To begin, we discuss why Apache Storm does not perform optimally in a multi-core environment. Storm can be ran in local mode where it emulates execution on a cluster. This mode exists so that it is possible to develop and debug topologies without needing access to a cluster. However, there are several reasons why the local mode is not as performant as it could be.

4.1.1 Tuple Processing Overhead

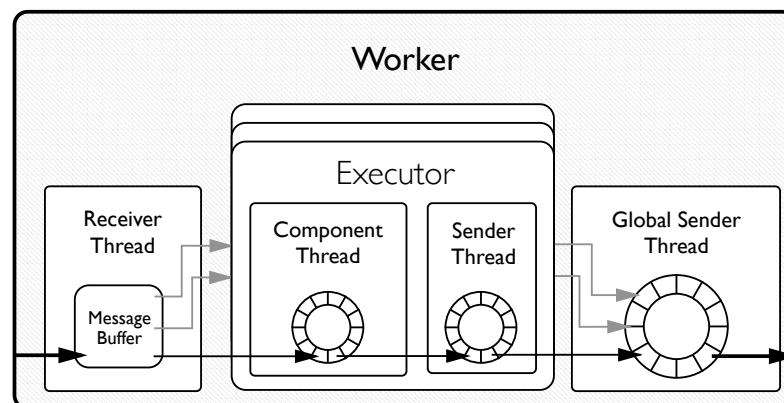


Figure 4.1: Tuple processing in Apache Storm.

Figure 4.1 shows how tuple processing is implemented inside a Storm worker process. The tuple is read from a message buffer by the receiver thread of the worker and put on a receive queue of the target executor. The tuple is then picked up by the component thread of the executor for task execution.

After the component thread has executed the task it puts the tuple on the executor send queue. There, it is picked up by the executor sender thread which puts the tuple on the global send queue of the worker. Finally, the global sender thread of the worker serialises the tuple and sends it downstream.

Alternatively, if the tuple is forwarded to an executor in the same worker process it is put on the receive queue of the corresponding executor directly after task execution.

The queues used in tuple processing are implemented as ring buffers using the Disruptor library [22]. Detailed background on how Disruptor works and its performance benchmarks can be found in [35]. In summary, due to less write contention, lower concurrency overhead, and being more cache-friendly the Disruptor pattern can offer latency of inter-thread messages lower than 50 nanoseconds and a throughput of over 25 million messages per second [35].

There is significant overhead required to emulate sending tuples to executors in other worker nodes. For one, there is the overhead from the tuple passing through the three queues of a worker. The authors of LMAX Disruptor showed that a three step pipeline can have half the throughput of a single consumer-producer pipeline [36].

Furthermore, to emulate over-the-network messages Storm uses a Hashmap of Linked-BlockingQueues which according to [35] has several orders of magnitude lower performance than the Disruptor.

4.1.2 Thread Overhead

Acker Bolt

The Acker bolt ensures that tuples propagate through the topology even if a failure in processing occurs. It provides the so-called at-least-once semantics. In Storm it is included in every topology. It can be disabled via the configuration file in which case it is mostly idle not receiving any tuples. However, it can still use up resources especially if it waits for tuples using a busy waiting strategy.

Heartbeats & Timers

Every worker has a heartbeat thread that simulates sending heartbeat messages to the Nimbus node. It does this by writing to a local cache which is persisted to a file by a write on every heartbeat. Since the write is implemented using the `java.io` package the write is blocking – the thread cannot continue until the write is completed. While heartbeats are essential in cluster mode to signal the node being alive, there is no need for them in local mode.

Zookeeper Emulation

More overhead is produced by a local Zookeeper server which emulates the Zookeeper nodes of a cluster. Running the Zookeeper server is a massive addition to the list of overheads as shown in the following paragraphs. The purpose of Zookeeper is to maintain state of running topologies and nodes of the cluster. As we will show in the following sections maintaining this state on multi-core is not necessary.

During profiling we found that a topology with one worker and three executors was being executed with 55 threads (not including system JVM threads and threads created by the profiler). Table 4.1 shows a breakdown of what the individual threads were used for.

| Spout Parallelism | # of Threads |
|----------------------------------|--------------|
| Main Thread | 1 |
| Worker Sender & Receiver Threads | 2 |
| Acker & System Component Threads | 2 |
| Executor Component Threads | 3 |
| Executor Sender Threads | 5 |
| Various Timers & Event Loops | 14 |
| Zookeeper Server | 28 |

Table 4.1: Storm thread usage: topology with three executors.

To find out what state the threads were actually in at any given time the topology was executed for three minutes and a JVM thread dump was recorded every second. The average results of this experiment can be observed in table 4.2 and the state distribution over time can be seen in Figure 4.2.

| Spout Parallelism | # of Threads |
|-------------------|--------------|
| RUNNABLE | 8 |
| TIMED WAITING | 22 |
| WAITING | 25 |

Table 4.2: Storm thread states: topology with three executors

Even though three minutes may seem to be a short amount of time the fact that there is almost no variation shows that it is sufficient. As can be seen from the table, most of the threads were either in state WAITING or TIMED WAITING. According to the Java documentation on thread states [37] these two states are used for threads that are waiting for an action from a different thread and cannot be scheduled by the scheduler until that action is executed.

On average there were eight threads in state RUNNABLE which JVM uses to mark threads which are executing on the JVM and are possibly waiting for resources

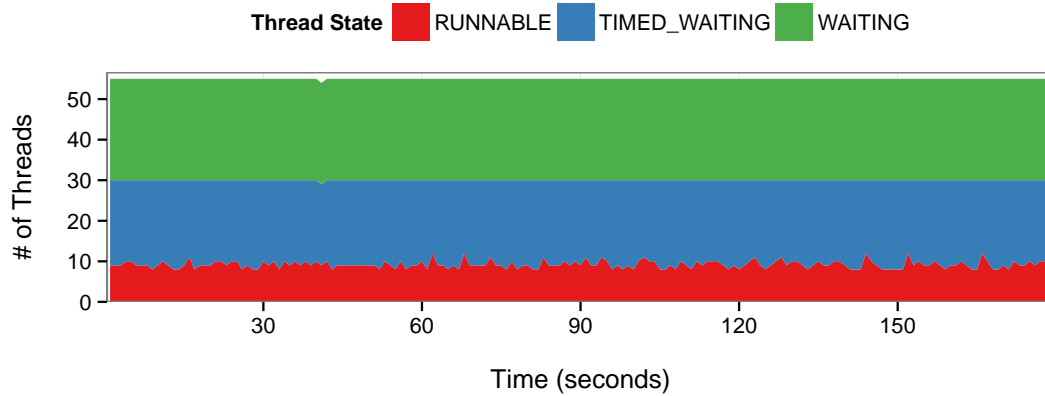


Figure 4.2: Thread state distribution over time.

from the operating system (OS) such as processor [37]. Hence, these are threads directly competing to be scheduled by the OS. This means that for three components running in parallel there are five threads doing potentially unnecessary work.

In the subsequent sections we will show that these threads were in fact unnecessary and we will discuss how the number of threads was reduced. In fact, to execute the same topology on Storm-MC requires only 5 threads.

4.2 Storm-MC Design

The design we adopted for porting worker nodes is to only have one worker executing all the executors of a topology. This design simplified the communication model and allowed removal of unnecessary abstractions.

Additionally, the source code for the Nimbus service was merged with the worker source code. This was done because there is no need to run Nimbus and worker specific code in parallel. Once the Nimbus source code sets up the topology, all the work is done by the worker source code. Hence they can be executed serially.

4.2.1 Porting Nimbus

Nimbus in Apache Storm performs as a server that clients can send topologies to for execution. In Storm-MC, we opted for a different design. Storm-MC is designed as a standard Java library that applications can import and use to create and execute topologies in the main method of their application. This means that Storm-MC does not support running multiple topologies at the same time. However, to do that one only needs to run the topology in a separate process. This is because, unlike when executing on a cluster, different topologies do not need to share any state and it is more natural to execute them as separate processes. This design decision

has the added benefit of each process having its own part of main memory thus reducing cache conflicts as shown in [39] and providing higher security by not having different topologies share memory space.

The interaction with the Nimbus service in Storm is usually through a shell script with a path to a JAR file of the topology and the main class to execute. This shell script was ported over to Storm-MC but instead of communicating with a service it spawns a new separate process that executes the topology.

Unlike Apache Storm, Storm-MC does not support topology scheduling. Since within one process there is always only one topology running at a time and the hardware configuration of the machine does not change, the parallelism is clearly defined by the number of executors per component specified by the topology configuration.

One way to implement scheduling could be to pin threads to specific cores. Unfortunately, Java does not provide support for CPU affinity; the assignments are handled automatically by the JVM. Potentially, this could be achieved by using C or C++, both of which support CPU affinity, but this was not implemented in Storm-MC.

The role of Nimbus in Storm-MC has effectively been reduced to validating the topology and its configuration and passing the topology along to the worker source code which handles topology execution.

4.2.2 Porting Worker Nodes

In Apache Storm, a worker node runs a supervisor daemon, which in turns launches worker processes which run executors which execute tasks. There are thus three ways to control the parallelism of a component in Apache Storm: setting the number of workers, setting the number of executors per component, and setting the number of tasks per executor.

In Storm-MC, however, there is only one worker wrapper which runs all executors and their tasks. Furthermore, only one task executes within an executor since tasks execute serially. Hence, the parallelism of a component is controlled by only one variable: number of executors per component. This represents the number of threads that will function as the component within a topology. This design has several benefits:

- All the communication occurs within one worker wrapper.
- The supervisor daemon can be removed as there is no need to synchronise or monitor workers.
- There is no need to simulate over-the-network message passing.
- Message passing between executor threads within a worker stays the same as in Apache Storm.

A comparison of an Apache Storm worker node and its Storm-MC equivalent is shown in Figure 4.3.

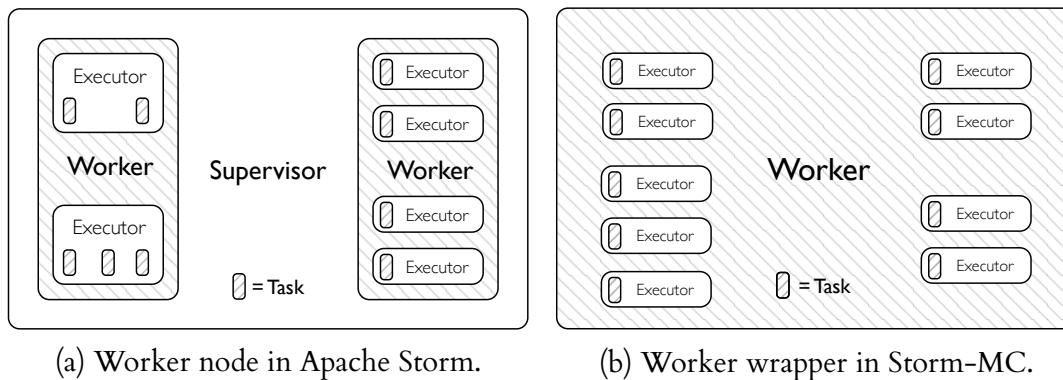


Figure 4.3: Comparison of a worker in Storm and Storm-MC

The role of the worker wrapper is to launch executors and provide them with a shared context through which they can communicate. This is done with a map of Disruptor queues which the executors use as receive queues to pick tuples from.

Moreover, the worker wrapper contains a map of components and streams. This map specifies which bolts subscribe to a stream a component produces. Executors use this map to figure out who they should send tuples to.

Additionally, the worker wrapper has a timer which components can use to get tick tuples at regular intervals. As mentioned before, bolts can use tick tuples to trigger events at regular intervals. For example, one might want to sort a window of tuples based on some criteria every five minutes.

mention metrics here?

Finally, a worker contains a configurable-size thread pool `ExecutorService`. Storm-MC executors can use this service via the topology context to launch background tasks on a shared thread pool.

4.2.3 Removing State

Storm-MC is completely stateless. The cluster state that was managed by Zookeeper in Apache Storm was completely stripped away. In Storm, workers use the Zookeeper cluster state to communicate with Nimbus and vice versa. For example, when Nimbus creates topology assignments it informs workers via the cluster state. In Storm-MC, we adopted a more functional approach where worker is just a function invoked by the Nimbus part of the source code.

While this is not something that is visible to a user of Storm-MC, it required a great effort as all the code that referenced the Zookeeper library had to be refactored.

4.2.4 Removing Serialisation

Similarly to removing the Zookeeper state, great amount of work was put into removing the dependency of Storm-MC on Apache Thrift. This was mostly done to reduce code bloat and remove an unnecessary dependency since there is no serialisation required in a multi-core environment.

Moreover, code generated by Thrift does not use standard Java camelCase naming conventions but instead uses underscore_case. For example, Thrift generates method names such as `get_component_common`.

Removing Thrift required refactoring all the data types generated automatically by Thrift. This not only reduced the size of the codebase significantly but also made the code more readable and self-documenting than the code generated by Thrift.

4.3 Implementation Details

Most of the implementation of Storm-MC was ported over from Apache Storm with adjustments where necessary. The problem with describing implementation of ported software is that there is a lot of functionality that needed to be changed but the changes required were usually not substantial. This is the case with Storm-MC as well. This section tries to outline the implementation of Storm-MC.

4.3.1 Topology Submission

Topologies are built using an instance of the `TopologyBuilder` class which uses the builder pattern – same as in Apache Storm. While the basis of this class was reused from Storm, the internals had to be refactored so they work with data types used by Storm-MC. Once a topology is built it is submitted to an instance of `LocalCluster` class. This class is used in Storm for emulating the cluster on a local machine. Storm-MC adapted the class for backwards compatibility. This way, code created for Storm needs minimal adaptation to work on Storm-MC. A topology is submitted for execution via the `submitTopology` method which takes three arguments: the name of the topology, a Java Map with configuration, and a topology built by `TopologyBuilder`.

Users can submit a configuration file written in YAML. This is done by setting a JVM property called `storm.conf.file`. This file can for example define the capacity of the Disruptor queues, what waiting strategy should the components employ when there are no tuples to pick up, and what hooks they want executed every time a tuple is processed. Furthermore, users can define their own topology validator by implementing the `ITopologyValidator` interface.

4.3.2 Tuple Processing

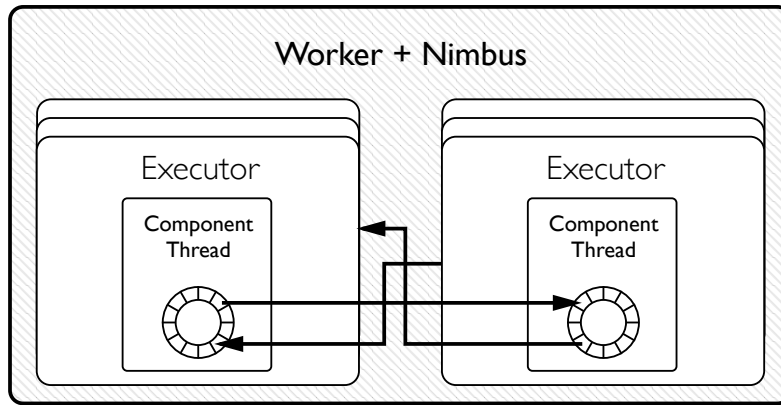


Figure 4.4: Tuple processing in Storm-MC.

The implementation of tuple processing in Storm-MC is depicted in Figure 4.4. As can be seen from the figure, the queues used for remote message sending present in Storm were stripped away and there is only one Disruptor queue for every executor. Once an executor is done processing a tuple it puts it directly on the receive queue of its downstream bolts.

In Apache Storm, every executor runs two threads, one for tuple processing and one for tuple sending. In Storm-MC, however, there is only one thread per executor. Hence, the number of queues a tuple needs to pass through in a component is lowered as compared to Apache Storm.

Tuple processing in Storm-MC is a variant of multiple producer single consumer problem. In general, multiple producer single consumer problems are hard to optimise since there needs to be some form of synchronisation between multiple producers trying to produce an entry at the same time. In Disruptor queues this is implemented using the atomic Compare-and-Swap (CAS) operation. This operation ensures that even if multiple threads attempt to modify a variable only one of them succeeds and all threads involved are able to tell whether it was them that succeeded. Hence, one thread will succeed at claiming the next entry of a Disruptor queue and others will have to retry.

Alternatively, locks can be used to synchronise access but lock-free queues using CAS are considered to be more efficient than locks because they do not require a kernel context switch [35]. However, even with CAS a processor must lock its instruction pipeline to ensure atomicity and employ a memory barrier to make the changes visible to other threads.

We investigated other data structures besides Disruptor queues that could be used for tuple exchange in Storm-MC but Disruptor queues are considered state of the art in low-latency parallel systems. This is one area that Apache Storm did really well and we were not able to further optimise.

Other options we considered were `ArrayBlockingQueue` and `LinkedBlockingQueue`

both of which are in the Java standard library. However, the Disruptor shows superior throughput and latency compared to these options as shown in [36].

4.3.2.1 Waiting Strategies

There are four different waiting strategies an executor can employ while waiting for a tuple to become available:

BlockingWaitStrategy

This strategy uses a lock and a condition variable. The thread waits on the condition variable and is signalled once a new tuple becomes available. This strategy wastes the minimum number of CPU cycles when an executor is waiting.

SleepingWaitStrategy

This strategy initially spins for hundred iterations, then uses `Thread.yield()`, and finally uses `LockSupport.parkNanos(1L)` to sleep. Thus after quiet periods this strategy might introduce latency spikes.

YieldingWaitStrategy

This strategy initially spins for hundred iterations and then uses `Thread.yield()`. This strategy is a good compromise between CPU utilisation and great performance.

BusySpin WaitStrategy

In this strategy the thread is in a so-called tight loop, where it checks whether a new entry is available and only breaks out of the loop if there is a new entry. This strategy has the best performance but works well only if the number of CPU cores is higher than the number of active threads since it maximises CPU utilisation.

The default strategy used in Storm-MC is `BlockingWaitStrategy` but users can change the strategy in the configuration file.

4.3.2.2 Tuple Pools

Once a component wants to send a new tuple to its downstream components it needs to initialise a Java Tuple object. Here, we saw room for improvement since this might need to be done at very high rates, possibly million times per second.

Hence, a tuple pool was implemented where executors could place tuples after they were done with them so the tuples could be reused by other executors. However, access to this pool also had to be synchronised and hence accessing the pool introduced higher latency than simply initialising new tuples. Moreover, Java garbage collector actually does a good job of re-claiming unused memory and thus the idea of using a tuple pool was abandoned.

This problem could potentially be circumvented by using different constructs to achieve tuple passing but backwards compatibility with programs written for Apache Storm was deemed more important than the potential gain in performance.

4.3.3 Executors

The implementation of an executor processing a tuple is as follows:

1. The executor repeatedly tries to pick a tuple from its Disruptor queue. If there are no tuples to be picked up it employs a waiting strategy between trials.
2. Once it picks up a tuple, it processes it as per the component it represents.
3. If it emits a new tuple it attempts to send it downstream by repeatedly trying to claim an entry of the downstream ring buffer.
4. Once it successfully claims an entry it proceeds back to step 1.

4.4 Differences between Apache Storm and Storm-MC

The codebase of Apache Storm is fairly large – 54,985 lines of code as reported by cloc [38]. Thus we had to prioritise features that were ported over to Storm-MC. Table 4.3 presents a list of Storm features and shows which were ported over to Storm-MC and which were not.

| Feature | Apache Storm | Storm-MC |
|---------------------------|--------------|----------|
| Multi-language Topologies | ✓ | ✓ |
| Hooks | ✓ | ✓ |
| Metrics | ✓ | ✓ |
| Tick Tuples | ✓ | ✓ |
| Multiple Topologies | ✓ | ✗ |
| Topology Scheduling | ✓ | ✗ |
| Trident API | ✓ | ✗ |
| Built-in Metrics | ✓ | ✗ |
| Nimbus as a Server | ✓ | ✗ |

Table 4.3: Feature comparison of Apache Storm and Storm-MC.

A feature that we deemed very important is support for multi-language topologies. Thus, Storm-MC allows you to define components in other languages such

as Python or Ruby and connect them into a Java topology. An example of a component defined in Python can be seen in Listing 4.

Storm-MC has support for task hooks just like Storm. Task hooks allow you to capture a number of events and execute custom code when the event occurs at a registered component. Hooks are created by subclassing `BaseTaskHook`. They can, for example, be used to update a web server with the latest performance metrics.

Additionally, Storm-MC has support for topology metrics. This way, components can record metrics such as number of tuples processed or a count of event occurrences. These metrics can then be automatically consumed by a bolt that subclasses the `MetricsConsumerBolt` class.

As mentioned before, Storm-MC supports tick tuples which can be used to trigger component-local events at regular intervals.

Apache Storm supports an alternative high-level API called Trident which then gets converted into spouts and bolts by the Storm library. Trident was omitted from Storm-MC but it would be possible to implement it on top of the current API.

Moreover, Apache Storm collects JVM metrics with a bolt called `SystemBolt`. This bolt is added automatically to all Storm topologies. This bolt is not included in Storm-MC topologies automatically but users can choose to add this bolt on their own.

Chapter 5

Evaluation

This chapter evaluates Storm-MC. We describe the metrics used to evaluate performance of Storm-MC (5.1), list the configuration used for benchmarking (5.2), compare Storm-MC to Apache Storm executing in local mode on a set of different topologies (5.3), and finally talk about challenges encountered while designing Storm-MC (5.4).

5.1 Evaluation Metrics

The system was evaluated on the following metrics:

Throughput

The number of tuples processed by every component of the topology in the given time is recorded.

CPU Utilisation

The CPU utilisation is recorded every second throughout program execution and the average is computed.

Resident Memory Size

The resident size of main memory is recorded every second throughout program execution and the average is computed.

5.2 System Configuration

5.2.1 Software Setup

All performance benchmarks were ran using the following software packages:

- Apache Storm version 0.9.2 ¹

¹<https://github.com/mrknmc/storm/releases/v0.9.2-fix>

- Storm-MC version 0.1.6 ²
- A fork of IBM Storm Email Benchmarks version 0.1.12 ³
- Storm-benchmark version 0.1.0 ⁴

The Apache Storm source code had to be adapted to include a workaround for a deadlock bug present in version 0.9.2. This bug caused a topology to exit with threads left in Zombie state under certain conditions. This prevented Storm from logging the benchmark metrics after execution. Hence a workaround was added so the results were logged.

Version 0.1.6 is the latest version of Storm-MC as of this moment. The first release was version 0.1.0 which was production-ready but since then there were 6 minor versions fixing bugs as they were discovered during testing.

IBM open sourced a suite of benchmarks which they used to compare Apache Storm to their real-time stream system IBM InfoSphere Streams [40]. These benchmarks were adapted and used to benchmark Storm-MC against Apache Storm. The current version is 0.1.12.

Finally, a number of spout and bolt components were used from the storm-benchmark project which Apache Storm developers use to benchmark Storm.

Since Storm-MC reuses package names from Apache Storm, the same benchmark is directly executable by both libraries. This saved a lot of time and furthermore there is no need to maintain two benchmarks suites. An example topology submission to Storm-MC and Apache Storm, respectively would look as follows:

```
java -Dstorm.home=storm-mc -cp storm-multicore-0.1.6.jar
...
-Dstorm.jar=storm-email-benchmark.jar
com.ibm.streamsx.storm.email.benchmark.FileReadWordCount wordcount
java -Dstorm.home=apache-storm -cp storm-core-0.9.2-incubating.jar
...
-Dstorm.jar=storm-email-benchmark.jar
com.ibm.streamsx.storm.email.benchmark.FileReadWordCount wordcount
```

In this example the Storm home directory was set, the corresponding library was added to the Java classpath, the JAR file containing the benchmark was submitted, and the main class and topology name were specified.

5.2.2 Hardware Setup

The machine used for benchmarking is the Informatics Student Compute server (student.compute.inf.ed.ac.uk). The server has the following hardware components:

²<https://github.com/mrknmc/storm-mc/releases/0.1.6>

³<https://github.com/mrknmc/benchmarks/releases/0.1.12>

⁴<https://github.com/manuzhang/storm-benchmark>

Processor: Intel® Xeon® E5-2690 v2 @ 3.00 GHz

The machine has two sockets with the same processor each. The processor has 10 physical cores with Hyper-Threading Technology which means it can handle up to 20 threads in parallel. Thus with two sockets, there is potential to execute up to 40 threads in parallel.

Main Memory

The machine has 378 GB of main memory. Since data stream processing uses windows to store only up to a certain amount of memory this was more than enough to conduct the benchmarks.

5.2.3 Storm Configuration

As mentioned before, when submitting a topology the programmer needs to submit a configuration file as well. To ensure that the performance difference between Apache Storm and Storm-MC was not caused by different configuration, the default configuration file from Storm 0.9.2 was used to benchmark both projects. Most notably, the size of the ring buffer used by executors (`topology.executor.receive.buffer.size`) was set to 1024 and the wait strategy employed by executors when there are no tuples to pick up (`topology.disruptor.wait.strategy`) was set to `BlockingWaitStrategy`.

5.3 Results

To assess performance of Storm-MC, 3 different benchmarks were executed, each with a different focus. The benchmarks were executed for a constant period of time – five minutes – after which the system was killed and metrics were collected. To avoid any performance differences caused by varying amounts of heap memory required by the tested systems, the programs were run with the following flag: `-Xmx10240M`. This flag sets the maximum amount of heap memory used by the JVM to 10 GB which was more than enough for all benchmarks.

The parallelism of components was varied from one to six and average CPU utilisation and resident memory size were recorded by the Unix `top` program [41]. Maximum CPU utilisation with 40 threads is 4,000%. Resident memory size is the amount of non-swapped physical memory a task has used. This metric can be deceiving as it depends on how OS manages memory but it is the only fairly reliable memory metric reported by `top` that can be used for Java programs. We only employ this metric to show proportional difference in memory usage between Storm and Storm-MC.

5.3.1 WordCount Topology

The first topology tested for performance is a variant of the aforementioned WordCount topology. This topology has a spout `FileReadSpout` generating random sentences, which sends messages to a `SplitSentenceBolt` bolt which splits the sentences on whitespace and sends individual words to a `CountBolt` which counts word frequencies. Recall, that this topology is shown graphically in Figure 3.1. Since the components do not store any data in memory or make any I/O calls this topology is mostly CPU-bound.

The number of tuples processed by each component in Storm-MC and Apache Storm is shown in tables 5.1 and 5.2, respectively. As can be seen from the tables, not only was CPU utilisation in Storm-MC lower, Storm-MC often processed more than twice as many tuples per component than Apache Storm. The number of tuples processed by `CountBolt`, the last component of the topology, is also shown in Figure 5.1. Since this topology is serial, the number of tuples processed by `CountBolt` is a good indicator of total throughput.

| Parallelism | FileReadSpout | SplitSentenceBolt | CountBolt | CPU Utilisation | Resident Size |
|-------------|---------------|-------------------|-------------|-----------------|---------------|
| 1 | 25,767,502 | 25,767,502 | 225,815,174 | 217.9% | 690.8M |
| 2 | 34,403,678 | 34,403,127 | 301,493,247 | 414.6% | 759.1M |
| 3 | 45,731,188 | 45,732,988 | 400,767,999 | 611.5% | 798.4M |
| 4 | 52,285,327 | 52,283,540 | 458,187,555 | 805.5% | 804.1M |
| 5 | 55,326,941 | 55,325,167 | 484,844,652 | 998.7% | 806.0M |
| 6 | 56,747,319 | 56,744,629 | 497,285,149 | 1,195.3% | 824.8M |
| 10 | 40,341,798 | 40,336,962 | 353,490,567 | 1,967.4% | 2.7G |

Table 5.1: Storm-MC: Component throughput in WordCount Topology.

| Parallelism | FileReadSpout | SplitSentenceBolt | CountBolt | CPU Utilisation | Resident Size |
|-------------|---------------|-------------------|-------------|-----------------|---------------|
| 1 | 12,583,377 | 12,579,132 | 110,233,966 | 294.5% | 2.2G |
| 2 | 16,800,475 | 16,796,695 | 147,194,709 | 481.7% | 2.8G |
| 3 | 22,120,695 | 22,107,696 | 193,735,106 | 687.1% | 2.6G |
| 4 | 20,720,637 | 20,711,756 | 181,500,586 | 895.3% | 2.6G |
| 5 | 17,177,688 | 17,164,209 | 150,412,037 | 1,129.3% | 2.5G |
| 6 | 17,402,418 | 17,388,691 | 152,374,303 | 1,342.1% | 2.3G |
| 10 | 12,229,523 | 12,211,100 | 107,002,632 | 2,136.7% | 2.8G |

Table 5.2: Apache Storm: Component throughput in WordCount Topology.

Furthermore, it can be seen that after the parallelism is increased beyond three the throughput of Apache Storm tails off and starts going down. This can be attributed to the number of threads ran by Apache Storm. For Storm-MC this tailing off oc-

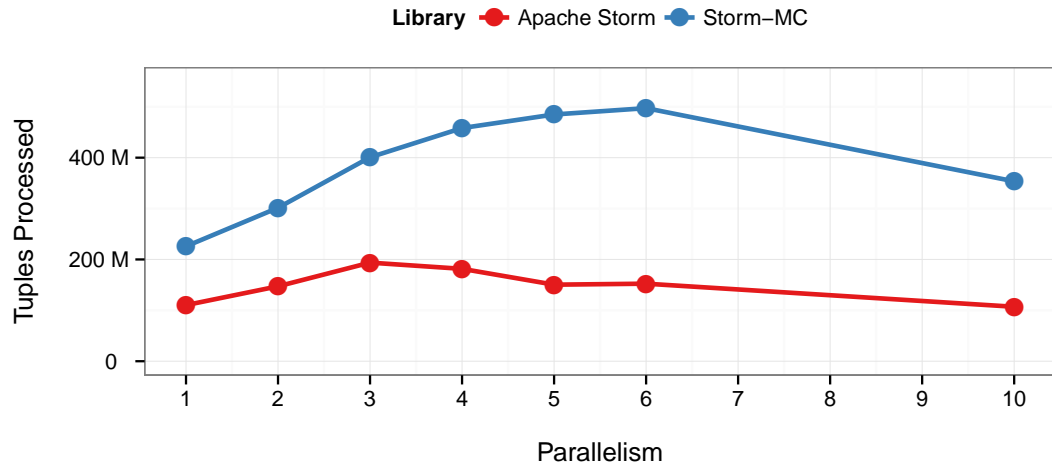


Figure 5.1: CountBolt throughput in Apache Storm and Storm-MC

curs with parallelism of six where the overhead of multiple producers possibly trying to publish to the same queue becomes apparent. Moreover, with parallelism set to 6, Storm-MC executes with 20 threads which is close to the number of physical cores of the machine. However, it should be noted that even with parallelism equal to 10, Storm-MC still processes more than three times as many tuples as Storm.

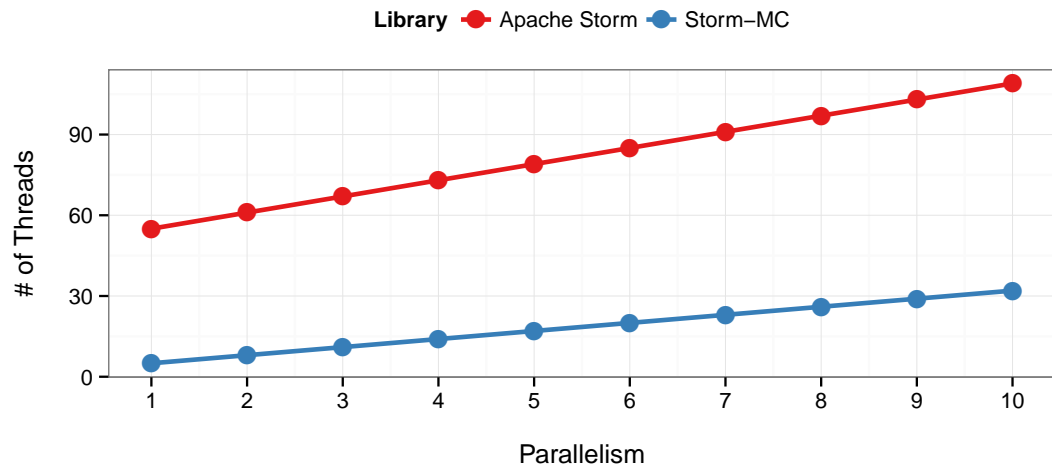


Figure 5.2: Number of threads used by Apache Storm and Storm-MC

The number of threads required to execute a topology is a linear function of the parallelism for both Storm and Storm-MC. However, as shown in Figure 5.2, the number of threads required by Storm increases more rapidly than Storm-MC. For example, with parallelism set to 10, Storm creates 109 threads whereas Storm-MC creates only 32. More formally, the number of threads required by both systems can be expressed as:

$$\begin{aligned}
49 + 2 \times \sum_c^{\text{components}} \text{parallelism}(c) & \quad \text{for Apache Storm.} \\
2 + \sum_c^{\text{components}} \text{parallelism}(c) & \quad \text{for Storm-MC.}
\end{aligned}$$

Figure 5.3: Number of threads used by Apache Storm and Storm-MC.

N.B.: This is a general formula that applies to all topologies, not just WordCount.

Of note, the resident size used by Storm-MC is also less than half of the resident size used by Storm for cases with parallelism less than 10.

5.3.2 Enron Topology

Next, Enron topology from the IBM benchmarks was tested for performance. In this topology, serialised emails from the Enron email dataset are read from a file by a `ReadEmailsDecompressSpout` spout. They are then deserialised by `AvroDeserializeBolt` bolt, filtered by `NewFilterBolt` bolt, modified by `ModifyBolt` bolt, and finally metrics are recorded by a `NewMetricsBolt` bolt. Additionally, every instance of the `NewMetricsBolt` bolt sends its local average email throughput to a global (excluded from the parallelism setting) `GlobalMetricsBolt` bolt every four seconds. This bolt then records the global average email throughput. Figure 5.4 shows this topology graphically.

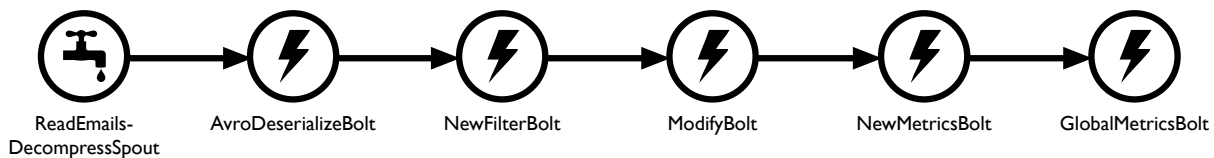


Figure 5.4: Enron topology.

Similarly to the WordCount topology, this topology is serial in nature. However, whereas the spout in WordCount topology keeps a small number of sentences in memory, the Enron topology has a spout that produces tuples by reading from a file. Thus, this benchmark is mostly I/O-bound. The average email throughput in Storm-MC and Apache Storm is shown in tables 5.3 and 5.4, respectively.

As can be seen from the tables, the difference in throughput in Enron Topology is less staggering than in WordCount. This is due to the fact that the throughput is limited by the file reads of the spout. However, as the parallelism increases the improvement in throughput of Storm-MC becomes more apparent as shown in Figure 5.5. This figure also shows that the throughput is fairly volatile. This is due to the fact that the file is loaded into main memory in chunks and hence the throughput drops when the spout is trying to read from a file in between the loads.

| Parallelism | Emails Processed | CPU Utilisation | Resident Size |
|-------------|------------------|-----------------|---------------|
| 1 | 3,285,742 | 297.7% | 806.8M |
| 2 | 6,696,612 | 482.1% | 756.1M |
| 3 | 8,493,772 | 729.5% | 341.4M |
| 4 | 11,102,969 | 1036.9% | 326.0M |
| 5 | 12,630,475 | 1311.0% | 260.8M |
| 6 | 14,082,501 | 1590.3% | 334.0M |

Table 5.3: Storm-MC: Email Throughput in Enron Topology.

| Parallelism | Emails Processed | CPU Utilisation | Resident Size |
|-------------|------------------|-----------------|---------------|
| 1 | 2,943,709 | 406.6% | 1.94G |
| 2 | 4,832,874 | 945.1% | 2.93G |
| 3 | 5,623,028 | 1,427.4% | 3.32G |
| 4 | 6,238,395 | 1,891.2% | 3.56G |
| 5 | 6,105,155 | 2167.4% | 3.65G |
| 6 | 7,242,298 | 2388.6% | 4.09G |

Table 5.4: Apache Storm: Email Throughput in Enron Topology.

As before, the resident size used by Storm-MC is significantly lower than that of Storm.

5.3.3 RollingSort Topology

The RollingSort topology was ported over from the aforementioned storm-benchmark project. This topology only includes one spout sending tuples to one bolt. The RandomMessageSpout spout produces hundred character-long strings of random digits from zero to eight. The SortBolt bolt then stores a rolling window of hundred of such messages and sorts them every 10 seconds. The graphic representation of this topology can be seen in Figure 5.6.

This benchmark is considered to be memory-bound: the bolt stores a window of tuples in memory and performs a non-linear time sort. The results of running this benchmark on Storm-MC and Apache Storm can be seen in tables 5.5 and 5.6, respectively.

While Storm-MC still outperforms Apache Storm, the difference in performance is marginal. This is due to the fact that the topology only has two components and the application is mostly memory bound. Storm-MC provides speed improvements for topologies that are mostly CPU-bound and have several components working serially such as WordCount. Furthermore, Storm-MC beats Storm significantly when the parallelism is high, as in previous topologies.

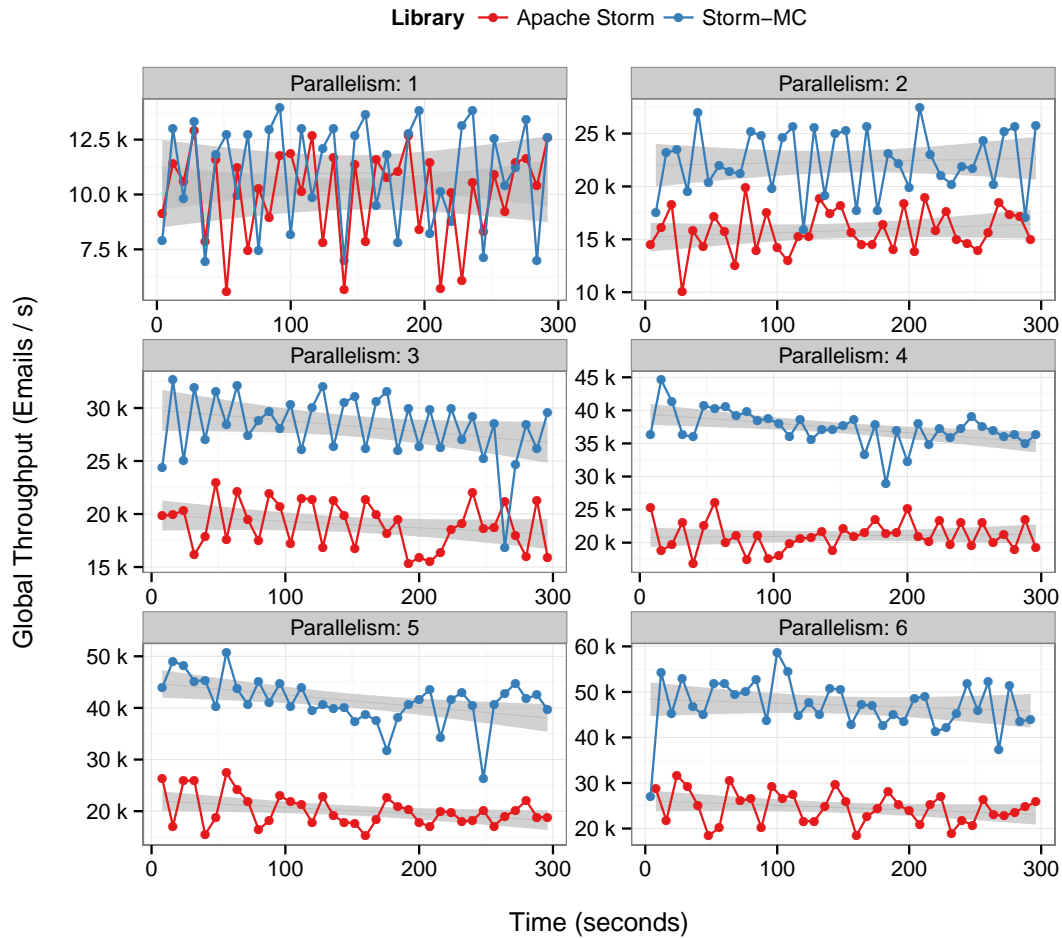


Figure 5.5: Global email throughput over time with standard error.

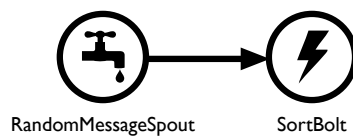


Figure 5.6: RollingSort topology.

5.4 Challenges

In this section we discuss challenges we encountered while porting Apache Storm to multi-core. We also try to provide a critical analysis of the project.

Backwards Compatibility

A big challenge while working on this project was ensuring that the final system is backwards compatible with Apache Storm. Doing this ensured that existing applications developed for Apache Storm can be executed with Storm-MC. On the other hand, this was sometimes limiting the possible performance improvements.

| Parallelism | RandomMessageSpout | SortBolt | CPU Utilisation | Memory Usage |
|-------------|--------------------|-------------|-----------------|--------------|
| 1 | 249,143,444 | 249,142,400 | 186.2% | 504.3M |
| 2 | 444,261,351 | 444,259,400 | 352.0% | 401.7M |
| 3 | 350,861,061 | 350,859,800 | 514.7% | 382.9M |
| 4 | 412,429,850 | 412,428,600 | 675.2% | 314.2M |
| 5 | 470,813,184 | 470,811,300 | 835.8% | 423.2M |
| 6 | 498,957,255 | 498,954,600 | 989.6% | 235.1M |
| 7 | 519,744,352 | 519,741,700 | 1,149.0% | 637.9M |
| 8 | 532,285,376 | 532,283,800 | 1,302.9% | 618.1M |
| 9 | 501,519,539 | 501,517,700 | 1,430.4% | 579.4M |
| 10 | 555,468,830 | 555,467,000 | 1,651.6% | 564.7M |

Table 5.5: Storm-MC: Component throughput in RollingSort Topology.

| Parallelism | RandomMessageSpout | SortBolt | CPU Utilisation | Memory Usage |
|-------------|--------------------|-------------|-----------------|--------------|
| 1 | 173,906,935 | 173,900,300 | 267.3% | 3.0G |
| 2 | 226,583,924 | 226,579,200 | 468.3% | 3.0G |
| 3 | 310,949,455 | 310,943,000 | 634.6% | 2.9G |
| 4 | 362,675,336 | 362,663,600 | 815.2% | 2.8G |
| 5 | 409,470,032 | 409,462,100 | 969.4% | 2.7G |
| 6 | 435,471,042 | 435,459,600 | 1,139.6% | 2.6G |
| 7 | 395,386,336 | 395,309,900 | 1,327.7% | 2.6G |
| 8 | 366,680,402 | 366,553,300 | 1,509.1% | 2.7G |
| 9 | 359,091,633 | 358,942,200 | 1,689.5% | 2.7G |
| 10 | 313,912,451 | 313,811,300 | 1,889.5% | 2.7G |

Table 5.6: Apache Storm: Component throughput in RollingSort Topology.

Unfamiliarity with Clojure

One of the main challenges while working on this project was learning a new programming language - Clojure. Since most of the implementation of Apache Storm is written in Clojure, this language had to be studied and its concepts well understood for us to be able to write code that worked with the existing codebase. By the end of the project writing Clojure has become second nature to us but initially progress was slow.

Lack of Documentation

Even though Apache Storm is a popular project documentation is available only for the high level concepts used within Storm. The implementation details are often obscured away in hard to understand functions. Since the documentation is lacking our knowledge of Storm had to be obtained by reading the source code of an initially unfamiliar language. By the end of the

project the Storm-MC code became well documented and we might attempt back-porting to Apache Storm.

Chapter 6

Conclusion

This final chapter concludes with a summary of contributions of this project (6.1) and discusses future work that could stem from this project (6.2).

6.1 Summary of Contributions

The primary contribution of this project is Storm-MC – a library aimed at data stream processing applications. The benefits of using Storm-MC are twofold:

- It offers the same easy-to-use API as Apache Storm.
- It is tailored to multi-core environments.

Since Storm-MC uses the same API as Apache Storm, applications written with Storm in mind can be ported to use Storm-MC with minimum amount of effort. Thus if an application requires parallelism satisfiable by a single multi-core machine, it can be executed on one machine instead of a cluster.

Moreover, the Storm API allows programmers to create data stream processing applications on multi-core with an unprecedented ease. All of this comes with the superior performance Storm-MC offers compared to running Apache Storm in local mode, as shown in Section 5.3.

6.2 Future Work

Storm-MC could be improved in a number of ways:

Storm-MC as a Server

Storm-MC could be updated to allow server-like execution. This could have several benefits such as being able to execute multiple topologies at the same time with a thin wrapper that could control their execution just like with

Apache Storm. This was not implemented as part of this project as we assumed most of the time users are executing only one topology per machine.

Higher Level Abstractions

Defining components of a Storm-MC topology is fairly simple. Users only need to define how components are connected, how they process tuples, and what tuples they emit. However, this could be taken even further with the user only specifying high-level functions and the Storm-MC library figuring out how to distribute the work. In Apache Storm this is implemented in Trident which was not ported as part of this project.

Automatic Parallelism

Sometimes when configuring a topology it may be difficult to predict the rate at which spouts are going to produce tuples. If the rate is underestimated consumers could be lagging behind producers. On the other hand, if the rate is overestimated consumers could be idle, not doing any useful work. Thus it might be advantageous to have an automatic parallelism setting which could add or remove consumers based on the current tuple rate.

It may seem that this would be trivial to implement with a pool of threads representing one component. However, there are several problems that need to be considered. For example, fields grouping guarantees that tuples with the same field values go to the same executor. Changing the parallelism at runtime breaks this guarantee.

Alternatively each executor could use a pool of threads. This comes with its own set of problems: the executor object would have to provide synchronised access to the pool which would only increase latency.

Performance Comparison with Distributed Storm

The benchmarks in this report compared Storm-MC to Apache Storm running in local mode. It would be interesting to see how Storm-MC compares to Apache Storm running on a cluster. One could compare the number of nodes required in a cluster to the number of cores required in a multi-core server to achieve certain throughput for a given topology. This could provide insight into when it becomes advantageous to deploy the topology to a cluster.

Appendices

Appendix A

Listings

```

1  public class RandomSentenceSpout extends BaseRichSpout {
2      SpoutOutputCollector _collector;
3      Random _rand;
4
5      public void open(Map conf, TopologyContext context,
6          SpoutOutputCollector collector) {
7          _collector = collector;
8          _rand = new Random();
9      }
10
11     public void nextTuple() {
12         Utils.sleep(100);
13         String[] sentences = new String[]{
14             "the cow jumped over the moon",
15             "an apple a day keeps the doctor away",
16             "four score and seven years ago",
17             "snow white and the seven dwarfs",
18             "i am at two with nature" };
19         String sentence = sentences[_rand.nextInt(sentences.length)];
20         _collector.emit(new Values(sentence));
21     }
22
23     public void ack(Object id) {}
24
25     public void fail(Object id) {}
26
27     public void declareOutputFields(OutputFieldsDeclarer declarer) {
28         declarer.declare(new Fields("word"));
29     }
30
31 }

```

Listing 2: RandomSentenceSpout.java: Definition of a spout that emits a randomly chosen sentence from a predefined collection of sentences.

```

1 public static class SplitSentence extends ShellBolt
2     implements IRichBolt {
3
4     public SplitSentence() {
5         super("python", "splitsentence.py");
6     }
7
8     public void declareOutputFields(OutputFieldsDeclarer declarer) {
9         declarer.declare(new Fields("word"));
10    }
11
12    public Map<String, Object> getComponentConfiguration() {
13        return null;
14    }
15 }

```

Listing 3: SplitSentence.java: Definition of a bolt that executes a Python script.

```

1 import storm
2
3
4 class SplitSentenceBolt(storm.BasicBolt):
5
6     def process(self, tup):
7         words = tup.values[0].split(" ")
8         for word in words:
9             storm.emit([word])
10
11
12 SplitSentenceBolt().run()

```

Listing 4: splitsentence.py: Definition of a bolt that splits sentences on whitespace in Python.

```
1 public static class WordCount extends BaseBasicBolt {  
2     Map<String, Integer> counts = new HashMap<String, Integer>();  
3  
4     public void execute(Tuple tuple, BasicOutputCollector collector) {  
5         String word = tuple.getString(0);  
6         Integer count = counts.get(word);  
7         if (count == null)  
8             count = 0;  
9         count++;  
10        counts.put(word, count);  
11        collector.emit(new Values(word, count));  
12    }  
13  
14    public void declareOutputFields(OutputFieldsDeclarer declarer) {  
15        declarer.declare(new Fields("word", "count"));  
16    }  
17 }
```

Listing 5: WordCount.java: Definition of a bolt that counts word frequencies.

Bibliography

- [1] J. Dean and S. Ghemawat, “Mapreduce: A flexible data processing tool,” *Communications of the ACM*, vol. 53, no. 1, pp. 72–77, 2010.
- [2] (2015). Apache Hadoop, [Online]. Available: <https://hadoop.apache.org> (visited on 03/25/2015).
- [3] (2015). Apache Storm, [Online]. Available: <https://storm.apache.org> (visited on 03/20/2015).
- [4] E. Solovey. (2015). Handling five billion sessions a day – in real time, [Online]. Available: <https://blog.twitter.com/2015/handling-five-billion-sessions-a-day-in-real-time> (visited on 03/20/2015).
- [5] K. A. Kumar, J. Gluck, A. Deshpande, and J. Lin, “Hone: ”scaling down” hadoop on shared-memory systems,” *Proc. VLDB Endow.*, vol. 6, no. 12, pp. 1354–1357, Aug. 2013, ISSN: 2150-8097. DOI: 10.14778/2536274.2536314. [Online]. Available: <http://dx.doi.org/10.14778/2536274.2536314>.
- [6] B. Wheeler. (Jul. 2011). Tiler sees opening in clouds, [Online]. Available: http://www.linleygroup.com/newsletters/newsletter_detail.php?num=4732 (visited on 03/26/2015).
- [7] P. Lotfi-Kamran, B. Grot, M. Ferdman, S. Volos, Y. O. Koçberber, J. Picorel, A. Adileh, D. Jevdjic, S. Idgunji, E. Özer, and B. Falsafi, “Scale-out processors,” in *39th International Symposium on Computer Architecture (ISCA 2012), June 9–13, 2012, Portland, OR, USA*, IEEE Computer Society, 2012, pp. 500–511, ISBN: 978-1-4673-0475-7. DOI: 10.1109/ISCA.2012.6237043. [Online]. Available: <http://dx.doi.org/10.1109/ISCA.2012.6237043>.
- [8] (2014). Apache Storm: Powered By, [Online]. Available: <https://storm.apache.org/documentation/Powered-By.html> (visited on 03/25/2015).
- [9] L. Golab and M. T. Özsu, “Issues in data stream management,” *ACM Sigmod Record*, vol. 32, no. 2, pp. 5–14, 2003.
- [10] (2015). Apache Spark, [Online]. Available: <https://spark.apache.org> (visited on 03/20/2015).
- [11] (2015). Apache S4, [Online]. Available: <http://incubator.apache.org/s4> (visited on 03/20/2015).
- [12] M. M. Gaber, A. Zaslavsky, and S. Krishnaswamy, “Mining data streams: A review,” *SIGMOD Rec.*, vol. 34, no. 2, pp. 18–26, Jun. 2005, ISSN: 0163-5808. DOI: 10.1145/1083784.1083789. [Online]. Available: <http://doi.acm.org/10.1145/1083784.1083789>.
- [13] A. Bifet, G. Holmes, R. Kirkby, and B. Pfahringer, “Moa: Massive online analysis,” *J. Mach. Learn. Res.*, vol. 11, pp. 1601–1604, Aug. 2010, ISSN: 1532-

4435. [Online]. Available: <http://dl.acm.org/citation.cfm?id=1756006.1859903>.
- [14] G. Holmes, A. Donkin, and I. Witten, “Weka: A machine learning workbench,” in *Proc Second Australia and New Zealand Conference on Intelligent Information Systems*, Brisbane, Australia, 1994.
- [15] S. Akhter and J. Roberts, *Multi-core programming*. Intel press Hillsboro, 2006, vol. 33.
- [16] A. Prat-Pérez, D. Dominguez-Sal, J.-L. Larriba-Pey, and P. Trancoso, “Producer-consumer: The programming model for future many-core processors,” in *Proceedings of the 26th International Conference on Architecture of Computing Systems*, ser. ARCS’13, Prague, Czech Republic: Springer-Verlag, 2013, pp. 110–121, ISBN: 978-3-642-36423-5. DOI: 10.1007/978-3-642-36424-2_10. [Online]. Available: http://dx.doi.org/10.1007/978-3-642-36424-2_10.
- [17] N. Marz. (2015). About Me, [Online]. Available: <http://nathanmarz.com/about/> (visited on 03/20/2015).
- [18] (2015). Apache Zookeeper, [Online]. Available: <http://zookeeper.apache.org> (visited on 03/15/2015).
- [19] (2015). Apache Thrift, [Online]. Available: <https://thrift.apache.org> (visited on 03/15/2015).
- [20] (2015). Esoteric Software Kryo, [Online]. Available: <https://github.com/EsotericSoftware/kryo> (visited on 03/20/2015).
- [21] (2015). Netty, [Online]. Available: <http://netty.io> (visited on 03/15/2015).
- [22] (2015). LMAX Disruptor, [Online]. Available: <https://lmax-exchange.github.io/disruptor/> (visited on 03/20/2015).
- [23] (2015). Apache Kafka, [Online]. Available: <http://kafka.apache.org> (visited on 03/15/2015).
- [24] *Apache HBase*, <http://hbase.apache.org>, 2015. [Online]. Available: <http://hbase.apache.org> (visited on 03/27/2015).
- [25] F. Chang, J. Dean, S. Ghemawat, W. C. Hsieh, D. A. Wallach, M. Burrows, T. Chandra, A. Fikes, and R. E. Gruber, “Bigtable: A distributed storage system for structured data,” *ACM Trans. Comput. Syst.*, vol. 26, no. 2, 4:1–4:26, Jun. 2008, ISSN: 0734-2071. DOI: 10.1145/1365815.1365816. [Online]. Available: <http://doi.acm.org/10.1145/1365815.1365816>.
- [26] A. Chatzistergiou and S. D. Viglas, “Fast heuristics for near-optimal task allocation in data stream processing over clusters,” in *Proceedings of the 23rd ACM International Conference on Conference on Information and Knowledge Management*, ser. CIKM ’14, Shanghai, China: ACM, 2014, pp. 1579–1588, ISBN: 978-1-4503-2598-1. DOI: 10.1145/2661829.2661882. [Online]. Available: <http://doi.acm.org/10.1145/2661829.2661882>.
- [27] K. Chandrasekaran, S. Santurkar, and A. Arora, “Stormgen – A domain specific language to create ad-hoc storm topologies,” in *Proceedings of the 2014 Federated Conference on Computer Science and Information Systems, Warsaw, Poland, September 7-10, 2014.*, M. Ganzha, L. A. Maciaszek, and M. Paprzycki, Eds., 2014, pp. 1621–1628, ISBN: 978-83-60810-58-3. DOI: 10.15439/2014F278. [Online]. Available: <http://dx.doi.org/10.15439/2014F278>.

- [28] T. Dimson and M. Ganjoo, “Hailstorm: Distributed stream processing with exactly once semantics,” 2014.
- [29] C. Ranger, R. Raghuraman, A. Penmetsa, G. Bradski, and C. Kozyrakis, “Evaluating mapreduce for multi-core and multiprocessor systems,” in *High Performance Computer Architecture, 2007. HPCA 2007. IEEE 13th International Symposium on*, Ieee, 2007, pp. 13–24.
- [30] N. Marz. (Oct. 2014). History of apache storm and lessons learned – thoughts from the red planet – thoughts from the red planet, [Online]. Available: <http://nathanmarz.com/blog/history-of-apache-storm-and-lessons-learned.html> (visited on 03/15/2015).
- [31] (2015). Rabbit MQ, [Online]. Available: <http://www.rabbitmq.com> (visited on 03/15/2015).
- [32] (2015). The Streaming APIs, [Online]. Available: <https://dev.twitter.com/streaming/overview> (visited on 03/24/2015).
- [33] E. Gamma, R. Helm, R. Johnson, and J. Vlissides, *Design patterns: Elements of reusable object-oriented software*. Pearson Education, 1994, pp. 97–106.
- [34] L. Aniello, R. Baldoni, and L. Querzoni, *Adaptive online scheduling in storm*, <http://www.orgs.ttu.edu/debs2013/presentations/DEBS13-Paper88-Querzoni.pdf>, Jul. 2013. [Online]. Available: <http://www.orgs.ttu.edu/debs2013/presentations/DEBS13-Paper88-Querzoni.pdf>.
- [35] M. Thompson, D. Farley, M. Barker, P. Gee, and A. Stewart, “Disruptor: High performance alternative to bounded queues for exchanging data between concurrent threads,” May 2011. [Online]. Available: <http://lmax-exchange.github.io/disruptor/files/Disruptor-1.0.pdf>.
- [36] (2015). LMAX Disruptor Wiki, [Online]. Available: <https://github.com/LMAX-Exchange/disruptor/wiki/Performance-Results> (visited on 03/20/2015).
- [37] (2015). Java Thread Documentation, [Online]. Available: <http://docs.oracle.com/javase/7/docs/api/java/lang/Thread.html> (visited on 03/22/2015).
- [38] A. Danial, *Cloc*, <http://cloc.sourceforge.net>, 2014. [Online]. Available: <http://cloc.sourceforge.net> (visited on 03/26/2015).
- [39] D. Chandra, F. Guo, S. Kim, and Y. Solihin, “Predicting inter-thread cache contention on a chip multi-processor architecture,” in *Proceedings of the 11th International Symposium on High-Performance Computer Architecture*, ser. HPCA ’05, Washington, DC, USA: IEEE Computer Society, 2005, pp. 340–351, ISBN: 0-7695-2275-0. DOI: 10.1109/HPCA.2005.27. [Online]. Available: <http://dx.doi.org/10.1109/HPCA.2005.27>.
- [40] (2015). IBM – InfoSphere Streams, [Online]. Available: <http://www.ibm.com/software/products/en/infosphere-streams> (visited on 03/24/2015).
- [41] J. Warner, *Top manual page*, <http://linux.die.net/man/1/top>, 2015. [Online]. Available: <http://linux.die.net/man/1/top> (visited on 03/24/2015).