



R-Storm

Resource-Aware Scheduling in Storm

Gheorghică Mutu
SAI I
Event Based Systems



Abstract

Big data has led to the popularity of real-time distributed stream processing systems like Apache Storm.

Storm lacks intelligent scheduling, as the default round-robin approach is inefficient.

R-Storm is a system that implements resource-aware scheduling within Storm.

R-Storm aims to increase overall throughput by maximizing resource utilization and minimizing network latency.

R-Storm can satisfy soft and hard resource constraints and minimize network distance between components.



Introduction

Challenges in processing large amounts of data quickly in the age of data dominance.

In 2012, 2.5 exabytes of data were created every day, and this has increased to 2.3 zettabytes in 2014.

Storm is an open source distributed real-time computation system.

R-Storm is the first system to implement resource-aware scheduling in Storm.

R-Storm outperforms default Storm in overall throughput and resource utilization.

Efficient scheduling of multiple topologies.



Storm: Distributed Real-time Data Processing Framework

Storm is a distributed processing framework for real-time data processing.

Storm topologies are computation graphs that process live data in real-time.

Storm jobs fragment input data into chunks processed by tasks.

Unlike MapReduce, Storm topologies run forever until they are killed.

Key terms in Storm: Tuples (basic unit of data), Streams (unbounded sequence of tuples), Components (processing operators), Tasks (instantiation of Spout or Bolt), Executors (threads that execute tasks), Worker Process (process spawned by Storm that runs executors).



Storm Components: Spouts and Bolts

Components in Storm are called Spouts and Bolts.

Spout: Source of data streams, emits unbounded number of tuples downstream.

Example: Spout can receive live financial data and transform it into tuples for processing.

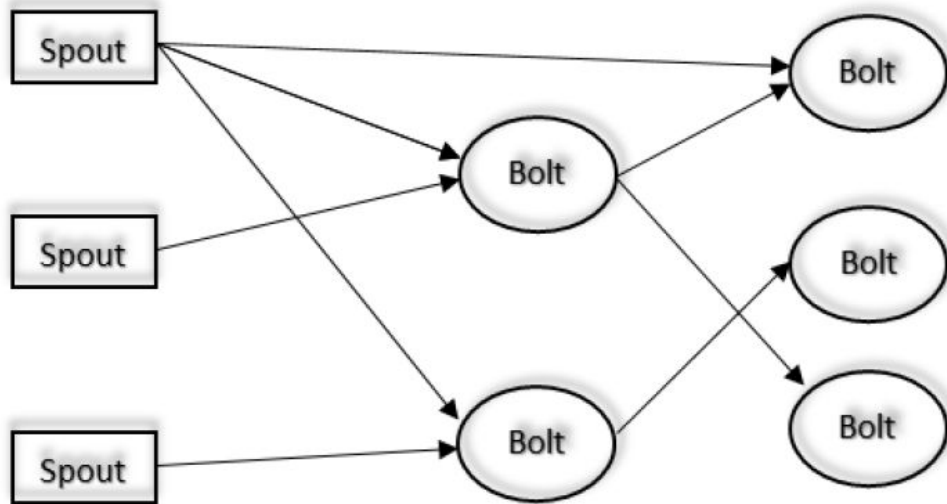
Bolt: Consumes, processes, and potentially emits new data streams.

Bolts can consume input streams from Spouts or other Bolts, perform processing, and emit new streams.

Bolts can filter, aggregate, join, query databases, and perform user-defined functions.

Multiple Bolts can work together to compute complex stream transformations, such as trending topics in tweets from Twitter.

An Example of Storm topology





Storm Cluster: Master & Worker Nodes

Master Node

Schedules tasks among worker nodes and maintains an active membership list.

Nimbus daemon communicates with Zookeeper for consistency and fault tolerance.

Worker Node

Executes tasks assigned by the Master Node.

Supervisor daemon listens for tasks and can have multiple worker processes for efficient processing.



Storm Topology and Parallelization

Links between components indicate how tuples are passed in a Storm topology.

Components can be parallelized to improve throughput.

User needs to specify parallelization hint for each component to set concurrent tasks.

Tasks from a component can be executed at different physical locations to optimize performance.

Default Storm scheduler places tasks of bolts and spouts on worker nodes in a round-robin manner for distributed processing.

Intercommunication of tasks & Storm machine

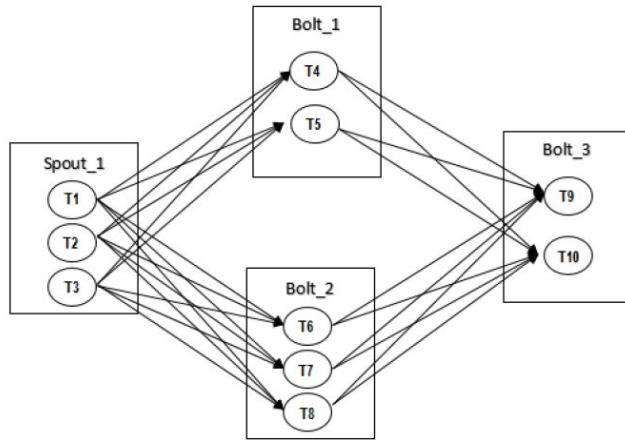


Figure 2: Intercommunication of tasks within a Storm topology

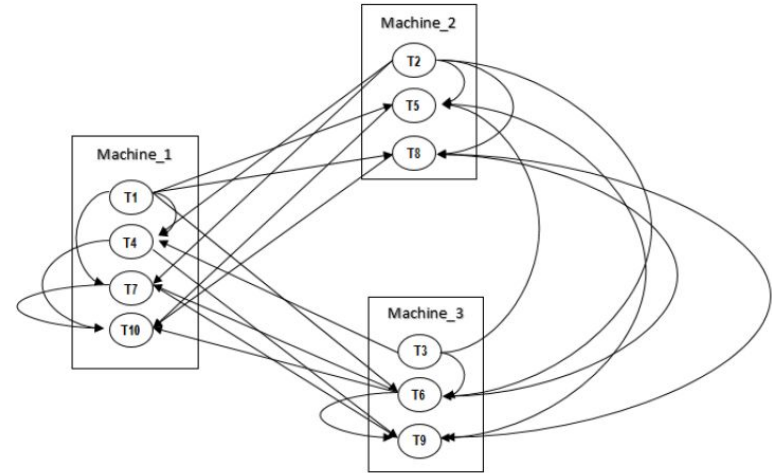


Figure 3: An example Storm machine



Resource-aware Task Assignment

Problem: Assigning tasks to machines while meeting resource requirements.

Three types of resources: CPU, memory, and bandwidth.

Hard constraints must be fully satisfied, soft constraints can be partially satisfied.

User-specified constraints.

Linear programming optimization problem.

Assumes limited budget for resources on each node, and task-specific resource requirements.

Considers resource utilization trade-offs for performance improvement.



Resource-aware Task Assignment

Tasks have CPU, bandwidth, and memory requirements.

Nodes have budget constraints for CPU, bandwidth, and memory.

Goal is to maximize resource utilization, minimize network latency, and stay within budget.

Formulated as a complex variation of Knapsack optimization problem.

Challenges: multiple knapsacks (clusters and nodes), multidimensional knapsack problem (3 resources), and quadratic knapsack problem (successive task assignment).

Efficient approximation algorithms can be utilized for computational feasibility.



Resource-aware Task Assignment

Problem: Quadratic Multiple 3-Dimensional Knapsack Problem (QM3DKP)

Previous research applied knapsack problem variations to specific contexts

Existing algorithms for knapsack problems have limitations in terms of computational complexity and real-time scheduling requirements

Need for a scheduling algorithm that can handle multiple nodes, resource requirements, and scheduling decisions quickly

Introducing R-Storm as a solution for scheduling tasks within Storm with simplicity and low overhead.

Resource-aware Task Assignment

$$\begin{aligned} & \text{Maximize} \quad \sum_{i \in \text{clusters}} \sum_{j \in \text{nodes}} Q_{\theta_{ij}} \\ & \text{s.t.} \quad \sum_{i \in \text{clusters}} \sum_{j \in \text{nodes}} c_{\tau_{ij}} \leq W_1, \\ & \text{and} \quad \sum_{i \in \text{clusters}} \sum_{j \in \text{nodes}} b_{\tau_{ij}} \leq W_2, \\ & \text{and} \quad \sum_{i \in \text{clusters}} \sum_{j \in \text{nodes}} m_{\tau_{ij}} \leq W_3. \end{aligned} \tag{1}$$

which is simplified into

$$\begin{aligned} & \text{Maximize}_{\{\mathcal{N}' \subseteq \mathcal{N}\}} \sum_{\theta \in \mathcal{N}'} Q_{\theta_i} \quad \text{subject to} \\ & \sum_{\tau_i \in \mathcal{N}'} c_{\tau_i} \leq W_1, \quad \sum_{\tau_i \in \mathcal{N}'} b_{\tau_i} \leq W_2, \quad \sum_{\tau_i \in \mathcal{N}'} m_{\tau_i} \leq W_3. \end{aligned} \tag{2}$$



R-Storm Scheduling Algorithm

Optimal solution to resource-aware scheduling in Storm is difficult and computationally infeasible.

Need simpler and effective algorithms for solving knapsack problems, specifically the Quadratic Multiple 3-Dimensional Knapsack Problem.

Observations about Storm's environment: data centers with server racks connected by switches.

Insights on communication latency: inter-rack communication is slowest, inter-node communication is slow, inter-process communication is faster, and intra-process communication is the fastest.

Scheduling algorithm designed based on these insights.

Two core parts, task selection and node selection.

R-Storm Scheduling Algorithm

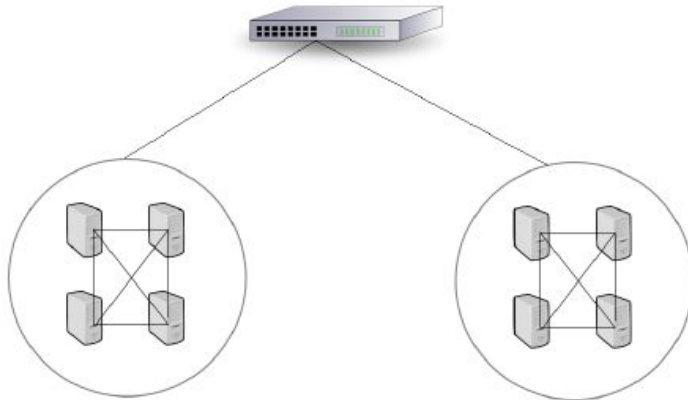


Figure 4: Typical cluster layout

Figure 5: An example node selection in a 3D resource space

Algorithm 1 R-Storm Schedule

```
1: procedure SCHEDULE
2:    $taskOrdering \leftarrow TASKSELECTION()$ 
3:   for each Task  $\tau$  in  $taskOrdering$  do
4:     Node  $n \leftarrow NODESELECTION(\tau, cluster)$ 
5:     SCHEDULE( $\tau, n$ );
6:   end for
7: end procedure
```



Algorithm Overview

Algorithm for scheduling tasks in a Storm topology on a cluster of nodes based on resource availability.

Nodes represented by vectors, where each vector represents resource availability.

Algorithm finds the node closest in Euclidean distance to the resource demand vector of a task, while ensuring hard constraints are not violated.

Tasks of components that communicate with each other have highest priority to be scheduled in close network proximity.

No hard resource constraints are violated.

Resource wastes on nodes are minimized.

Algorithm Overview

Figure 5 shows a visual example of a selected minimum-distance node to a given task in the 3D resource space, while the hard resource constraint (i.e. Z axis) is not violated.

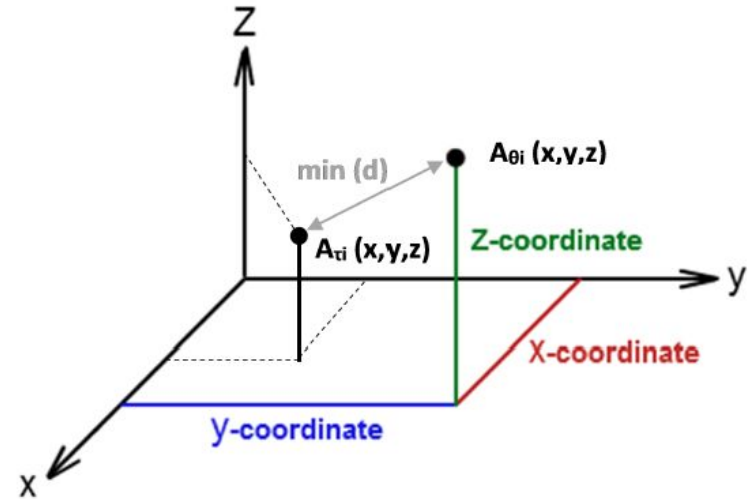


Figure 5: An example node selection in a 3D resource space



Algorithm Overview: Task Selection

Task selection is an important step in scheduling tasks in a Storm topology.

Traversal of the directed graph representing the Storm topology starts from spouts using breadth-first search (BFS) traversal.

BFS traversal creates a partial ordering of components, where adjacent components are placed in close succession.

Partial ordering of components is then used to create a partial ordering of tasks.

Ordering tasks based on the partial ordering of components ensures that tasks from adjacent components are scheduled close together, fulfilling the first desired property of the scheduling algorithm.



Algorithm Overview: Task Selection

Algorithm 2 Topology Traversal

```
1: procedure   BFSTOPOLOGYTRAVERSAL(Component
   root)
2:   queue; % queue of Components
3:   visted; % list of Components
4:   if root == null then
       return null
5:   end if
6:   queue.add(root)
7:   visted.add(root)
8:   while queue is not empty do
9:     Component com ← queue.remove()
10:    for each Component n in com.neighbor do
11:      if visted does not contain n then %check
whether visited or not
12:        queue.add(n);
13:        visted.add(n)
14:      end if
15:    end for
16:  end while
  return visted
17: end procedure
```

Algorithm 3 Task Selection

```
1: procedure TASKSELECTION
2:   components ← BFSTOPOLOGYTRAVERSAL(root)
3:   while taskOrdering does not contain all tasks do
4:     for each Component c in components do
5:       if c has tasks then
6:         Task  $\tau$  ← c.getTask()
7:         taskOrdering.add( $\tau$ )
8:         c.removeTask( $\tau$ )
9:       end if
10:    end for
11:  end while
  return taskOrdering
12: end procedure
```



Algorithm Overview: Node Selection

Node selection is the next step after task selection in scheduling tasks in a Storm topology.

For the first task in the topology, the server rack or sub-cluster with the most available resources is selected, and the first task is scheduled on the node in that server rack with the most available resources (Ref Node).

For the remaining tasks, nodes are selected based on the Distance procedure in Algorithm 4, where the bandwidth attribute is defined as the network distance from the Ref Node to each node.

Selecting nodes based on the Distance procedure ensures that tasks are scheduled tightly around the Ref Node, minimizing network latency for tasks communicating with each other.

Algorithm Overview: Node Selection

This process is visually depicted in Figure 5.

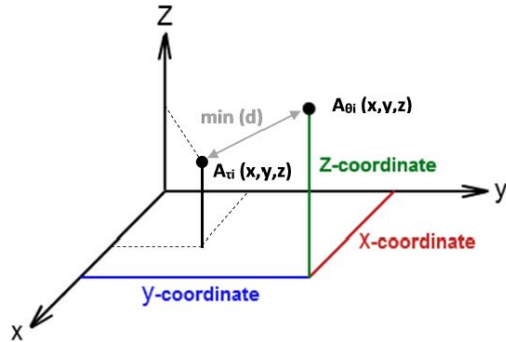


Figure 5: An example node selection in a 3D resource space

Algorithm 4 Node Selection

```

1: procedure NODESELECTION(Task  $\tau$ , Cluster  $cluster$ )
2:    $\mathcal{A}_\theta$ : set of all nodes  $\theta_i$  in the n-dimensional
      space   $\%(n=3 \text{ in our context})$ 
3:    $\mathcal{A}_\tau$ : set of all tasks  $\tau_i$  in the n-dimensional
      space   $\%(n=3 \text{ in our context})$ 
4:    $A_{\theta_i}$ : n-Dimensional vector of resource availability on
      each node  $\theta_i$ , such that  $A_{\theta_i} \in \mathcal{A}_\theta$ 
5:    $A_{\tau_i}$ : n-Dimensional vector of resource availability for
      each task  $\tau_i$ , such that  $A_{\tau_i} \in \mathcal{A}_\tau$ 
6:   if global  $refNode == \text{null}$  then
7:     ServerRack  $s \leftarrow$ 
      findServerRackWithMostResources( $cluster$ )
8:      $refNode \leftarrow$  findNodeWithMostResources( $s$ )
9:   end if
10:  select  $A_{\theta_j}$  such that:
       $A_{\theta_j} = \min d(A_{\tau_i}, A_{\theta_j}) \quad \forall A_{\theta_j} \in \mathcal{A}_\theta$  given
       $d(\lambda, \lambda^2) \leftarrow \text{DISTANCE}(\tau_i, \theta_j), H_{\theta_j} > H_{\tau_i} \%$ for all hard
      resource constraints
      return  $\theta_j$ 
11: end procedure
12: procedure DISTANCE(Task  $\tau_i$ , Node  $\theta_j$ )
13:    $distance \leftarrow weight_m * (m_{\tau_i} - m_{\theta_j})^2 + weight_c * (c_{\tau_i} -$ 
       $c_{\theta_j})^2 + weight_b * (\text{newtorkDistance}(refNode, \theta_j))$ 
      return  $\sqrt{distance}$ 
14: end procedure

```

Implementation

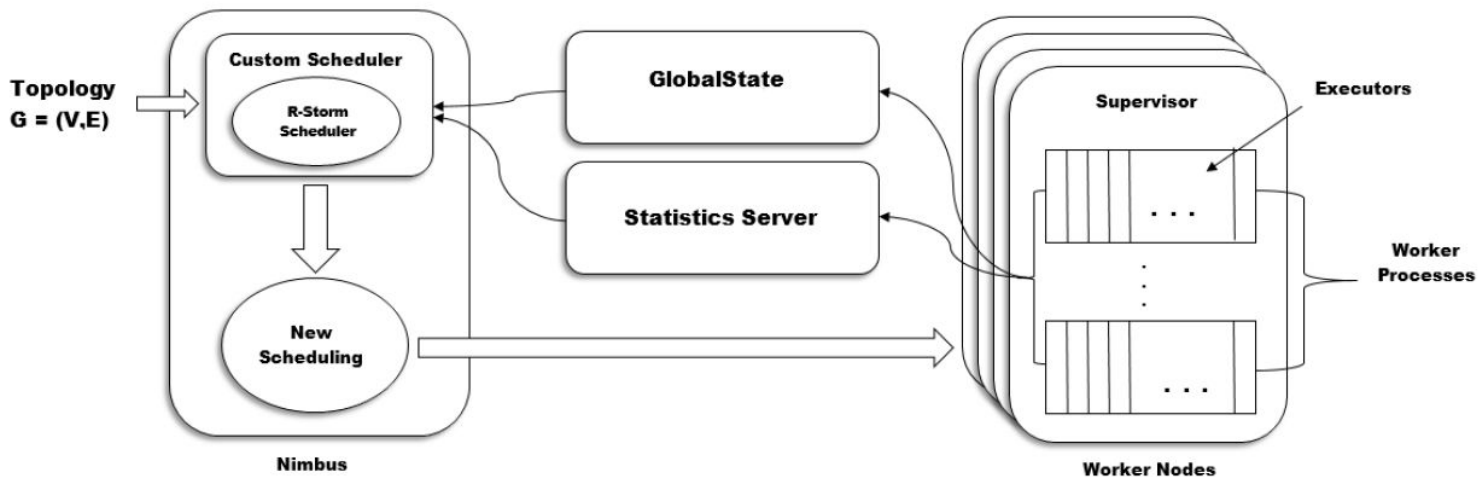


Figure 6: R-Storm Architecture Overview.



Implementation

R-Storm has three modules:

- **StatisticServer**: collects statistics in the Storm cluster for evaluative purposes
- **GlobalState**: stores important state information regarding scheduling and resource availability
- **ResourceAwareScheduler**: custom scheduler that implements the core R-Storm scheduling algorithm and initializes the other two modules.

There is also a list of APIs for the user to specify the resource demand of any component and the resource availability of any physical machine.



Evaluation: Experimental Setup

The aim is to evaluate the scheduling algorithm in a simulated real-world environment.

Emulab, a network testbed, was used to run experiments.

The Emulab setup consists of a total of 13 machines, with one designated as the master node and the other 12 as worker nodes.

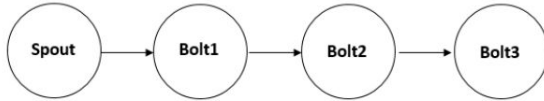
To emulate inter-rack communication latency, two VLANs were created, with each VLAN holding 6 machines.

The latency cost of inter-rack communication is 4ms for a round trip time.

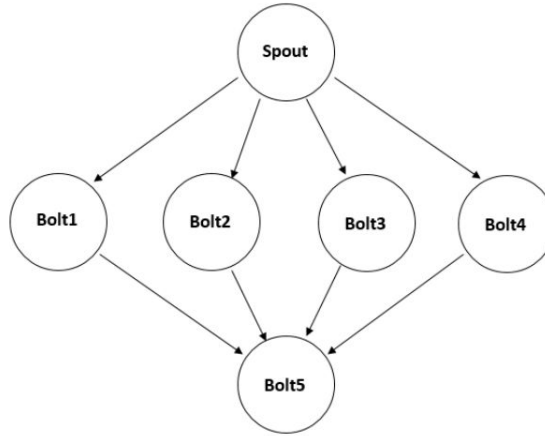
Each machine runs on Ubuntu 12.04 LTS with a single 3GHz processor, 2GB of RAM, 15K RPM 146GB SCSI disks, and is connected via 100Mbps network interface cards.



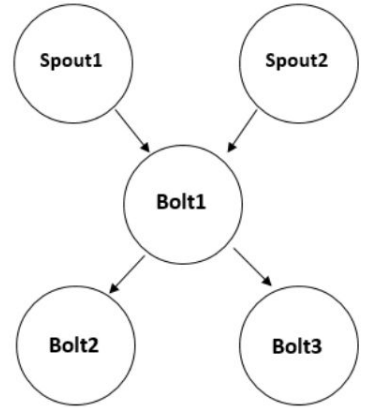
Evaluation: Micro-benchmark topologies



(a) Layout of Linear Topology



(b) Layout of Diamond Topology



(c) Layout of Star Topology

Figure 7: Layout of Micro-benchmark Topologies



Evaluation: Network Resource Bound

Micro-benchmark topologies configured to be network resource-bound

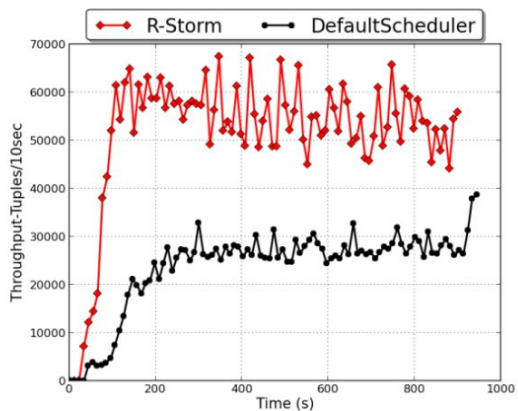
Overall throughput limited by network speed, not computation time

R-Storm outperforms Storm's default scheduler

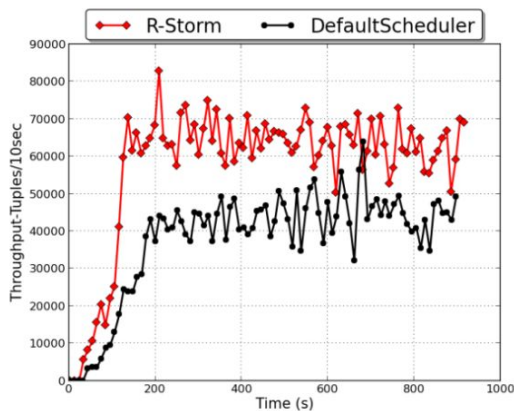
Results show around 50%, 30%, and 47% higher throughput for Linear, Diamond, and Star Topologies respectively

R-Storm minimizes network communication latency by colocating communicating tasks on the same machine or server rack

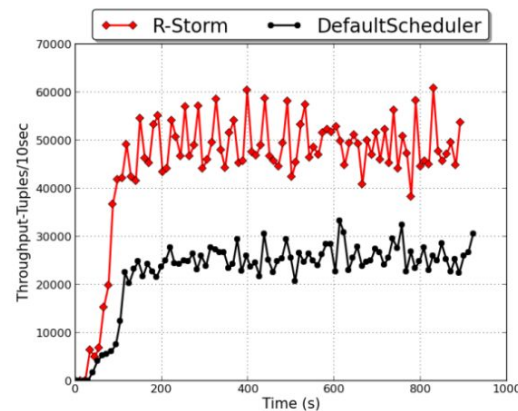
Evaluation: Network Resource Bound



(a) Linear Topology



(b) Diamond Topology



(c) Star Topology

Figure 8: Experimental results of Network-bound Micro-benchmark Topologies



Evaluation: Computation Time Bound

The overall throughput of micro-benchmark topologies is limited by computation time spent at each component.

Each component is configured to conduct a significant amount of arbitrary processing.

R-Storm efficiently schedules tasks closely together to maximize CPU utilization and minimize resource waste.

The cluster used consists of 12 machines partitioned into two racks.

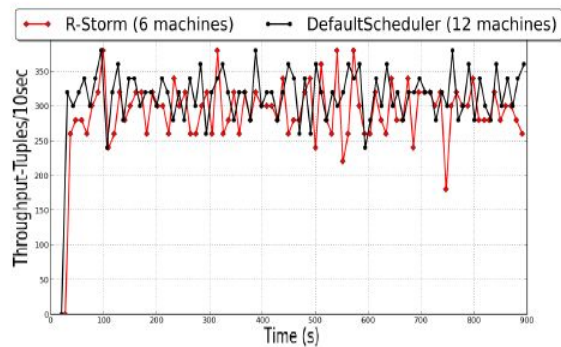
Storm's default scheduler schedules executors on all 12 machines regardless of the actual computation each executor will use.

By using R-Storm, the user can provide hints on how much CPU computation each instance of a component in a topology needs.

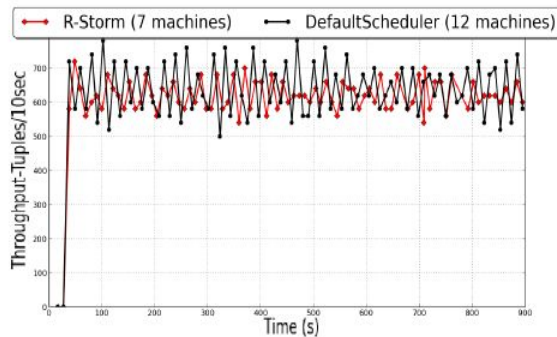
R-Storm only uses as many machines as needed to satisfy user-specified resource requirements.

R-Storm achieves the same level of throughput while using only a portion of the total number of machines.

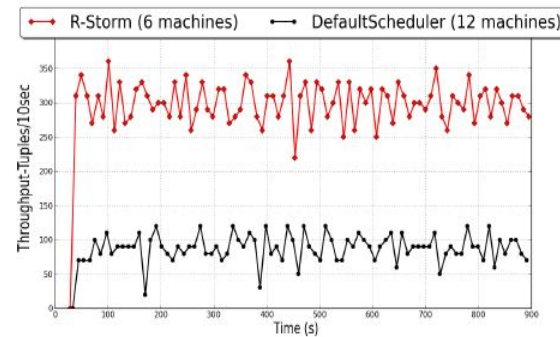
Evaluation: Computation Time Bound



(a) Linear Topology



(b) Diamond Topology



(c) Star Topology

Figure 9: Experimental results of Computation-time-bound Micro-benchmark topologies

Evaluation: Computation Time Bound

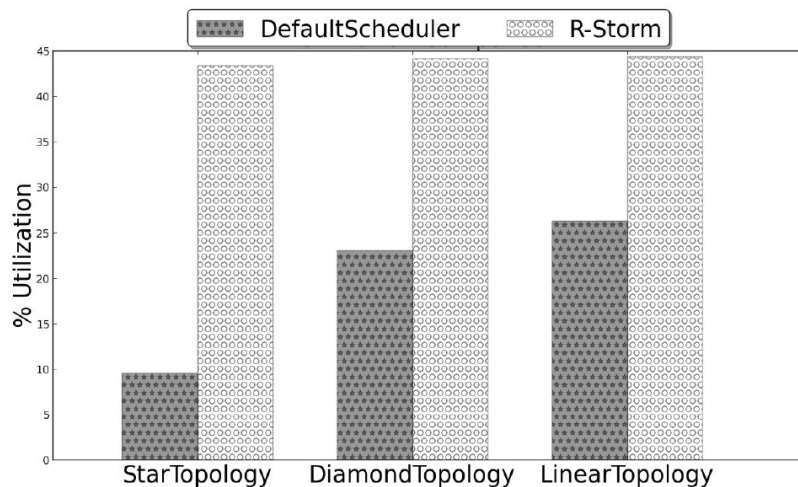


Figure 10 displays a comparison of the average CPU utilization of machines used in the cluster when scheduling using Storm's default scheduler versus R-Storm.

Figure 10: CPU Utilization Comparison



Evaluation: Yahoo Topologies

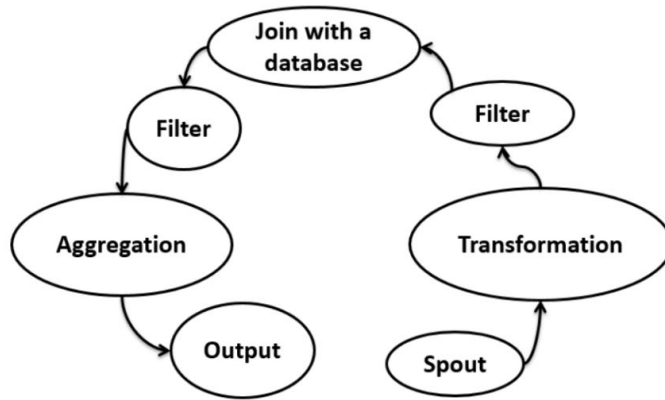
Two topologies from Yahoo! Inc. were evaluated: Page Load and Processing.

These topologies process event-level data from advertising platforms for real-time analytical reporting.

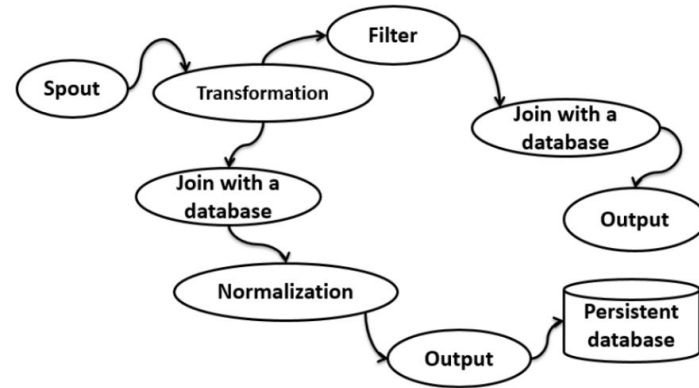
Experimental results show that R-Storm outperformed Storm's default scheduler.

On average, the Page Load and Processing Topologies have 50% and 47% better overall throughput, respectively, when scheduled by R-Storm as compared to Storm's default scheduler.

Evaluation: R-Storm on Industry Topologies



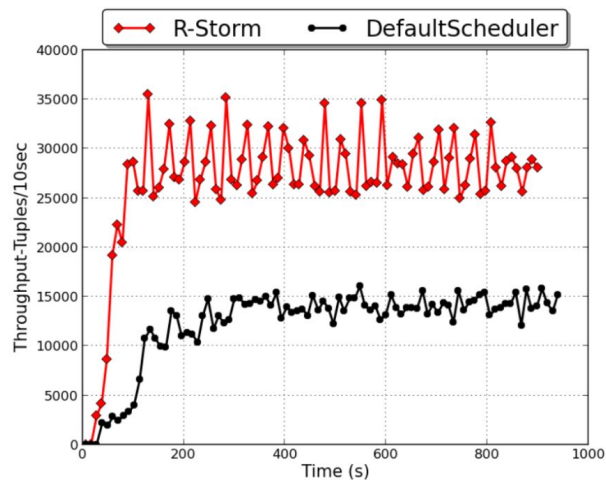
(a) Layout of Page Load Topology



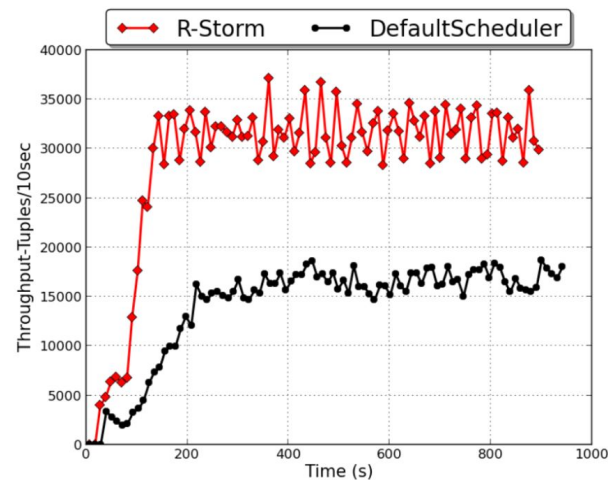
(b) Layout of Processing Topology

Figure 11: Production Topologies Modeled After Typical Industry Topologies

Evaluation: R-Storm on Industry Topologies



(a) Experiment results of Page Load Topology



(b) Experiment results of Processing Topology

Figure 12: Experiment results of Industry Topologies



Evaluation: Multi-topology Performance

The experiment evaluates the performance of R-Storm when scheduling multiple topologies in a cluster.

A 24-machine cluster is used, separated into two 12-machine subclusters.

Both the Yahoo! PageLoad and Processing topologies are scheduled by R-Storm and Default Storm.

R-Storm outperforms default Storm in terms of throughput for both topologies.

The Processing topology experiences terrible performance when scheduled by default Storm, with an average overall throughput near zero.

The average throughput of the PageLoad topology is around 53% higher when scheduled by R-Storm compared to default Storm.

The average throughput of the Processing topology is orders of magnitude higher when scheduled by R-Storm compared to default Storm.

Evaluation: Multi-topology Performance

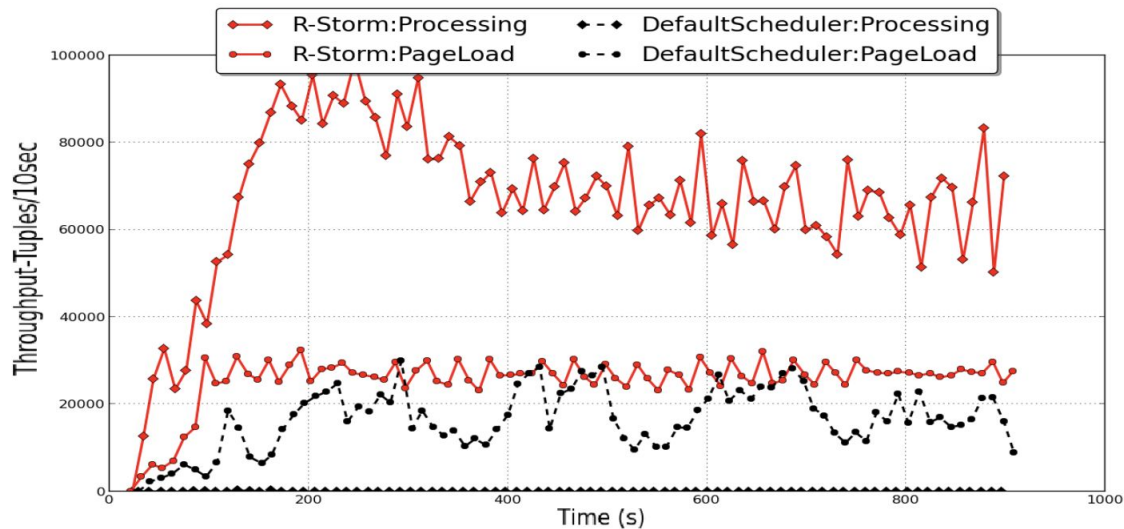


Figure 13: Throughput comparison of running multiple topologies.



Conclusion

Storm is an open-source platform for processing streams of data in real-time with a scheduling mechanism that is inadequate, scheduling tasks in a round-robin fashion with disregard to resource demands and availability.

R-Storm implements resource-aware scheduling within Storm, satisfying both soft and hard resource constraints and minimizing network distance between components that communicate with each other.

R-Storm achieves 30-47% higher throughput and 69-350% better CPU utilization than default Storm for micro-benchmark topologies and 50% higher overall throughput for Yahoo! Storm topologies.

R-Storm performs much better when scheduling multiple topologies than default Storm.

The concepts and algorithms used in R-Storm are applicable to other distributed data stream processing systems that have a DAG-based data processing model and should yield similar improvement as in Storm.