

Wang M, Jayaraman P, Solaiman E, Chen L, Li Z, Jun S, Georgakopoulos D, Ranjan R. [A Multi-layered Performance Analysis for Cloud-based Topic Detection and Tracking in Big Data Applications](#). *Future Generation Computer Systems* 2018. DOI: 10.1016/j.future.2018.01.047

Copyright:

© 2018. This manuscript version is made available under the [CC-BY-NC-ND 4.0 license](#)

DOI link to article:

<https://doi.org/10.1016/j.future.2018.01.047>

Date deposited:

18/02/2018

Embargo release date:

03 March 2019



This work is licensed under a [Creative Commons Attribution-NonCommercial-NoDerivatives 4.0 International licence](#)

A Multi-layered Performance Analysis for Cloud-based Topic Detection and Tracking in Big Data Applications

Meisong Wang¹, Prem Prakash Jayaraman¹, Ellis Solaiman¹, Lydia Y. Chen¹,
Zheng Li¹, Song Jun¹, Dimitrios Georgakopoulos¹, Rajiv Ranjan¹

^a*School of Computer Science, Australian National University, ACT, Australia*

^b*School of Computer Science, Newcastle University, Newcastle Upon Tyne, UK*

^c*Department of Electrical and Information Technology, Lund University, Sweden*

^d*Faculty of Science, Engineering and Technology, Swinburne University of Technology,
Melbourne, Australia*

^e*Zurich Research Laboratory, IBM, Zurich, Switzerland*

^f*Department of Computer Science, Chinese University of Geosciences, Wuhan*

Abstract

In the era of the Internet of Things and social media; communities, governments, and corporations are increasingly eager to exploit new technological innovations in order to track and keep up to date with important new events. Examples of such events include the news, health related incidents, and other major occurrences such as earthquakes and landslides. This area of research commonly referred to as Topic Detection and Tracking (TDT) is proving to be an important component of the current generation of Internet-based applications, where it is of critical importance to have early detection and timely response to important incidents such as those mentioned above. The advent of Big data though beneficial to TDT applications also brings about the enormous challenge of dealing with data variety, velocity and volume (3Vs). A promising solution is to employ Cloud Computing, which enables users to access powerful and scalable computational and storage resources in a "pay-as-you-go" fashion. However, the efficient use of Cloud resources to boost the performance of mission critical applications employing TDT is still an open topic that has not been fully and effectively investigated. An important prerequisite is to build a performance analysis capable of capturing and explaining specific factors (for example; CPU, Memory, I/O, Network, Cloud Platform Service, and Workload) that influence the performances of TDT applications in the cloud. Within this

paper, our main contribution, is that we present a multi-layered performance analysis for big data TDT applications deployed in a cloud environment. Our analysis captures factors that have an important effect on the performance of TDT applications. The novelty of our work is that it is a first kind of vertical analysis on infrastructure, platform and software layers. We identify key parameters and metrics in each cloud layer (including Infrastructure, Software, and Platform layers), and establish the dependencies between these metrics across the layers. We demonstrate the effectiveness of the proposed analysis via experimental evaluations using real-world datasets obtained from Twitter.

Keywords: Cloud-based TDT, Big Data, Performance Analysis, Cloud Computing

1. Introduction

The advent of Big Data applications that are fueled by numerous data sources such as social media and the Internet of Things, has created new opportunities for individuals, communities, governments, and corporations to make
5 use of this new and potentially important data that is continuously being generated. This area of research commonly referred to as Topic Detection and Tracking (TDT) is becoming a critical component of the current generation of Internet-based applications. An example where TDT research is of critical importance is in developing the capability to provide early detection and then
10 timely response to potential landslides using data obtained from sensors, and from social media outlets such as Twitter. The Achilles heel for TDT applications thus far has been limited access to real-time data which has an impeding effect on the accuracy of the application. The Big Data era has the potential to enhance the development of TDT applications by satisfying the requirement
15 of acquiring large volumes of data from variety of sources at high velocity. Traditional TDT techniques are incapable of coping with Big Data challenges best characterized by the 3V features, which are Variety, Velocity, and Volume. Volume means that the amount of data is so large that traditional storage devices

cannot store it (e.g. Every day, around 2.5 quintillion bytes of data is created,
20 which means that 90% of the data in the world was created in the last two
years [?]). Variety refers to the many sources and types of data, which creates
problems for storing, mining and analysing the data. Velocity means being able
to deal with the massive and continuous speed at which data flows from sources
like sensors, social media, and various networks to the cloud to be processed
25 and stored.

Recently, cloud computing techniques have emerged as reliable, effective and
practicable means for tackling the problems confronting TDT in the Big Data
era. For instance, there are a number of cloud storage frameworks both com-
mercial and free such as Amazon S3 that can be used to store large amounts
30 of data (*Volume*). Some NO-SQL databases can be used to store, process and
analyse various types of data (*Variety*). In addition, parallel computing frame-
works such as Apache Spark can be effectively used to significantly enhance the
speed of processing Big Data, and even to meet real-time analysis requirements,
which consequently addresses the "*Velocity*" problem. Another benefit that
35 Cloud computing can offer is the scalability that can satisfy the requirement of
processing data which is rapidly increasing in volume.

1.1. Motivation and Research Problem

Although Cloud computing creates clear advantages for TDT applications
(for processing and analysing Big Data) such as those identified above, it also
40 generates new challenges, and one of the most important challenges is how to
optimise the cloud resources to support mission critical TDT applications. An
important first step is to study and analyse the performance of cloud-based
TDT (CTDT) applications. Developing analysis capabilities that can capture
the performance of CTDT applications is not a trivial task given 1) the multi-
45 layered nature of cloud computing (IaaS, SaaS, and PaaS), 2) different metrics
required to capture the performance of TDT applications when compared with
other cloud-based applications such as e-commerce and customer relationship
management systems, and 3) dependencies between each of the metrics across

cloud layers. Existing TDT analysis techniques [?] [?] [?], capture the
50 performance of processing and analysing Big Data in clouds, but cannot be
applied accurately to model the performance of CTDT applications due to the
lack of consideration for all layers (end-to-end) that constitute a typical CTDT
application (i.e. IaaS, SaaS, PaaS, etc.).

1.2. Overview of Methods and Contributions

55 In this paper, we present a first kind of vertical multi-layered (infrastructure,
platform, and software) performance analysis which captures and analyses the
key metrics that have an important effect on the performance of CTDT big data
applications. The main contributions of this paper are:

- We clearly identify the key performance metrics that impact the perfor-
60 mance of CTDT applications with respect to each cloud layer (i.e. IaaS,
PaaS and SaaS).
- We then analyse and establish the dependencies between these metrics.
The aim of the analysis is to be able to capture the performance of CTDT
applications in order to be able to effectively optimise resources for such
65 applications deployed in clouds.
- We conduct comprehensive experimental evaluations using real-world datasets
obtained from Twitter to validate the effectiveness of the identified metrics
and their dependencies.

The paper is organized as follows: Section ?? summarizes a comprehen-
70 sive survey of existing work related to the optimization of CTDT applications,
and also existing work related to performance analysis for Cloud resource opti-
mization; Section ?? illustrates our performance analysis framework in detail;
To evaluate our performance analysis, we apply it to a specific case, which is
a Naïve Bayesian based CTDT application in Section ??; In Section ?? we
75 present experimental results based on a CTDT applications that we implement;
Conclusions and future directions are in Section ??.

2. Related Work

Studies that are related to our work can be divided into: 1) development and implementation of cloud-based TDT applications using machine learning techniques; and 2) performance analysis for platform-as-a-service TDT applications running on clouds using frameworks such as MapReduce. As we shall see, none of these studies can be used to efficiently analyse the end-to-end performance of CTDT applications. The first class of studies mainly focuses on how to develop and implement a CTDT application using various machine learning algorithms (e.g. state vector machine, Naive Bayes etc.). However, these works lack an analysis of factors that influence the performance of the CTDT application. On the other hand, the second set of studies are heavily focused on PaaS-based approaches such as Map Reduce and lack consideration for metrics such as performance of the distributed machine learning algorithms and related dependencies across the cloud layers layers (IaaS, PaaS and SaaS). These factors are important, and when not considered often lead to inaccuracy of performance modelling results. This will have significant consequences on mission critical CTDT applications that are dependent on fast, scalable and accurate analysis of events. For example, consider a landslide scenario. Under normal conditions, the sensors deployed in the field monitoring the activity of the land (e.g. movement of earth) produce data at a constant rate, and data coming from social media streams is relatively less constant. However, in case of an abnormal situation, the sensor data rate and social media data increases significantly resulting in increased volume. The challenge here is that a CTDT application running in the cloud needs to be able to optimise the cloud resources to cater for such diverse situations (normal and abnormal). Failing to do so will result in mission critical applications failing to meet their goals; i.e. detecting and alerting their users to important events [?] [?]. In cloud computing terminology, this is generally referred to as Quality of Service (QoS) guarantees enforced by service level agreements (SLA) [?].

Table ?? presents a summary of CTDT applications focusing on develop-

Table 1: Characteristics of cloud-based TDT applications.

Study	Performance Model	Performance Guarantee	Performance Metrics	IaaS	PaaS	SaaS
[?]	No	No	No	Yes	No	No
[?]	No	No	No	Yes	Yes	No
[?]	No	No	No	Yes	No	No
[?]	No	No	No	Yes	Yes	No
[?]	No	No	No	Yes	No	No

ment and implementation. As described earlier, the first class of CTDT applications lack performance analysis and evaluation capabilities, and provide no performance guarantees (QoS or SLA). This means that they cannot be used to
110 develop QoS guarantees for mission critical CTDT big data applications.

A summary of platform-as-a-service CTDT applications is shown in Table ???. As stated earlier, the focus of this related work is to develop a performance model for map-reduce or similar distributed framework-based TDT applications. We compare these approaches by using the taxonomy presented below:

- 115 1) HDFS: Are factors of HDFS taken into consideration?
- 2) Memory: Whether this work considers effects of memory.
- 3) Task Scheduler: Whether this work consists of scheduling mechanisms of MapReduce tasks.
- 4) Real Environment: Whether this work is based on a real environment or
120 another approach such as simulator.
- 5) Simulator: Whether this work is based on a simulator.
- 6) Greedy Algorithm: Whether this work uses greedy algorithms to calculate or estimate the execution time of MapReduce tasks. This is a separate research problem as different Map/Reduce scheduling strategies will lead to vary-
125 ing run-time performance (e.g., Mapper/Reducer response time). However, analysing how different scheduling strategies affect run-time performance is not the focus of this paper. In our model this is an input parameter available through workload benchmarking.
- 7) Network: Whether this works considers the impact of the network.

Table 2: Characteristics of related performance models.

Study	HDFS	Memory	ML	Task Scheduler	Real Environment	Simulator	Greedy Algorithm	Network
[?]	Yes	Yes	No	No	Yes	No	No	Yes
[?]	Yes	No	No	Yes	No	No	No	No
[?]	No	No	No	Yes	No	Yes	Yes	No
[?]	No	No	No	Yes	No	No	Yes	No
[?]	Yes	No	No	No	Yes	No	No	No
[?]	Yes	No	No	Yes	Yes	Yes	No	No
[?]	No	No	No	Yes	No	Yes	Yes	No
[?]	Yes	Yes	No	Yes	Yes	No	No	Yes
[?]	No	No	Yes	No	No	Yes	No	No
[?]	No	No	Yes	No	No	No	No	No
[?]	Yes	No	No	No	Yes	No	No	Yes

130 In summary, both classes of CTDT applications surveyed, lack the ability to represent the key metrics that influence the performance of CTDT applications across cloud layers. In order to support QoS guarantees (which we believe will be an essential part of future CTDT applications), it is essential to understand the impact of the application’s components on each layer in order to optimise
135 and orchestrate cloud resources. To the best of our knowledge, we are the first to present a performance analysis that considers the performance metrics within all end-to-end layers of a typical CTDT application, as well as the dependencies between each of those metrics.

3. Multi-layered Performance Model for CTDT Big Data Applications

3.1. Background

Let us consider a disease detection CTDT system. Such a system could potentially use a combination of MapReduce, HDFS and Amazon or Spark Streaming, HDFS and Windows Azure or Storm, HDFS and Google Compute

145 Engine. The goal of such a CTDT application is to provide timely and accurate notification to its users allowing them to respond to adverse events such as earthquakes or diseases outbreak. Current CTDT approaches depend on QoS guarantees provided by the cloud provider, which are limited and restrictive. For instance, it limits QoS to IaaS resources such as CPU, Memory and Storage
 150 [?]. However, to support CTDT applications such as the ones described earlier, there is a need to go beyond a simple QoS guarantee strategy to a more end-to-end approach, i.e., the QoS must satisfy constraints such as *events detected within x minutes of occurrence and notification delivered with y minutes*. We need to acknowledge that factors exerting substantial effects on the performance
 155 of a CTDT application come from different layers (SaaS, PaaS, and IaaS). For example, consider a typical Batch Processing architecture (e.g., MapReduce) presented in Figure ???. From the figure, we can see that several factors from different layers can affect the performance of a system. In the *machine learning libraries* layer, the accuracy and precision of the classification techniques
 160 such as the *Support Vector Machine (SVM)* and the *Naive Bayesian* model depends on the underlying input data sets (e.g., Tweets). However, in this work we validate the performance analysis technique in context of *Naive Bayesian* classification algorithm. Moreover in a MapReduce-based TDT application, the optimal number of Map Tasks is also essential for achieving the highest speed
 165 of a system. In addition, an appropriate scheduling method equally has a pivotal role to play in the speed of a system. For a CTDT application using a master-slave distributed file system (e.g., HDFS), single failure is obviously a catastrophe in terms of speed. In IaaS layer, for example, whether the applied memory is sufficient has a significant influence on the speed of a Spark-based
 170 TDT application.

In summary, we cannot ignore factors from any layer. Also in addition to considering factors from all layers we need to identify dependencies between these factors and how they can influence the performance of big data applications. Because, commonly, the cooperative effect of more than one metric has
 175 more effect or at least has equal effect on performance. Finally, the developed

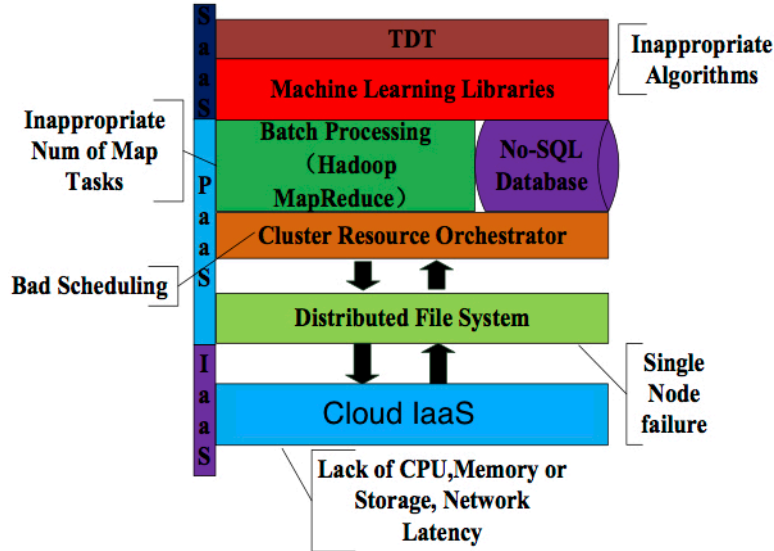


Figure 1: Factors which affect the performance in different layers.

analysis needs to cater to a range of CTDT big data applications rather than being constrained to a specific class.

3.2. Metrics influencing the Performance of CTDT applications

To capture the performance of CTDT applications, the first step would be to identify and determine which metrics should be used to measure the performance of a CTDT application at each layer. There are different performance metrics in terms of different practical needs. Be that as it may, there can be certain common important metrics that can be applied to most TDT applications such as speed, accuracy, price (for commercial applications), etc. Regardless of economic terms, speed is the factor of first-rate importance in a CTDT application particularly for mission critical disaster detection systems such as epidemic detection, earthquake detection, fire detection, etc. Consider earthquake detection for instance. Detecting the earthquake and warning citizens even a fraction of a minute earlier may save many lives. Furthermore, accuracy is another important metric for CTDT applications. A speedy but inaccurate traffic congestion detection system aiming to inform travellers about traffic jams or even

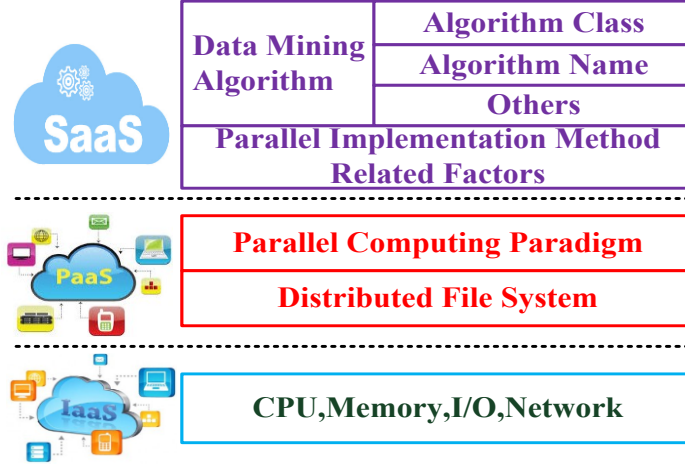


Figure 2: Architecture of a Performance Analysis Framework for CTDT Applications.

suggest alternative routes, for example, would mean nothing because it provides outdated or fraudulent information that misleads travellers, and could even lead to more traffic jams. We select speed and accuracy as two key metrics in our performance analysis. Speed can be measured by calculating the execution time of a CTDT application while accuracy differs in different kinds of CTDT applications in terms of different data mining algorithms adopted. Within this paper, and for the purposes of our experiments, “precision” is used to describe the accuracy of classification algorithms whereas “perplexity” is used to measure the accuracy of clustering algorithms.

3.3. CTDT Big Data Applications: Performance Analysis Framework

Figure ?? illustrates our proposed performance analysis framework. We develop a generic framework that could be easily adopted to model a range of CTDT big data applications that could include several technologies at each of the IAAS, PAAS and SAAS layer.

Data Mining Algorithm means the group of factors related to the data mining algorithm adopted. Different kinds of data mining algorithms [?] have different effects on both accuracy and speed. For instance, as we discussed before, measuring the accuracy of a clustering algorithm based CTDT application

210 requires the calculation of perplexity. In contrast, for a classification based one, we need to compute the precision. *Algorithm Class* means the type of data mining algorithm (e.g., Classification or Clustering) while *Algorithm Name* means the exact algorithm used (e.g., K-means, LDA, Naive Bayesian, etc.). Even in the same class, different algorithms might influence the performance of a system in different ways. For example, K-means and Canopy are both clustering
 215 algorithms, yet their influences on the speed of the system are substantially different, as K-means can be executed in more than one iteration whilst Canopy has only one iteration. *Others* refers to factors that might be important but beyond the scope of our existing work (providing scope for improvement). *Parallel Implementation Method* represents factors related to different paralleling
 220 methods of conversion of sequential data mining algorithms into parallel ones, such as MapReduce or MPI. *Parallel Computing Paradigm* means the different kinds of parallel computing frameworks adopted and relevant factors such as MapReduce (e.g., the factor of *Number of Mappers or Reducers*), Storm, Spark,
 225 etc. *Distributed File System* refers to factors related to the distributed system such as Hadoop Distributed File System. In the *IaaS* layer, we consider CPU, memory, I/O and Network related factors.

From the above architecture, it is obvious that our performance analysis defines several groups of factors rather than specific factors. Because different
 230 CTDT applications might adopt different implementation methods, such as different parallel computing paradigms (MapReduce or Storm). Our performance analysis can now be applied to almost all MapReduce-based TDT applications. In the next step, we will illustrate how to use it for a MapReduce-based Flu Detection system.

235 4. Using the Multi-layered Performance Analysis Framework to understand MapReduce-based TDT applications

In this section, we demonstrate how the proposed multi-layered performance analysis framework could study the impact of key identified parameters for

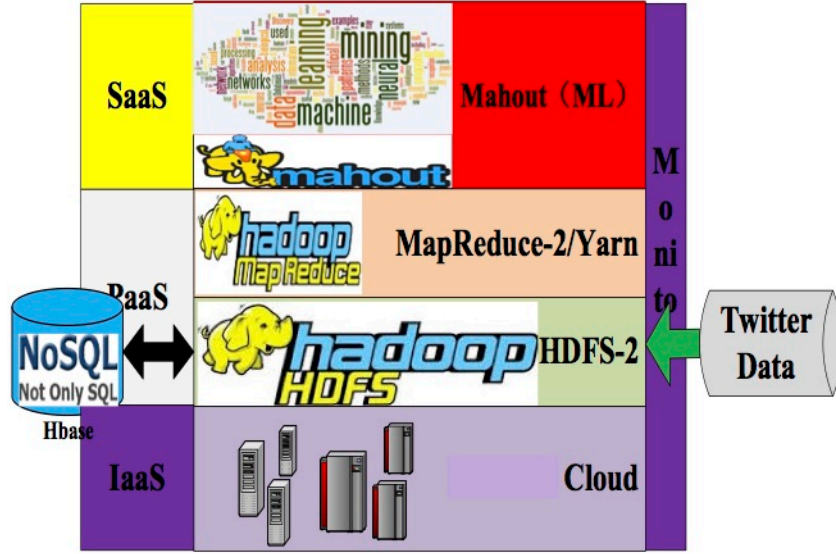


Figure 3: Architecture of Distributed Disease Detection System.

MapReduce-based TDT application.

4.1. MapReduce based TDT Application Architecture

We present the architecture of a MapReduce [?] based TDT application in Figure ?. The disease detection TDT application in this scenario uses data from Twitter to detect Flu-related events by analysing the tweets. First, we store the Twitter data in HDFS (Hadoop Distribute File System). In our work, the data was provided by COSMOS project (<https://www.cs.cf.ac.uk/cosmos/>).

We run MapReduce-2 and the HBase Database. On top of Hadoop, we employ Mahout [?] which is a distributed and scalable machine learning library. One of the advantages of Mahout is that most of its ML algorithms can be executed as a Map-Reduce job. The disease detection application (also known as an “epidemic detection” application) [?] is built on a combination of clustering, classification and topic detection algorithms.

Table 3: Factors in the IaaS Layer.

Name	Explanation
T	The capacity of a single node (the number of floating point operations FLOPs per second).
B	The bandwidth (Mbps).
P	The number of computers.
$Data$	Including data size and information in data.

4.2. Modelling of the Disease Detection System

As discussed in Section 3, speed is an important performance metric, therefore we will discuss how to model the speed of a MapReduce TDT application (Diseases Detection System).

The execution time of a MapReduce TDT process is actually a MapReduce data mining process consisting of one or more MapReduce jobs. In the MapReduce paradigm most jobs are executed in a sequential way, therefore calculating the execution time of a MapReduce data mining process can be divided into two parts: calculating the execution time of a single MapReduce job and calculating how many MapReduce jobs contained in a MapReduce data mining process. To calculate the execution time of a single MapReduce job, we need to identify the process of a single MapReduce job.

The calculation of a single MapReduce job involves capturing the performance in IaaS, PaaS and SaaS layers. Factors relevant to the IaaS Layer can be seen in Table ???. We explain the details of PaaS and SaaS using the example of using Naive Bayes' classification for predicting disease types.

4.2.1. IaaS Layer Analysis Factors

As discussed earlier, the analysis can be divided into three independent layers, which have dependencies on each other. We will illustrate the practical use of the analysis based on the diseases detection application in terms of the three layers. Factors relevant to the IaaS Layer can be shown in Table ??:

4.2.2. PaaS Layer Analysis Factors

For the PaaS Layer, by adopting Hadoop MapReduce and HDFS, the factors of the performance analysis are listed as shown in Table ??.

Table 4: Factors in the PaaS Layer.

Name	Explanation
T_{total}	The execution time (seconds) of a single MapReduce job.
T_{map}	The execution time (seconds) of a mapper task.
$T_{shuffle}$	The execution time (seconds) of a shuffle task.
T_{reduce}	The execution time (seconds) of a reducer task.
P_{start}	The percentage of the finished mapper tasks when the “shuffle” starts.
W_{map}	The product of the amount of workload for a single mapper task and it is related to the set of the blocksize of the HDFS, the spilt of the MapReduce, the data size (IaaS).
N_{map}	The number of the mapper tasks.
W	The workload of the whole input data size (MB or GB).
C_{umap}	Coefficient describing the relationship between the node (TaskTracker) and mapper.
$C_{ureduce}$	Coefficient describing the relationship between the node (TaskTracker) and Reducer.
$T_{ureduce}$	The execution time (seconds) of a single Reduce task.
N_{reduce}	The number of reducer tasks.
W_{reduce}	The workload for a single reducer task.
T_{umap}	The execution time (seconds) of a single mapper task.
$W_{uoutmap}$	The workload of the single mapper task.
W_{outmap}	The workload of all the mapper task.
B_{HDFS}	Blocksize of HDFS.
$N_{replication}$	The number of replication of data in HDFS.
MaxMemory of Map & Reduce task	Maximum Memory allocated to mapper or reducer Task can use, it can affect the execution time of a single task.

In fact, a MapReduce-based Data mining algorithm consists of one or several MapReduce jobs. Now we can calculate the execution time of each MapReduce job. The total execution time of a MapReduce job can be computed according to Equation (??).

$$T_{total} = T_{map} \times StartPercent + T_{shuffle} + T_{reduce} \quad (1)$$

Execution time of a map task can be computed using Equations (??) and (??).

$$T_{map} = T_{umap} \times N_{map} / P \quad (2)$$

$$T_{umap} = W_{map} \times C_{umap} / T \quad (3)$$

C_{umap} depends on several factors, such as the CPU speed, memory size, and

available network bandwidth, etc. From the above formula, we can see that
 285 by increasing the number of nodes (i.e. number of Map and Reduce instances),
 end-to-end execution time of a MapReduce job (and the CTDT application)
 can be reduced. Unfortunately, it is not always the case, due to that C_{umap} will
 change with the changing of other parameters, such as CPU, Memory, Number
 of Mapper, etc.

290 In the MapReduce based Hadoop framework the N_{map} parameter (number
 of Map Tasks) is determined by setting: “dfs.block.size”, “mapred.map.tasks”,
 “mapred.min.split.size”, “input data size”, “goal number of mapper”, and “mapred.max.split.size”.
 How to compute N_{map} will be illustrated in the following equations.

The execution time of a reduce task can be computed by using Equations
 295 (??), and (??).

$$T_{reduce} = T_{ureduce} \times N_{reduce} / P \quad (4)$$

$$T_{ureduce} = W_{reduce} \times C_{ureduce} / T \quad (5)$$

The coefficient $C_{ureduce}$ is similar to C_{umap} , and the only difference is that
 $C_{ureduce}$ is for Reduce tasks (Reducer). The formula to compute the execution
 time of a shuffle task is shown in (??) and (??).

$$T_{shuffle} = W_{uoutmap} \times (N_{map} \bmod P) / B \quad (6)$$

$$W_{uoutmap} = W_{outmap} / N_{map} \quad (7)$$

300 The number of Map tasks is determined by the following parameters: size
 of block in HDFS “dfs.block.size”, the goal number “mapred.map.tasks”, the
 minimum size of splitting data for each mapper “mapred.min.split.size” and
 the maximum size of splitting data for each mapper “mapred.max.split.size”.

Table 5: Factors in the SaaS Layer.

Name	Explanation
Execution Time	The whole execution time of Bayes' classification includes time taken for training, testing and learning.
Precision	The accuracy of the classification.
Njob	Depends on the PaaS level the number of Mapper and Reducer parameter/factor.
Class of ML Algorithm	Classification.
Name of ML Algorithm	Naïve Bayes'.
Complement (Boolean value)	Training process is based on C Naïve Bayesian or Standard Naïve Bayesian.
RunSequential (Boolean value)	MapReduce way or sequential way.

4.2.3. SaaS Layer Analysis Factors

We consider the Naïve Bayes' classification ML algorithm [?] to aid our discussion of performance modelling at this layer. Except speed, we will also discuss the accuracy ("Precision" for Classification algorithms). See Table ?? which details features of this ML algorithm.

The total execution time of the whole classification process can be computed as shown below in Equation (??).

$$T_{bayes} = N_{job} \times T_{total} \quad (8)$$

"RunSequential" is a special parameter which determines if the Naïve Bayesian training process has to be executed in a MapReduce way. If this is set to "true", the training set will be executed in a sequential way. This can be a typical situation for a TDT application where the training data is not large enough to be processed in parallel by exploiting the MapReduce distributed parallel programming abstractions. While the training can be done sequentially on one cluster node, the actual classifying (testing) phase can be implemented in the MapReduce way. In other words, the performance analysis has to capture such complex configuration decisions if it has to guarantee end-to-end execution times.

As noted earlier, Precision depends on the underlying ML algorithms and the type of data sets under consideration. Hence, we need to undertake various experimental studies to verify which ML algorithm leads to best possible pre-

cision for a given dataset. Even for a given classification (ML) algorithm, the precision can change due to the changing of other parameters. For instance, the parameter “complement” determined if the MapReduce Naïve Bayesian classifier is trained by using “Complementary Naïve Bayesian”. This could lead to different precision as compared to standard Naïve Bayesian.

Some parameters in this layer might have influence on factors from other layers or require the assistance from factors of other layers to cooperate in order to affect the speed, or precision of the system. For instance, the “RunPartial” (MapReduce-based Random Forest) will determine if the MapReduce job will be executed in memory. If the MapReduce job is executed in memory, the job will be memory-intensive and more memory (IaaS resources) might lead to less execution time of Random Forest MapReduce-based jobs. One of the advantages of our performance analysis is that it can capture or reveal these inter-layer dependencies. Next, we will discuss such specific dependencies in relation to different algorithms such as Random Forest, Naïve Bayesian.

The basic theory of Naïve Bayesian classifiers is to group an unclassified item into a class where such an item has the highest probability related. For instance, $x = a_1, a_2 \dots a_m$ is an unclassified dataset and each a is a feature of x while $C = y_1, y_2 \dots y_n$ is a set of all classes and each y represents a class. Take the disease detection application for instance, x is an unprocessed tweet and y_i means a sort of known epidemic such as “flu”, “measles” or “Ebola”.

Traditionally a naïve Bayesian classifier process is sequential, which means it will not scale to processing of large volumes of Big Data. To efficiently process big data, it is better that naïve Bayesian classification algorithm should be parallelized. We adopt the MapReduce programming model to parallelize naïve Bayesian classification algorithm and explain the key steps and analyzing factors in the following.

There are two main steps in the training part of naïve Bayesian classifier: 1) Counting the ClassPrior $P(y_i)$ for each class. 2) Counting the conditional probability for each attribute per class $P(a|y_i)$ (in text classification, the attribute can be the word).

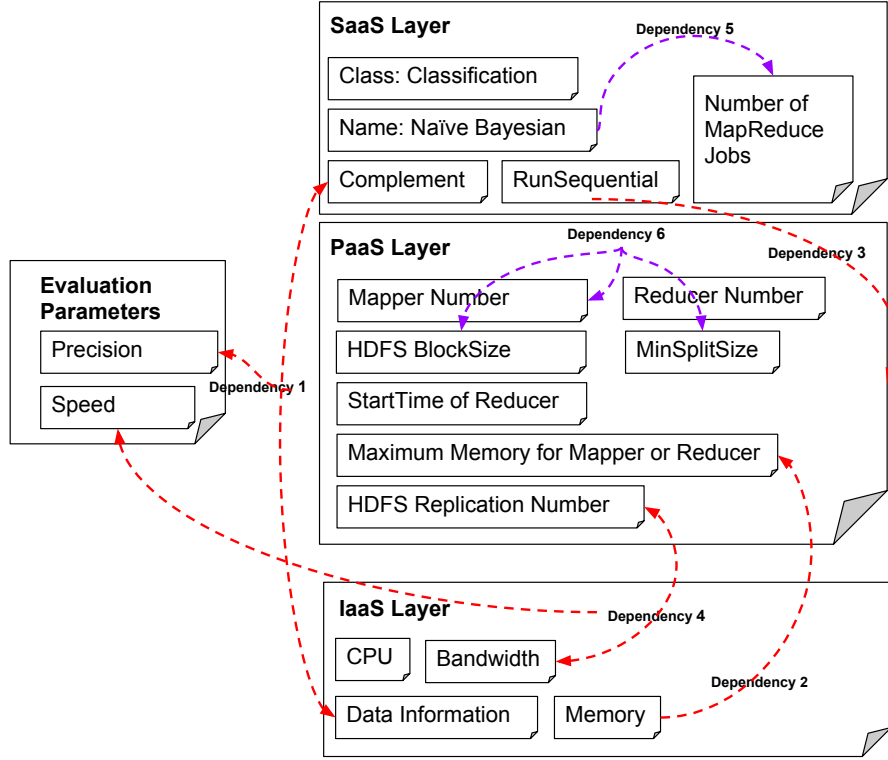


Figure 4: Dependency Across Layers

As a consequence, it is necessary for a naïve Bayesian classifier based on MapReduce to have two main MapReduce jobs to undertaking the above 2 steps: first for counting the Classifier and second for computing the conditional probability.

The practical implementation of naïve Bayesian classifier in MapReduce varies, especially for the training part. In this paper, we study Naïve Bayesian implementation based on the MapReduce framework by exploiting Mahout (an open-source scalable machine learning library) as an ML engine. Even in Mahout, the specific implementation of naïve Bayesian classifier has been changed since its initial release.

4.3. Dependency across layers

Here, we use the training process of Naïve Bayes in Mahout to illustrate the dependencies across layers as shown in figure ?? . Red lines represent dependencies from different layers while purple lines represent dependencies from the same layer.

- 1) Dependency1: Dependency between “Complement” and “Data Information”. Means that the influence of “Complement” on “Precision” might be affected by “Data Information”. For instance, the Complementary Naïve Bayesian method is more effective for classifying unbalanced data than the balanced data.
- 2) Dependency2: Dependency between “Memory” and “Maximum Memory for a Mapper or a Reducer”. It means that “Memory” in IaaS layer has to cooperate with “Maximum Memory for a Mapper or a Reducer” in order to manage speed of executing of analytics tasks. Specifically, without tuning the “Maximum Memory for a Mapper or a Reducer”, Memory available at the IaaS layer cannot affect the performance of the underlying TDT application. On the other hand “Maximum Memory for a Mapper or a Reducer” has an upper bound. For instance, if the total available Memory (IaaS) is 1000MB and the total number of Reducers and Mappers are 10, the “Maximum Memory” cannot be over 100MB, and otherwise the MapReduce job execution will not commence.
- 3) Dependency 3: Dependency between “RunSequential” and MapReduce. As we mentioned before, the MapReduce training process can be executed only when this parameter is false.
- 4) Dependency 4: Dependency between “Bandwidth” and “HDFS Replication Number”. When the replication of data is not enough, the node might need to copy required data from another node to process. In this situation, “Bandwidth” has a more significant role to play in the speed of the system, for the reason that low Bandwidth might lead to the slow speed of copying data from one node to another.

Table 6: Clusters adopted in experiments in CSIRO ICT Cloud.

Cluster	Specification
Cluster 1	1 node, pseudo-distributed Hadoop 2, HDFS 2, 1 CPU,
Cluster 2	2 nodes, 1 master node (Namenode, ResourceManager, JobTracker), 1 slave node (DataNode, NodeManager, TaskTracker), Hadoop2.4.1, Mahout 1.0, 2 CPU cores (2.40GHz)
Cluster 3	3 nodes, 1 master node, 2 slave nodes, Hadoop2.4.1, Mahout 1.0, 3 CPU cores (2.40GHz)
Cluster 4	4 nodes, 1 master node, 3 slave nodes, Hadoop2.4.1, Mahout 1.0, 4 CPU cores (2.40GHz)
Cluster 5	5 nodes, 1 master node, 4 slave nodes, Hadoop2.4.1, Mahout 1.0, 5 CPU cores (2.40GHz)
Cluster 6	6 nodes, 1 master node, 5 slave nodes, Hadoop2.4.1, Mahout 1.0, 6 CPU cores (2.40GHz)
Cluster 7	7 nodes, 1 master node, 6 slave nodes, Hadoop2.4.1, Mahout 1.0, 7 CPU cores (2.40GHz)
Cluster 8	8 nodes, 1 master node, 7 slave nodes, Hadoop2.4.1, Mahout 1.0, 8 CPU cores (2.40GHz)
Cluster 9	9 nodes, 1 master node, 8 slave nodes, Hadoop2.4.1, Mahout 1.0, 9 CPU cores (2.40GHz)
Cluster 10	10 nodes, 1 master node, 9 slave nodes, Hadoop2.4.1, Mahout 1.0, 10 CPU cores (2.40GHz)

- 5) Dependency5: It is an inner Dependency within the SaaS layer. The parameter “Number of MapReduce jobs” is determined by the parameter “Name of Algorithm”. In a Naïve Bayesian performance analysis it is 3 for training set (in the new edition of Mahout) and 1 for the testing (classifying) part, and for another classification algorithm (e.g. Random Forest), it might require a different number of tasks at training and testing steps/phases
- 6) Dependency6: An inner Dependency within PaaS layer, the number of mappers is affected by the size of the HDFS block and the min splitting size of input data.

5. Experimentation and Evaluation

5.1. Experimental Environment

The environment of our experiments is based on a CSIRO ICT Cloud which is built with OpenStack. There are 10 clusters adopted in our experiment, shown in Table ???. As explained previously, the data for our experiments is collected from Twitter in order to detect outbreaks of flu.

We did all the experiments that required maximum 4 nodes at first, then we created snapshots of Clusters 1-4 and we did experiments on Cluster 5, Cluster 6, Cluster 7, Cluster 10. The reason we run our experiment under these different settings is that we will present the effect generated by the IaaS resources upon the speed of the system.

5.2. Experimental Results

5.2.1. IaaS Experiment (Number of VCPU cores)

Description of Experiment: In accordance with our performance analysis, when the CPU resource is enough, increasing of the CPU resource does not affect the execution time of a TDT application significantly. However, when the CPU resource is in shortage, for example, there is only a virtual machine with 1 core CPU in a cluster and the MapReduce-based data mining algorithm in a TDT application requires more than 5 Map tasks, the increasing of CPU resource might lead to the increasing of the speed of such a TDT application.

Because our system is built on CSIRO Cloud where we do not have the highest level of access privilege, we can only change the number of VCPU (Virtual CPU) cores. As we mentioned in this chapter, we built cluster 1-10 (Shown in Table ??). To eliminate the effect of memory, we kept the memory size of each Mapper or Reducer unchanged and the number of Mappers or Reducers unchanged. In this experiment, we chose Naive Bayesian as our algorithm, Figure ?? shows the result of the experiment.

The first figure shows that the speed increased with the increasing of the CPU core number, but the second figure shows that the speed was not affected by the increasing of the CPU core number. The first figure shows a group of experiments based on the Mapper number of "18" while the second figure represents a group of experiments based on the Mapper number of "1". Specifically, in the first group from cluster 1 to 10, the CPU resource of each cluster might not have the maximum required CPU resource, which led to the situation that all the Mappers might not be able to start at the same time. Consequently, with the increasing of CPU resources, the number of Mappers which could be

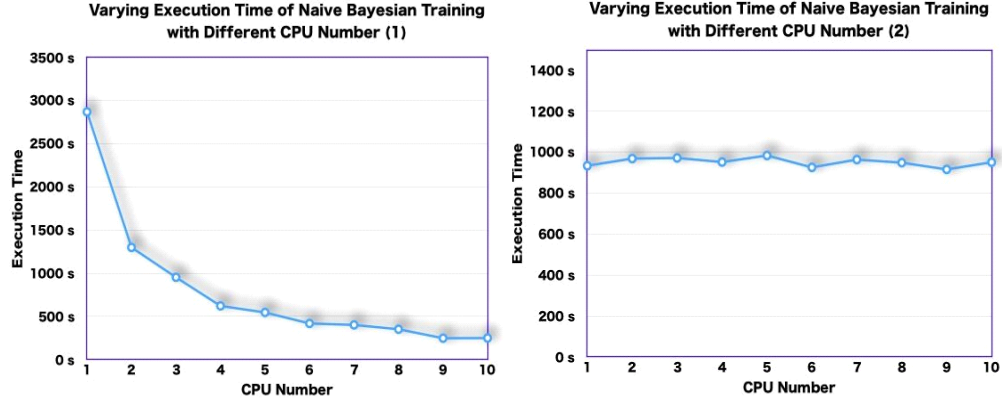


Figure 5: Result of execution time of Naïve Bayesian Trainings with Different Number of CPU Cores.

open in the meantime increased and this led to the increasing of the speed. The second group of experiments was based on 10% of the data of the first group (for saving time), 1 Mapper and 2 Reducer (the same as with the first group). The Mapper Number is "1" and we set the parameter "mapreduce.map.cpu.vcores" (number of virtual cores to request from the scheduler for each map task) as "1" and "mapreduce.reduce.cpu.vcores" (number of virtual cores to request from the scheduler for each reduce task) as "1".

Conclusion of Experiment: The number of CPU cores will affect the speed when the CPU resource is so little that it cannot start all the Mappers at the same time. When the CPU resource is sufficient, the increasing of CPU numbers cannot affect the speed significantly. The result also shows how the influence of different parameters: "Number of CPU" (IaaS), "Mapper Number" and "mapreduce.map.cpu.vcores" might affect the speed together. This has been identified in our performance analysis.

5.2.2. PaaS Experiment (Number of Mappers and Reducers)

Description of Experiment: As mentioned in our performance analysis there is an optimal number of Mappers for the speed of the system, and the number of Mappers might affect the speed significantly. We change the number of Mappers and make other factors fixed. The result can be shown in Figure

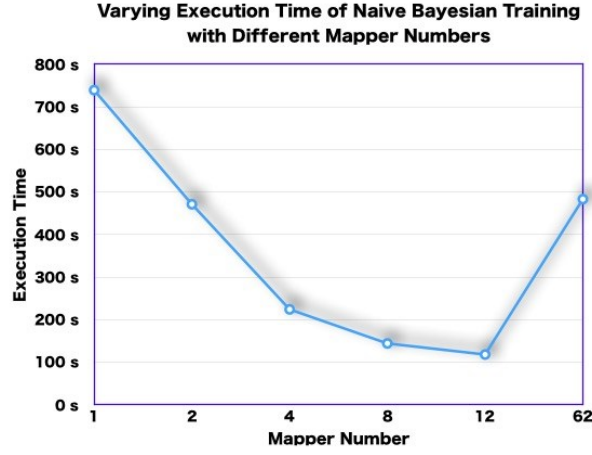


Figure 6: Execution Time of Naïve Bayesian Training with Different Mapper Numbers.

??.

Conclusion of Experiment: From the results, we can conclude that the number of Mappers can affect the speed of Naïve Bayesian training and Random Forest Training. Furthermore, there is an optimal number of Mappers for a MapReduce-based Naïve Bayesian and Random Forest (Training). There is an optimal value for mapper number that can achieve a minimum execution time.

5.2.3. SaaS Experiment

Description of Experiment: According to our performance analysis, the other parameter possessing a significant role to play in the performance of a TDT application based on Naïve Bayesian is “trainComplementary” which determines whether the Naïve Bayesian algorithm is executed as ”Complementary Naïve Bayesian” or “Standard Naïve Bayesian”. Complementary Naïve Bayesian is a Naïve Bayesian variant overcoming some weaknesses of the standard Naïve Bayesian. The Naïve Bayesian classifier tends to classify documents into a category possessing a great number of documents while the complementary uses data from all categories apart from the category that is worked on.

This parameter might affect the precision of the system, in accordance with our performance analysis. To evaluate the effect of this parameter, we conducted

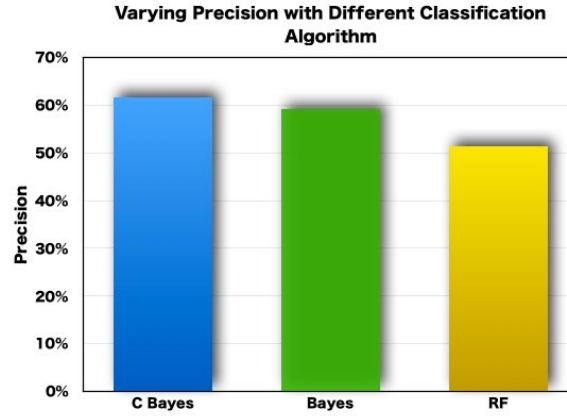


Figure 7: Varying Precision with Different Classification Algorithms.

the following experiment: keeping other parameters unchanged and seeing the result of precision in terms of different kinds of classification algorithms. In this experiment, we also compare the precision of Random Forest classifier with the same data. The result of this experiment is shown in Figure ??.

Conclusion of Experiment: The parameter “Complement” can control whether the classification algorithm is based on C Bayes or standard Bayes and indirectly affects the precision of the system. Furthermore, the parameter “name of classification algorithm” can affect the precision of the classification-based system, which means different classification algorithms have a different precision based on the same data.

5.3. Evaluation Summary

In conclusion, our experiments show that our performance analysis has achieved the following: 1) Our performance analysis is capable of capturing metrics which affect the performance of a CTDT application across all three layers. 2) Our performance model can illustrate the dependencies between these metrics.. 3) Our performance analysis can reflect on how these factors affect performance. 4) Our performance analysis can be used to predict the execution time of a TDT application under various conditions. As discussed in Section ?? to the best of our knowledge we are the first to present a performance analysis

that considers the performance metrics within all end-to-end layers of a typical
495 CTDT application, as well as the dependencies between each of those metrics.

6. Conclusions and Future Work

Cloud computing technology offers a possible solution to tackle new challenges of TDT (Topic Detection and Tracking) techniques in the Big Data era. However, this new combination of Cloud resources and TDT (CTDT) generates
500 a new issue – how to analyze the performance of CTDT to meet the demands posted by big data applications. Our performance analysis framework provides a practical and generic solution to analyse and model the performance of big data-based CTDT applications. We demonstrate the effectiveness of the performance analysis framework using the case study of MapReduce-based TDT
505 applications. Within our analysis, we have identified key parameters in each cloud layer, and established the dependencies between these metrics across the layers. We have also demonstrated and validated the correctness of parameters and their relationship across cloud layers via experimental evaluations using real-world datasets.

510 There are a number of issues that require further work. For example, we need to apply this performance analysis framework to more MapReduce-based TDT applications using different data mining algorithms, such as Random Forest, LDA, SVM, etc. Moreover, we will also extend this performance analysis framework to other classes of Big Data Applications, which can be based on
515 other types of programming paradigms such as Stream Processing. To achieve such generalization, we will extend the performance analysis framework to capture the limited sets of data flow and analytic patterns generated by different classes of Big Data Applications (e.g., real-time traffic modelling, and real-time energy modelling). For example, in context of the Stream Processing paradigm,
520 we will need to consider real-time analytic latency as the most important performance model parameter (at the PaaS layer), as compared to batch processing response time of the Hadoop programming paradigm. In our view, such exten-

sions will not affect formulations across the other layers of the Big Data stack including SaaS and IaaS.