

POLITECNICO DI MILANO
Facoltà di Ingegneria
Scuola di Ingegneria Industriale e dell'Informazione
Dipartimento di Elettronica, Informazione e Bioingegneria
Master of Science in
Computer Science and Engineering



Using Symbolic Execution to Improve the Runtime Management of Spark Applications

Supervisor:

PROF. LUCIANO BARESI

Co-Supervisor:

DR. GIOVANNI QUATTROCCHI

Master Graduation Thesis by:

DAVIDE BERTOLOTTI
Student Id n. 787366

Academic Year 2017-2018

Dedico questa tesi ai miei genitori
A mia moglie Egidia
e ai miei figli Fabio e Anna
che sono quanto di più meraviglioso e prezioso
abbia mai potuto desiderare

ACKNOWLEDGMENTS

Voglio ringraziare chi mi ha seguito nella preparazione e nella stesura di questa tesi, cioè il relatore Prof. Luciano Baresi e il correlatore Dr. Giovanni Quattrocchi, che si sono dimostrati disponibili, pazienti e prodighi di consigli.

Voglio anche ringraziare la mia famiglia: mia moglie Egidia, i miei figli Fabio e Anna ai quali ho sottratto attenzione e tempo che ho dedicato allo studio. Senza il loro sacrificio e supporto non sarei sicuramente riuscito a concludere questo mio impegno.

Davide

CONTENTS

Abstract	x
1 INTRODUCTION	1
1.1 Context	1
1.2 Problem and Motivation	4
1.3 Solution and Contribution	6
2 STATE OF THE ART	9
2.1 Big Data	9
2.1.1 Batch Processing: Hadoop	9
2.1.2 Batch Processing: Spark	13
2.1.3 Streams Processing: Flink	17
2.1.4 Streams Processing: Storm	19
2.2 Runtime Management of Big Data Applications	21
2.3 Elastic Resource Provisioning	22
2.4 Spark Resource Provisioning	25
2.4.1 Apache Hadoop Yarn	26
2.4.2 Apache Mesos	28
2.4.3 Spark on Yarn	31
2.4.4 Spark on Mesos	31
2.5 Virtualization and Containerization	34
2.6 Symbolic Execution	38
2.6.1 Symbolic Execution Engines	42
2.6.2 Path Selection	46
2.6.3 Symbolic Backward Execution	47
3 xSPARK	48
3.1 Architecture	51
3.2 Heuristic	53
3.3 Controller	54
4 PROBLEM STATEMENT AND SOLUTION OVERVIEW	57
4.1 Problem Statement	57
4.2 Solution Overview	59
4.2.1 xSpark _{SEEP}	61
5 IMPLEMENTATION	64
5.1 Overview	64
5.1.1 Background: Current xSpark Heuristic	64
5.1.2 Current xSpark Scheduling Limitation	65
5.2 Implementation Scope and Objective	65
5.2.1 Symbolic Execution	65
5.2.2 xSpark _{SEEP} vs. xSpark	66
5.2.3 A new Heuristic	72
5.3 Application Parameters	72
5.4 Application Profiling	73
5.5 PEP*	74

5.6	GuardEvaluator	75
5.7	Symbol Store	80
5.8	Heuristic	80
5.9	Symbols	86
5.10	Intercepting Actions Results	90
5.11	Python Tool	91
5.11.1	Core Functionality	91
5.11.2	Cloud Environment Configuration	91
5.11.3	Tool Configuration	92
5.11.4	Application Profiling: PEP and PEP* generation	96
5.11.5	Application Execution	96
5.11.6	Download & Requirements	98
5.11.7	xSpark-dagsymb commands	98
5.11.8	Example: Profile and Test CallsExample	101
6	EVALUATION	104
6.1	Test Environment	104
6.2	Tested Applications	104
6.3	Experiments	106
6.3.1	Results	106
6.4	Threats to Validity	112
7	CONCLUSION	114
7.1	Conclusion	114
7.2	Future Work	114
	BIBLIOGRAPHY	115

LIST OF FIGURES

Figure 2.1	map-reduce model	10
Figure 2.2	map-reduce model	11
Figure 2.3	hadoop map-reduce architecture	12
Figure 2.4	HDFS Node Strucure	13
Figure 2.5	Spark Standalone Architecture	15
Figure 2.6	Spark DAG Example	16
Figure 2.7	Apache Flink® - Stateful Computations over Data Streams.	17
Figure 2.8	Apache Flink® - Bounded & Unbounded Data Streams.	18
Figure 2.9	Apache Storm.	19
Figure 2.10	Apache Storm Components Abstraction.	20
Figure 2.11	Elasticity Matching Function derivation.	24
Figure 2.12	Resource Provisioning Chart.	24
Figure 2.13	Apache Hadoop YARN Architecture	27
Figure 2.14	Apache Mesos Architecture	29
Figure 2.15	Apache Mesos Resource Offer Example	30
Figure 2.16	Spark on YARN Client Mode	32
Figure 2.17	Spark on YARN Cluster Mode	32
Figure 2.18	Spark on Mesos Coarse Grained Mode	33
Figure 2.19	Spark on Mesos Fine Grained Mode	33
Figure 2.20	Full virtualization and paravirtualization	35
Figure 2.21	Architecture of virtual machines and containers	37
Figure 2.22	Example: which values of a and b make the assert fail?	39
Figure 2.23	Symbolic execution tree of function foo given in Figure 2.22. Each execution state, labeled with an upper case letter, shows the statement to be executed, the symbolic store σ , and the path constraints π . Leaves are evaluated against the condition in the assert statement. Image courtesy of Association for Computing Machinery [14].	40
Figure 2.24	Concrete and abstract execution machine models. Image courtesy of Association for Computing Machinery [14].	42

Figure 2.25	Concolic execution: (a) testing of function foo even when bar cannot be symbolically tracked by an engine, (b) example of false negative, and (c) example of a path divergence, where abs drops the sign of the integer at &x. Image courtesy of Association for Computing Machinery [14].	45
Figure 3.1	xSpark high level setup and execution flow	49
Figure 3.2	Example of profiling data from a Louvain application	49
Figure 3.3	Architecture of xSpark.	51
Figure 3.4	Set point generation for an executor controller.	55
Figure 3.5	xSpark chart aggregate-by-key.	55
Figure 4.1	Example Spark application with conditional branches and loops.	59
Figure 4.2	The four PEPs of the application of Figure 4.1.	60
Figure 4.3	Symbolic execution algorithm of SEEPEP.	61
Figure 4.4	$xSpark_{SEEP}$	62
Figure 5.1	Simplified solution components overview.	64
Figure 5.2	Structure of profile JumboJSON. Each profile is identified by its profile id's "0", "1", ..., "7" and, in general, is characterized by its own number of stages and key values. Left and right images show the same JumboJSON profile where profile "0" (left) and profile "1" (right) are expanded, to show their difference in terms of number of stages (6 stages for profile 0 and 3 for profile 1) and totalduration (76,200 ms for profile 0 ms and 43,000 ms for profile 1).	70
Figure 5.3	Information about jobs in json DAG profile.	71
Figure 5.4	Class diagram of Heuristic related classes.	82
Figure 5.5	Stage Cardinality Mismatch - Acceptable situation.	87
Figure 5.6	Stage Cardinality Mismatch - stage id extension.	87
Figure 5.7	Example of symbol formal names.	88
Figure 5.8	Credential template file.	93
Figure 5.9	Credential template file.	95
Figure 5.10	Example - Application Matlab Chart.	98
Figure 5.11	Example - Zoomed Details Application Matlab Chart.	99
Figure 5.12	Example - Worker Matlab Chart.	99
Figure 5.13	Example - Zoomed Details Worker Matlab Chart.	100
Figure 6.1	control.json $xSpark_{SEEP}$ configuration parameters.	105

LIST OF TABLES

Table 5.1	Example of Symbolic Memory Store contents.	69
Table 6.1	PromoCalls paths.	108
Table 6.2	Louvain paths.	108
Table 6.3	Results for PromoCalls	109
Table 6.4	Results for Louvain	110

LISTINGS

Listing 2.1	Spark word count application example.	17
Listing 5.1	Example of JSON profile.	66
Listing 5.2	Changes to SparkSubmit method "submit".	73
Listing 5.3	Example of Launcher Code .	73
Listing 5.4	Changes to class DAGScheduler.scala - reading PEPs.	75
Listing 5.5	Interface class IGuardEvaluator.	75
Listing 5.6	Class GuardEvaluatorPromoCallsFile, implementing the IGuardEvaluator interface.	75
Listing 5.7	Changes to class DAGScheduler.scala - Loading GuardEvaluator.	79
Listing 5.8	Changes to class DAGScheduler.scala - Initializing Symbol Store.	80
Listing 5.9	Changes to class ControlEventListener.scala - selecting the heuristic.	81
Listing 5.10	Class HeuristicSymExControlUnlimited.scala implementation.	81
Listing 5.11	Changes to companion object case class SparkStageWeightSubmitted of class SparkListener.scala	83
Listing 5.12	Changes to method submitMissingTasks of class DAGScheduler.scala	84
Listing 5.13	Changes to method onStageWeightSubmitted of class ControlEventListener.scala	85
Listing 5.14	Changes to class DAGScheduler.scala - method runJob.	88
Listing 5.15	Changes to class DAGScheduler.scala - new method resultComputed.	89

Listing 5.16	Changes to class <code>SparkContext.scala</code> - new method <code>resultComputed.</code>	89
Listing 5.17	Changes to class <code>RDD.scala</code> - modified methods <code>count, collect and reduce.</code>	90
Listing 5.18	Fragment of configuration class <code>Config.</code>	92

ACRONYMS

OS	operating system
JVM	Java Virtual Machine
FIFO	First In First Out
LIFO	Last In First Out
VM	virtual machine
HDFS	Hadoop Distributed File System
RDD	Resilient Distributed Dataset
JT	Job Tracker
TT	Task Tracker
RM	Resource Manager
NM	Node Manager
AM	Application Master
NN	Name Node
DN	Data Node
CLC	container launch context
DAG	Directed Acyclic Graph
I/O	input/output
LXC	Linux Containers
MPI	Message Passing Interface
API	application programming interface
SQL	Structured Query Language
vCPU	virtual CPU
SSD	Solid State Disk
TPC	Transaction Processing Performance Council
QoS	Quality of Service

ABSTRACT

THE need to crunch a steadily growing amount of data generated by the modern applications is driving an increasing demand of flexible computing power, often satisfied by cloud computing solutions. Cloud computing revolutionized the way computer infrastructures are abstracted and used by exploiting virtualization and providing an easy access to resources through the Internet [72].

Big data is a research field that deals with ways to analyze and extract information from data sets containing structured and unstructured data whose size is so large that makes the processing of data with traditional databases and applications very difficult or practically impossible. The processing of these data therefore requires the use of distributed frameworks specialized for the parallel execution of programs, such as Apache Hadoop [5] and Apache Spark [8]. These specialized frameworks are used to transform the applications in atomic parts that can be executed in a distributed cluster of physical or virtual machines.

Hadoop uses the Map-Reduce programming model that was firstly introduced by Google [40] and consists of two distinct phases: Map and Reduce. The former processes and *transforms* a block of data to produce key-value pairs (tuples) as intermediate outputs, the latter *aggregates* these tuples into a final output.

A more advanced solution is represented by Apache Spark [8], that allows for faster executions by avoiding or limiting the use of the persistent storage and it provides a more sophisticated programming model based on parallel execution plans (*PEPs*) which are represented as directed acyclic graph (DAG) of phases.

Big data applications pose new challenges in satisfying requirements on the *Quality of Service (QoS)* provided to end users. In the context of traditional applications (e.g. web) this problem has often been faced using self-adaptive systems that control runtime *KPIs* (Key Performance Indicators) against changes in the application context and workload.

Compared to traditional applications, where a single execution lasts ten to hundreds of milliseconds, big data applications might require minutes or hours to process large datasets, thus *QoS* must consider the execution of single runs. Therefore, in this domain the *QoS* is often defined by means of **deadlines**, that are the maximum allowed duration of a single application execution.

The execution time of big-data applications depends on many factors such as the amount of computing resources available (cpu's,

memory, storage) and the concurrent execution of other applications on the same cluster. Therefore, in the literature one can find several approaches [84, 42, 27] that focus on resource allocation or scheduling techniques in order to reduce deadline violations by using different techniques such as integer linear programming, machine learning and control theory.

All of these approaches use the knowledge of the *PEP* to reason, predict or estimate the application execution time and work under the assumption that the application *PEP* does not change under different input datasets or application parameters. Unfortunately, this is true only if the application code does not contain any conditional branch whose outcome depends on user input values or the result of previous calculations involving input data. Otherwise, different application *PEPs* could be needed to describe all the possible execution flows generated by the combinations of all the different branch outcomes in the application code. Without considering all these *PEPs* the analysis of the application could be significantly unprecise.

xSpark, developed at Politecnico di Milano, is an extension of the Apache Spark framework that offers fine-grained dynamic resource allocation using lightweight containers and control-theoretical planners. It allows users to set application deadlines (which is not possible using Spark) and uses the knowledge of the application *PEP*, retrieved during a profiling phase, to dynamically allocate resources to the application. xSpark does not consider loops or conditional branches in the application code and assumes that the *PEP* is unique.

$xSpark_{SEEP}$, the solution described in this thesis, extends xSpark capability to safely run multi-*PEP* applications, by exploiting symbolic execution. At each decisional branch outcome in the application, $xSpark_{SEEP}$ determines which *PEPs* are still valid and prunes the *PEPs* tree, removing the invalid *PEPs*, thus leaving only the valid ones in the *PEPs* tree. A heuristic is used to select the *PEP* to execute among the valid ones, in order to minimize the risk of missing the deadline while maximizing the CPU usage efficiency.

$xSpark_{SEEP}$ is the result of the integration of *SEEP*, a tool exploiting symbolic execution to discover all possible application *PEPs* produced by different inputs and parameters, with a modified version of xSpark. We tested $xSpark_{SEEP}$ with two applications, Promocalls and Louvain, that uses *GraphX*, a Spark library specialized for graph processing. The evaluation shows that *SEEP* is able to effectively extract all the *PEPs* generated by Spark applications and that $xSpark_{SEEP}$ effectively and efficiently controls the allocation of resources during the execution of *PromoCalls* and *Louvain*, keeping the execution times within considered deadlines with significantly smaller errors and consuming a lower amount of resources than the original version of xSpark.

Since the current solution focuses on controlling a single application, a future work could be directed at extending $xSpark_{SEEP}$ to control multiple concurrent applications.

SOMMARIO

IL bisogno di elaborare una quantità sempre crescente di dati generati dalle moderne applicazioni sta guidando una crescente domanda di potenza di calcolo flessibile, spesso soddisfatta dalle soluzioni di cloud computing. Il cloud computing ha rivoluzionato il modo in cui le infrastrutture informatiche vengono astratte e utilizzate sfruttando la virtualizzazione e fornendo un facile accesso alle risorse attraverso Internet [72].

I big data sono un campo di ricerca che si occupa dei modi per analizzare ed estrarre informazioni da set di dati contenenti dati strutturati e non strutturati, la cui dimensione è così grande che ne rende molto difficile o praticamente impossibile l'elaborazione con database e applicazioni tradizionali. L'elaborazione di questi dati richiede quindi l'uso di framework distribuiti specializzati per l'esecuzione parallela di programmi, come Apache Hadoop [5] e Apache Spark [8]. Questi framework specializzati sono utilizzati per trasformare le applicazioni in parti atomiche che possono essere eseguite in un cluster distribuiti di macchine fisiche o virtuali.

Hadoop utilizza il modello di programmazione Map-Reduce che è stato inizialmente introdotto da Google [40] e si compone di due fasi distinte: Map e Reduce.

Il primo processa e *trasforma* un blocco di dati per produrre coppie chiave-valore (tuple) come output intermedi, quest'ultimo *aggrega* queste tuple in un output finale.

Una soluzione più avanzata è rappresentata da Apache Spark [8], che consente esecuzioni più rapide evitando o limitando l'utilizzo dell'archiviazione persistente e fornisce un modello di programmazione più sofisticato basato su piani di esecuzione paralleli (*PEPs*) che sono rappresentati come un grafico aciclico orientato (DAG) delle fasi.

Le applicazioni di big data pongono nuove sfide nel soddisfare i requisiti sulla *Qualità del Servizio* (*QoS*) fornita agli utenti finali. Nel contesto delle applicazioni tradizionali (ad esempio le applicazioni web), questo problema è stato spesso affrontato utilizzando sistemi autoadattativi che controllano *KPIs* (Key Performance Indicators) a runtime rispetto ai cambiamenti nel contesto dell'applicazione e nel carico di lavoro.

Rispetto alle applicazioni tradizionali, dove una singola esecuzione dura da decine a centinaia di millisecondi, le applicazioni di big data potrebbero richiedere minuti o ore per elaborare grandi dataset, quindi il *QoS* deve considerare l'esecuzione di singole esecuzioni. Pertanto, in questo dominio il *QoS* viene spesso definito attraverso **deadline**, che sono la durata massima consentita della singola esecuzione dell'applicazione.

Il tempo di esecuzione delle applicazioni Big Data dipende da molti fattori come la quantità di risorse di calcolo disponibili (CPU, memoria, storage) e l'esecuzione simultanea di altre applicazioni nello stesso cluster. Pertanto, in letteratura si possono trovare diversi approcci [84, 42, 27] che si concentrano sull'assegnazione delle risorse o sulle tecniche di schedulazione al fine di ridurre le violazioni delle scadenze utilizzando diverse tecniche come la programmazione lineare intera, l'apprendimento automatico e la teoria del controllo.

Tutti questi approcci utilizzano la conoscenza del *PEP* per ragionare, prevedere o stimare il tempo di esecuzione dell'applicazione e lavorare partendo dal presupposto che il *PEP* dell'applicazione non cambi a seconda dei diversi dataset di input o parametri dell'applicazione. Sfortunatamente, questo è vero solo se il codice dell'applicazione non contiene alcun branch condizionale il cui risultato dipende dai valori di input dell'utente o dal risultato di calcoli precedenti che coinvolgono dati di input. In caso contrario, potrebbero essere necessari diversi *PEPs* delle applicazioni per descrivere tutti i possibili flussi di esecuzione generati dalle combinazioni di tutti i diversi risultati che derivano dal codice dell'applicazione. Senza considerare tutti questi *PEPs*, l'analisi dell'applicazione potrebbe essere notevolmente imprecisa.

xSpark, sviluppato al Politecnico di Milano, è un'estensione del framework Apache Spark che offre un'allocazione dinamica delle risorse a grana fine utilizzando contenitori leggeri e pianificatori teorici del controllo. Consente agli utenti di impostare le deadline delle applicazioni (cosa non possibile utilizzando Spark) e utilizza la conoscenza del *PEP* dell'applicazione, recuperato durante una fase di profilazione, per allocare dinamicamente le risorse all'applicazione. xSpark non considera loop o branch condizionali nel codice dell'applicazione e presuppone che il *PEP* sia univoco.

xSpark_{SEEPEP}, la soluzione descritta in questa tesi, estende la capacità di xSpark di eseguire in sicurezza applicazioni multi-*PEP*, sfruttando l'esecuzione simbolica. Ad ogni risultato di un branch decisionale nell'applicazione, *xSpark_{SEEPEP}* determina quali *PEPs* sono ancora validi e pota l'albero dei *PEPs*, rimuovendo i *PEPs* non più validi, lasciando così solo quelli validi nell'albero dei *PEPs*. Una euristica viene utilizzata per selezionare il *PEP* da eseguire tra quelli validi, al fine di ridurre al minimo il rischio di oltrepassare la deadline, massimizzando l'efficienza di utilizzo della CPU.

xSpark_{SEEPEP} è il risultato dell'integrazione di *SEEPEP*, uno strumento che sfrutta l'esecuzione simbolica per scoprire tutti i possibili *PEP* delle applicazioni prodotte da diversi input e parametri, con una versione modificata di xSpark.

Abbiamo testato *xSpark_{SEEPEP}* con due applicazioni, Promocalls e Louvain, che utilizza *GraphX*, una libreria Spark specializzata per l'elaborazione di grafi.

La valutazione mostra che SEEPEP è in grado di estrarre efficacemente tutti i *PEPs* generati dalle applicazioni Spark e che *xSpark_{SEEPEP}* controlla in modo efficace ed efficiente l'allocazione delle risorse durante l'esecuzione di PromoCalls e Louvain, mantenendo i tempi di esecuzione entro scadenze prefissate con errori significativamente minori e consumando una quantità di risorse inferiore rispetto alla versione originale di xSpark.

Poiché la soluzione attuale si concentra sul controllo di una singola applicazione, un lavoro futuro potrebbe essere diretto all'estensione di *xSpark_{SEEPEP}* per controllare più applicazioni concorrenti.

INTRODUCTION

LOUD computing has become a widely used form of service oriented computing, thanks to the provisioning of infrastructures and solutions as a service. Cloud computing has established a totally innovative paradigm of virtualization and use of IT infrastructures. Elasticity (on-demand resource provisioning), metered usage of resources and pay-per-use, no upfront capital investment, fast time-to-market and risk transfer represent most of the most interesting features of cloud computing.

The term "big data" is used to describe a research field that deals with ways to analyze and extract information from data sets containing structured and unstructured data whose size is so large that makes the processing of data with traditional databases and applications very difficult or practically impossible. Data sets can contain a large amount of data that can be structured, like in the traditional relational databases, semi-structured, like in the self-described XML or JSON documents, or unstructured, like in the logfiles collected mostly by web applications to monitor usage or other user's preferences. More properly, big data deals with those data that cannot be handled using traditional database and information technologies.

1.1 CONTEXT

Nowadays, every second 2,780,000 emails are exchanged, 3,690 Skype calls are made, 902 Instagram photos are posted, 73,116 Google searches are performed, 1,502 Tumblr posts are created, and 8,411 Tweets are sent [47]. This is only a small fraction af all the data which is continuously collected and analyzed.

Gartner [34] defines big data as data that contains greater variety arriving in increasing volumes and with ever-higher velocity. This is known as the three V's characterizing big data: Volume, Velocity, Variety [89].

Volume is important because the amount of data drives both the size of memory infrastructure and the computational effort for their analysis. Big data requires processing of high volumes of low-density, unstructured data. These can be data of unknown value, such as Twitter data feeds, clickstreams on a webpage or a mobile app, or sensor-enabled equipment. In some cases, this might be tens of terabytes of data, sometimes hundreds of petabytes.

Velocity is the fast rate at which data is received and acted upon. Normally, highest velocity data streams are directly stored into mem-

ory versus being written to disk. Some internet-enabled smart products operate in real time, or near real time, requiring a very fast evaluation and action.

Variety refers to the many available data types. Traditional data types were structured and are suitable to be stored and managed in a relational database. With the rise of big data, data arrives in an unstructured form. Unstructured and semistructured data types, such as text, audio, and video require an additional preprocessing effort to transform, derive meaning and attach metadata to them.

The systems used to process big volumes of data need to be scalable, tolerate faults, and provide high availability [72].

Scalability means the ability to maintain a near-linear progression between size of data to process and computational resources to perform the task. A big challenge to scalability is the overhead to keep the numerous data fragments of the intermediate results of the data processing steps in synch. Such overhead could drain so many resources to hinder scalability already above modest scale factors.

Fault tolerance is a technology challenge in big data, especially when processing involves many networked nodes and it becomes cumbersome to retain all the checkpoints/restarts to be enacted upon partial processing failures. Devising 100% reliable systems is not an easy task, however the systems can be architected so that the probability of failure falls within the allowed range.

By means of hardware virtualization, cloud computing services satisfies all the requested requisites needed to manipulate big data. Elasticity and redundancy provided by cloud computing also enable big data application high availability, scalability and fault tolerance. Big data also represent an unprecedented business opportunity for many companies which started to deliver big data applications as a service. According to SoftwareTestingHelp [18], these are the top 10 big data companies of 2019: IBM [45], HP Enterprise [44], Teradata [81], Oracle [65], SAP [73], Dell EMC [32], Amazon [13], Microsoft [62], Google [39], VMware [87].

Big data applications are used to transform, aggregate and analyze a large amount of data in an easy and efficient way. These applications are split in atomic parts that can be executed in a distributed cluster of physical or virtual machines by means of specialized frameworks. The limit to the level of parallelism that can be obtain is given by the number of machines and the amount of synchronization needed to keep the data fragments representing the results of intermediate operations synchronized (e.g. aggregations, grouping). This paradigm has been historically represented by the map-reduce programming model firstly introduced by Google [40]. Nowadays, more advanced solutions are available, such as Apache Spark [8] and Apache Tez [10] that provide a greater flexibility and allow building large-scale data processing applications using a *Directed Acyclic Graph* (DAG) or *Parallel*

Execution Plan (PEP) structure to keep track of intermediate results of operations on datasets and determine which parts of the application can be executed in parallel.

The Spark processing platform is very popular [68] and is distinguished by the generality of the solutions and the speed of data processing and sharing. Compared to Hadoop, applications can be run up to 100x faster when operating in memory and up to 10x when operating on disk. Spark is an award-winning platform as the fastest among the open source engines in data sorting operations [57].

Spark embeds fault-tolerance and can run on off-the-shelf hardware. It allows for faster executions by avoiding or limiting the use of the persistent storage and it provides a more sophisticated programming model than the traditional Map-Reduce, based on parallel execution plans (*PEPs*) which are represented as directed acyclic graph (DAG) of phases.

Programming in Spark is made easy by the use high-level functionalities provided by a very rich set of software libraries which compose its API. It embeds specialized libraries for: Spark Core (basic), Spark Streaming, GraphX (graph), MLlib (machine learning), Spark SQL. Spark provides also a seamless integration with POSIX-compliant file system, including HDFS and virtual storage systems, including Amazon S3 and Azure Blob Storage. Spark can be deployed on Hadoop Yarn or Apache Mesos managed clusters, or can run its own cluster manager.

In order to execute big data applications, Spark [92] divides the computation into different phases and split the input dataset into partitions that are stored in a distributed fashion and processed in parallel. Spark provides the API in to support applications written in Java, Python or Scala. It stores in memory the partial results of the operations in order to facilitate its reuse, avoiding recalculations of data sets. The operations performed by Spark on the data can belong to the class of *actions* or the class of the *transformations*. The *actions* return the results of the distributed calculations to the application, while the *transformations* perform the parallel transformation of the data.

Spark operations are gathered in *stages* and the stages are then grouped into *jobs*. A stage is composed by a series of transformations that do not require data shuffling, while a job identifies a set of stages between two actions.

a *Parallel Execution Plan (PEP)* to maximize parallelism during the execution of an application. In effect, a stage is, by definition, executed in parallel, but different stages can also be performed simultaneously. For this reason, Spark materializes as acyclic graphs of stages, while the complete plan of an application is simply the sequence of its works.

A very important measure related to IT applications is the *Quality of Service* (*QoS* in the remainder of this document).

The notion of *QoS* in big data application differ by application type. Interactive applications are usually assessed according to response time or throughput, and their fulfillment depends on the intensity and variety of the incoming requests. Big data applications might require a single batch computation on a very large dataset, thus *QoS* must consider the execution of a single run. In this domain *QoS* is often called deadline, or the desired duration of the computation. Many factors influence the duration of an application execution, surely resource allocation and scheduling greatly influence the duration.

A resource allocation problem arises when applications with different structures run in contexts with different amount of resources or size of input datasets. A scheduling problem arises when many applications run concurrently on the same hardware, so that each application cannot have the totality of the resources assigned to itself. Satisfying deadline-based *QoS* constraints is a problem related to resource allocation, since the amount of allocated resources determines the duration of the execution of Spark applications. The simplest option available on all cluster managers is static partitioning of the resources. This way, each application is given a maximum amount of resources it can use, and holds them for the whole execution time. Memory sharing across applications is currently not provided. Spark also provides a mechanism to dynamically adjust the resources assigned to a specific application according to the workload. Applications may give resources back to the cluster if they are no longer used and re-acquire them again when needed.

xSpark, developed at Politecnico di Milano, is an extension of Spark that offers fine-grained dynamic resource allocation using lightweight containers and enforces *QoS* constraints. xSpark estimates the execution times and allocate resource in order to meet user defined deadlines. A previous work on xSpark has addressed the scheduling problem and established a policy for managing the deadlines when multiple applications run simultaneously on the same hardware.

While we have mentioned availability, fault tolerance and availability as fundamental requirements of big data application, and *QoS* as a measure of the capability to meet a user-defined deadline when processing big data, in this thesis will focus on *QoS*.

1.2 PROBLEM AND MOTIVATION

The growing importance of big data applications has favoured the birth of many analysis techniques to estimate the execution time of Spark applications and perform a proper allocation of resources. For example, Islam et al. [48] propose a solution to statically allocate resources to deadline-constrained Spark applications while minimizing execution

costs. Sidhanta et al. [77] estimate the duration of Spark applications using a closed-form model based on the size of the input dataset and the size of the available cluster. Alipourfard et. al [2] use Bayesian optimization to generate performance models of Spark applications and compute the best configuration for their execution. Baresi et al. [16, 15] propose xSpark, an extension of Spark that exploits control theory and containers¹ to scale allocated resources elastically given the execution times of interest and the other applications on the same cluster.

All the previous works regarding the estimation of the execution times and dynamic provisioning of resources have always assumed the uniqueness of the execution plan for an application, given the computing resources available. This assumption is a simplification of real-world applications, since the actual execution plan is generally different across different program paths when the application code includes conditional branches and loops. Hence, the assumption of a unique *PEP* limits the precision of the analysis and prediction techniques.

xSpark offers optimized and elastic provisioning of resources in order to meet user defined deadlines. This is obtained by using nested control loops. A centralized loop is implemented on the master node and controls the execution of different stages of an application. Multiple local loops, one per executor, focus on task execution. xSpark exploits metadata provided by an initial profiling to create an enriched annotated *Directed Acyclic Graph* (DAG) or *Parallel Execution Plan* (*PEP*) of the application that holds information about the execution of the stages. At runtime, the annotated *PEP* is used to understand how much work has already been done and how much work remains to complete the application. Since all executions of the same application use the same *PEP*, we infer that an xSpark implicit requirement is that applications cannot contain branches or loops, which might be resolved in different ways at runtime.

The xSpark centralized control loop is activated before the execution of each stage and uses a heuristic to assign a stage deadline and calculate the required CPU cores needed to satisfy that deadline, using the information contained in the enriched *PEP* and the user-provided application deadline. Many factors can influence the actual performance and invalidate the prediction. Local control loops have the objective to counteract those imprecisions, by dynamically modifying the amount of CPU cores assigned to the executors during the execution of a stage. A control theory algorithm determines the amount of CPU cores to be allocated to the executor for the next control period. Docker is used to tune the number of CPU cores allocate to the executors, which are run inside lightweight containers [30]. xSpark is able to use less resources than native Spark and can complete executions with less than 1% error in terms of set deadlines.

¹ <https://www.docker.com>.

As mentioned earlier in this document, all these examples are based on the assumption that the application *PEP* is unique for each set of admissible input data and parameter, but this assumption is valid only for the most trivial applications. In fact, the application code can contain branches and loops that can be accessed in different ways depending on the input data and application parameters, and therefore generate different program flows and generate a different *PEP*. As an example we can take the case of a Spark program that evaluates an intermediate result that is then used to determine the outcome of a conditional expressions of a program branch.

Symbolic execution techniques are exhaustively employed in testing of software programs to help identify unsafe inputs that cause the programs to crash. The use of these techniques to deliver a full coverage of the possible dangerous inputs is severely limited by the exponential growth of the required computational resources as a consequence of *path explosion*, due to the need to explore all the possible execution paths and solve all the corresponding *path conditions* of a complex program that has many iterative loops and/or conditional branches in its code. As a consequence, we can rapidly resort to unsolvability. The work of this thesis relies on a lightweight symbolic execution approach [12] that values the solvability of symbolic evaluation-based tools on an *efficiency* base instead of a *soundness and completeness* – often practically unfeasible – base.

In this thesis work we investigate the resource allocation problem when running big data multi-*PEP* applications with deadline-based *QoS* constraints.

The main reason behind this research and thesis work is to address the problem of running multi-*PEP*, deadline-constrained Spark applications obeying the set deadline, in an efficient way.

1.3 SOLUTION AND CONTRIBUTION

To address the questions raised by the investigation of resource allocation problems related to running big data multi-*PEP* applications with deadline-based *QoS* constraints, we propose a solution that leverages *symbolic execution*.

The solution covered by this thesis involves the combined use of a light symbolic execution procedure and search-based test generation to deduce the *PEP* corresponding to a specified set of data and input parameters [12]. The proposed approach is called *SEEP* (*Symbolic Execution-driven Extraction of Parallel Execution Plans*). It leverages a lightweight symbolic execution of the program to extract a set of execution (path) conditions which is representative of the *PEPs* in the application. A search-based test generation algorithm is then used in combination with these set of conditions to create sample datasets that are used to execute each *PEP* and to profile the application. A

prototype tool, also called *SEEPEP*, supports this methodology by identifying the *PEP** of the applications (e.g. the set of *PEPs*) and their associated path conditions and profiling data. It also builds the *PEP** incrementally by using the concrete values of the symbolic variables to update the *PEPs* whose path conditions are satisfiable. This information can then be used to refine the actual *PEP* used to concretely execute the application.

We have integrated the tool into a new version of xSpark called *xSparkSEEPEP* with the aim of systematically test the benefit of *SEEPEP* and to understand how the dynamic selection of the best *PEP* can help a more efficiently resources allocation. *xSparkSEEPEP* contains a specialized component that is dedicated to the selection of the worst-case *PEP*. At each xSpark job boundary, this component injects into the xSpark scheduler the actual worst-case *PEP*, which is chosen among the *PEPs* that are still valid (i.e. their path condition still hold).

In light of the proposed solution, we have identified the following research questions:

RQ₁ : Can the execution of Spark applications be effectively controlled by *SEEPEP*?

RQ₂ : Can the resource allocation capabilities of xSpark be improved, given it used a single, constant *PEP*?

The aim of this work is to give a contribution in terms of knowledge about the application of symbolic execution to the theoretical solution of the identified problem, and a contribution in practical terms by providing:

- i) a modified xSpark platform to enable the runtime management of multi-*PEP* deadline-constrained big data applications and
- ii) a toolchain to identify the applications' execution paths, extract their associated path conditions and generate a path condition evaluator, submit the application and its metadata to the modified xSpark platform and collect performance data for evaluating the QoS of the execution.

RESULTS AND FUTURE WORKS

The solution was tested with two applications: *Promocalls*², a paradigmatic example developed at Deib Labs of Politecnico di Milano, and *Louvain*, a Spark implementation of the *Louvain* algorithm [19], that we downloaded from a highly rated GitHub repository³. *Louvain* exploits *GraphX*, a Spark graph processing library to represent large networks of users and analyze communities in these networks.

² <https://github.com/seepep/promocalls>

³ <https://github.com/Sotera/spark-distributed-louvain-modularity>

We evaluated *SEEP* by integrating it with xSpark to control the parallel execution of two example Spark applications.

The results of the tests confirm the validity of our claim: *xSpark_{SEEP}*, being aware of the different PEPs generated by Spark applications, helps analyze and control their performance/execution time, and thus misses fewer deadlines and allocates resources more efficiently than xSpark.

Future works will address the performance improvement of the profiling phase, by using branch-based selection criteria instead of the simple path-based solution adopted in our presented solution. Another direction to explore in a future work is the application of the identified solution to the control of multi-*PEP* applications with non-strict deadlines.

STATE OF THE ART

2.1 BIG DATA

THE term "big data" refers to a research branch that deals with methodologies to analyze and extract information from data sets containing data that can be structured, like in the traditional relational databases, semi-structured, like in the self-described XML or JSON documents, or unstructured, like in the logfiles collected mostly by web applications to monitor usage or other user's preferences. More properly, we call big data those that cannot be handled using traditional database and software technologies. IDC [46] has predicted that by 2020 one tenth of the world's data will be produced by machines. The organisation forecast that in five years time the amount of connected devices communicating over the internet will reach 32 billion and generate 10% of the world's data. This data is collected and analyzed.

2.1.1 *Batch Processing: Hadoop*

MapReduce is a software framework introduced by Google in order to support the distributed computation of large dataset in cluster of computers. The framework is inspired by map and reduce functions used in functional programming, even though their purpose in the MapReduce framework is not the same as in the original form. There are available MapReduce libraries written in different programming languages.

MapReduce is a programming framework that allows us to perform distributed and parallel processing on large data sets in a distributed environment.¹.

- MapReduce consists of two distinct tasks – Map and Reduce.
- As the name MapReduce suggests, reducer phase takes place after mapper phase has been completed.
- So, the first is the map job, where a block of data is read and processed to produce key-value pairs as intermediate outputs.
- The output of a Mapper or map job (key-value pairs) is input to the Reducer.
- The reducer receives the key-value pair from multiple map jobs.

¹ Extracted from: <https://www.edureka.co/blog/mapreduce-tutorial/>

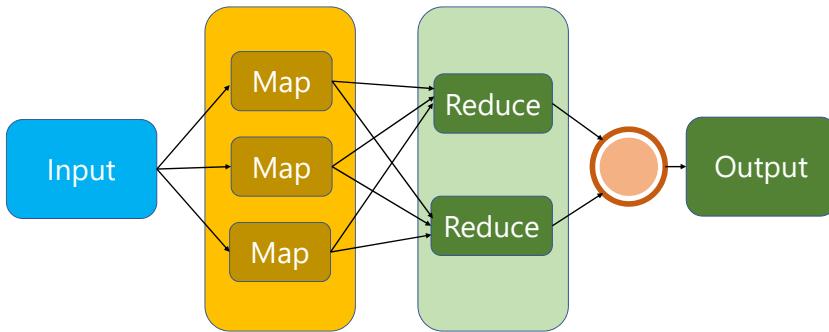


Figure 2.1: Abstract representation of map-reduce paradigm

- Then, the reducer aggregates those intermediate data tuples (intermediate key-value pair) into a smaller set of tuples or key-value pairs which is the final output.

There are open source implementation of the MapReduce framework, for example Apache Hadoop. The MapReduce framework is composed by different functions for each step:

1. Input Reader
2. Map Function
3. Partition Function
4. Compare Function
5. Reduce Function
6. Output Writer

The Input Reader reads the data from mass memory and splits the input in S different splits, with a fixed dimensions (e.g., 64 MB) that are successively distributed to M machines of the cluster that have the Map Function. The Input Reader has also the goal of generating a pair (key, value). The N machine of the cluster are divided in 1 master, whose goal is to detect idling slaves and assign them a task, and $N - 1$ slaves that receive the tasks assigned by the master node. In total, M Map tasks and R Reduce tasks are assigned. A slave that has been assigned the M -th task reads the content of the input, extracts the (key, value) pairs and send them to the Map function defined by the user,

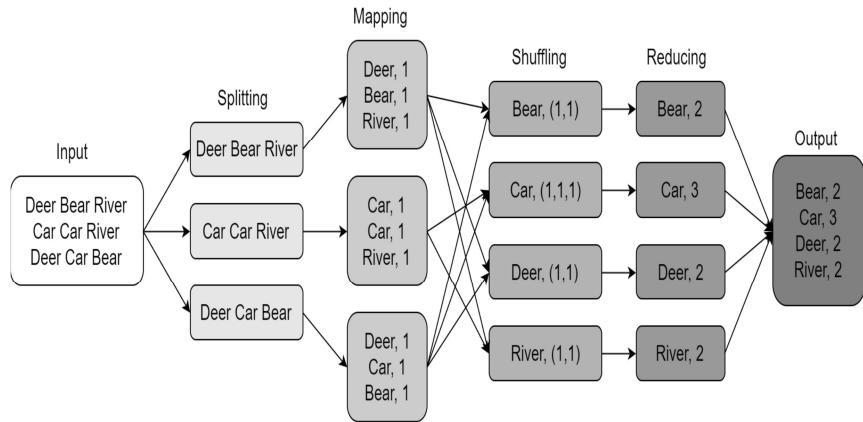


Figure 2.2: Map-Reduce word count job example. Counts the occurrence of each word in input

that generates zero or more (key, value) pairs as output. These pairs are buffered in memory. Periodically the buffered pairs are cached on disk and partitioned in R sections by the partition function. The addresses of the partitioned sections are sent to the master node which is responsible of rotating the location of the slaves that will process the Reduce function. Between the slave with the Map function and the one with the Reduce one, all the pairs are reordered in order to find the ones that point at the same value, and thus also have the same key. The so called shuffling phase is the process that is used to transfer data from mappers to reducers. Once all the keys that point to the same value have been found using the compare function, a merge operation is performed. The sorting operation is useful because in this way the reducer can know when a new reduce task should start. For each of the keys, the associated slave iterates on all the keys, takes the values with the same key and then applies the Reduce function defined by the user, generating one or more element in output. The Output Writer has the goal of writing the results back to mass storage. A sample word count application can be seen in figure 2.2 The input is a document containing words, our goal is to compute the number of occurrence of each of the words in the document. Each Map task applies its function on a line of the document, emitting for each of the words in the line a pair ('word', 1). For example if the input line is "Dear Bear River", it is split into ["Dear", "Bear", "River"] and then mapped into [("Dear", 1), ("Bear", 1), ("River", 1)]. After shuffling the map results, the Reduce task receives a word and a list containing as many ones as the times the word appeared in the document, the

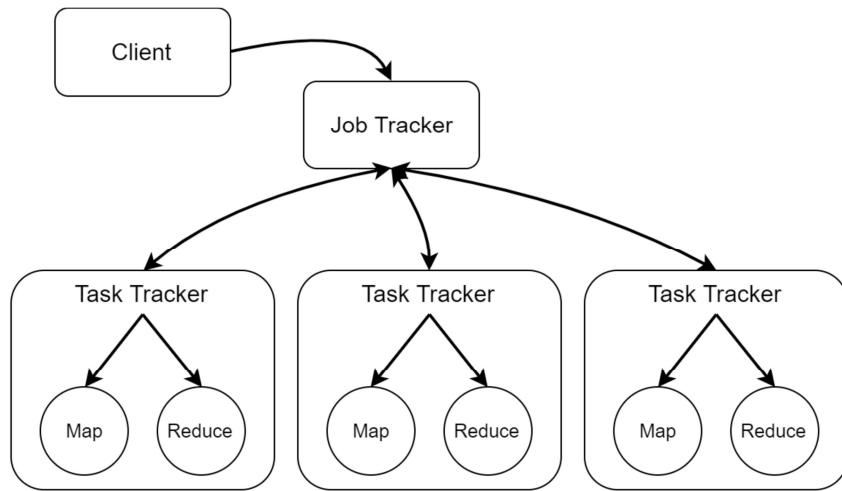


Figure 2.3: Hadoop Map-Reduce Architecture.

reduce function will simply sum the ones in the list, emitting as a result the pair ('word', 'count'). For example, a reducer can receive the key "Bear" with list of values (1, 1), this is reduced into ("Bear", 2). Reducers results are then collected and stored in mass storage.

Apache Hadoop² is an open-source framework for distributed storage and processing of big datasets using MapReduce programming model.

Apache Hadoop MapReduce cluster have a centralized structure composed by a single master Job Tracker (JT) and multiple worker nodes running Task Tracker (TT), as shown in figure 2.3. JT main goal is organizing the job tasks on the slave nodes and continuously monitor the Task Trackers by means of heartbeats. Heartbeats provide a way to retrieve information about the liveliness of the slaves and to inspect the progress of the executions of the different tasks. In order to be fault tolerant, if a task execution fails, it is re-executed possibly on a different slave. JT has the role of the cluster manager, so it needs also to check the admissibility of the submitted MapReduce jobs. TT have the objective of running the assigned task. They reply to heartbeats in order to affirm their liveliness and to update the master about the progress of the assigned tasks. They are configured with a fixed number of map and reduce task slots. As previously introduced, Apache Hadoop also offers a distributed file-system that stores data on different machine, providing an high aggregate bandwidth across the cluster. This functionality is called Hadoop Distributed File System (HDFS)³. It is highly fault tolerant and designed to be deployed on low

² url: <https://hadoop.apache.org/docs/>

³ *HDFS Architecture Guide*. url: https://hadoop.apache.org/docs/r1.2.1/hdfs_design.html

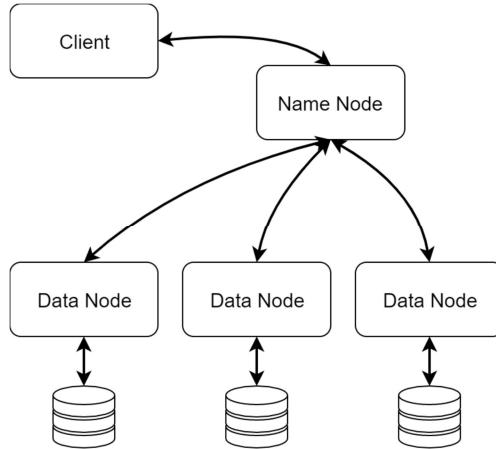


Figure 2.4: HDFS Node Structure.

cost hardware. HDFS exposes a filesystem namespace and allows user data to be stored in files and retrieved. Cluster structure is similar to the one of MapReduce cluster, with one master and multiple slaves, as we can see from figure 2.4. The master is composed by a single Name Node (NN), that manages the file system namespace and regulates access to the files by clients. NN executes filesystem operations such as opening, closing, renaming files and directory, but the most important operation performed is keeping track of the mapping between blocks and Data Node (DN). Indeed a file stored in HDFS is split into one or more blocks, and those blocks are stored in the Data Node (DN). Data Node (DN) represent the slaves, they are usually one per node, and manage the storage that is attached to the node they are running on. They are responsible for serving read and write operation requests from the clients, but also can perform block creation, deletion and replication. Block replication is a significant way to improve fault tolerance.

2.1.2 Batch Processing: Spark

Apache Spark is an open source framework for distributed computation [8], that provides an interface for programming entire clusters with implicit data parallelism and fault tolerance. With respect to the MapReduce paradigm, the in-memory multilevel primitives of Spark allow to have performances up to 100 times better in certain applications. Spark can work as standalone or on a cluster manager such as Apache Hadoop Yarn or Apache Mesos. It also needs a distributed storage and can natively use HDFS and other solutions.

Spark has been designed as a unified engine for distributed data processing. Its programming model resembles the one of MapReduce, but it is extended with a data sharing abstraction called Resilient Distributed Dataset (RDD). Using this abstraction, a wide range of processing workloads can be captured, including SQL, streaming, machine learning and graph processing. The generality of the Spark approach gives great benefits: firstly, applications are easier to develop since there is a unified API, secondly, it is a lot easier to combine processing tasks. With previous distributed computation framework we needed to write data to mass storage before using them in another engine. Spark, instead, can reuse the same data, often kept in memory. The programming abstraction at the foundation of Spark is the Resilient Distributed Dataset (RDD), that is a fault-tolerant collection of objects partitioned across the cluster and can be processed in parallel. The users create RDD's by applying operations called "transformations", such as map, filter and group-by, on the data. RDD's can be backed by a file obtained from an external storage. Spark evaluates RDD's in a lazy mode. This allows the construction of an efficient plan to execute the computation requested by the user. In particular, every transformation operation returns a new RDD, that is the representation of the result of the computation, but the computation is not executed immediately after the transformation request is encountered, but only when a Spark action is met. When an "action" is requested by the user code, Spark checks the entire graph of the transformation and uses it to create an efficient execution plan. For example, if there are many filters and maps in a row, Spark can merge them together and execute as a single operation. RDDs also offer an explicit support to data sharing among the computations that are ephemeral by default, but can be persisted to disk or memory for rapid reuse. This data sharing is one of the main differences between Spark and the previous computing models like MapReduce, because all the other operations that Spark can perform are similar to the ones of MapReduce. The data sharing capability allows huge speedups, up to 100 times, in particular when executing interactive queries and iterative algorithms. RDDs can also recover automatically from a failure. Traditionally, fault tolerance in distributed computing was achieved by means of data replication and checkpointing. Spark instead uses a different approach called lineage. Each RDD keeps track of its transformation graph used to generate the RDD and re-executes the transformation operations on the base data to recover every lost partition. The data recovery based on lineage is significantly more efficient than replication in case of data-intensive workload. In general, recovering lost partitions is faster than re-executing the entire program. Spark was designed to support different external systems for persistent storage, usually it is used in conjunction with a clustered file system like HDFS. Spark is

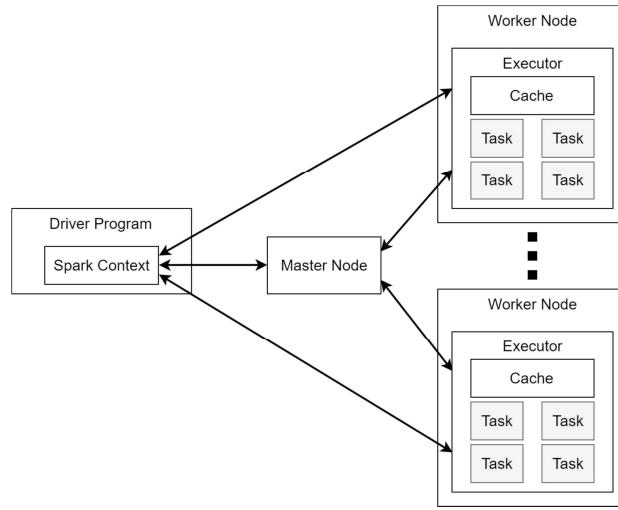


Figure 2.5: Spark Standalone Architecture.

designed as a storage-system-agnostic engine, to make it easy to run computations against data from different sources.

Different high-level libraries have been developed in order to simplify the creation of programs that can run in Spark framework.

- SQL and DataFrames: support for relational queries, that are the most common data processing paradigm
- Spark Streaming: implements incremental stream processing using a model called "discretized streams", input data is split into micro batches
- GraphX: graph computation interface
- MLlib: machine learning library, more than 50 common algorithms for distributed model training

Spark architecture follows the master/worker paradigm (figure 2.5). A master server accepts data and processing request, split them into smaller chunk of data and simpler actions that can be handled in parallel by the multiple workers. A Spark application is executed inside a driver program, that makes the user code executable on the computing cluster using a `SparkContext`. The driver program is responsible for managing the job flow and scheduling tasks that will run on the executors. The `SparkContext` will split the requested operations in tasks that can be scheduled for distributed execution on the workers. When a `SparkContext` is created, a new Executor process is created on each worker. An executor is a separate Java Virtual Machine (JVM) that runs for the entire lifetime of the Spark application, executes tasks using a

Details for Job 0

Status: SUCCEEDED
Completed Stages: 2

Event Timeline DAG Visualization

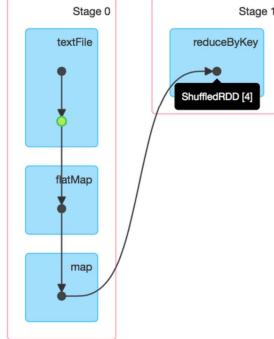


Figure 2.6: Spark DAG Example.

thread pool and store data for its Spark application. Communication between the `SparkContext` and the other components is performed using a shared bus. When an application is submitted to Spark, it is divided in multiple jobs. Jobs are delimited by `Spark` actions in the application code. `Spark` actions are those operations that return a value to the driver program after running a computation on the dataset. For each job, a `Directed Acyclic Graph` (DAG) is created to keep track of the `RDD`s that are materialized inside the job. DAG nodes represent the `RDD`s, meanwhile arcs represent transformations, that are the operations that create new datasets from existing ones. The application steps inside a single job are further organized into stages, that are delimited by operations require data reshuffling, that will inevitably break locality. `Spark` distinguishes between narrow transformations, that do not reshuffle data (e.g., `map`, `filter`), and wide transformations, that require data reshuffling (e.g., `reduceByKey`). Stages are also used to produce intermediate result that can be persisted to memory or mass storage to avoid re-computation. When all stages inside a job have been identified, `Spark` can determine which parallel tasks need to be executed for each stage, and schedule them for operation on the executors. `Spark` creates one task for each partition of the `RDD` received in input by a stage.

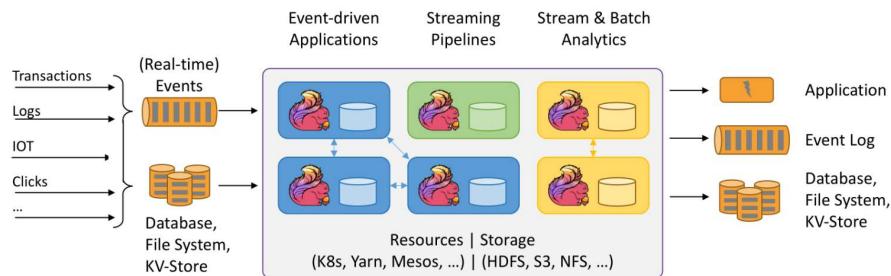
In Figure 2.6, we can see a simple DAG representing the single job of the word count application presented in Listing 2.1 [64]. The image was obtained from `SparkWeb` UI. Through a `textFile` operation, the input file is read from HDFS. Then a `flatMap` operation is applied to split each of the lines of the document into words. Following, another `map` is used to create ('word', 1) pairs. Finally, a `reduceByKey`

Listing 2.1: Spark word count application example.

```

1 sparkContext.textFile("hdfs :/. . . ")
2 .flatMap(line => line.split(" "))
3 .map(word => (word, 1)).reduceByKey(_ + _)
4 .saveAsTextFile("hdfs :/. . . ")

```

**Figure 2.7:** Apache Flink® - Stateful Computations over Data Streams.

operation is performed to count the occurrences of each word. The blue boxes represent the Spark operations that the user calls in his code, while the dots represent the RDDs that are created as a result of these operations. Operations are grouped into stages, represented by the boxes with a red border. The job has been divided into two stages because the `reduceByKey` transformation requires the data to be shuffled. The green dot represents a cached RDD, in particular the data read from HDFS has been cached, in this way future computations on this RDD can be done faster since data will be read from memory instead of HDFS. The default deployment of Spark is in standalone mode, that is using its embedded cluster manager.

2.1.3 Streams Processing: Flink

Apache Flink [3] is a framework and distributed processing engine for stateful computations over unbounded and bounded data streams. It can run in any commonly used cluster environments, computes at in-memory speed and manages data at any scale.

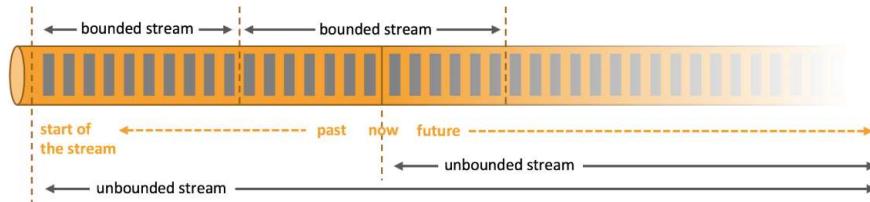


Figure 2.8: Apache Flink® - Bounded & Unbounded Data Streams.

Flink's architecture is suitable for processing unbounded and bounded data. Any type of data is generated as a stream of events. Sensors data, server logs or user interactions on a website or mobile application, credit card transactions, all of these data are generated as streams of data.

Data can be processed as unbounded or bounded streams [4].

Unbounded streams have a start but no defined end. They do not terminate and provide data as it is generated. Unbounded streams must be continuously processed, i.e., events must be promptly handled after they have been ingested. It is not possible to wait for all input data to arrive because the input is unbounded and will not be complete at any point in time. Processing unbounded data often requires that events are ingested in a specific order, such as the order in which events occurred, to be able to reason about result completeness.

Bounded streams have a defined start and end. Bounded streams can be processed by ingesting all data before performing any computations. Ordered ingestion is not required to process bounded streams because a bounded data set can always be sorted. Processing of bounded streams is also known as batch processing.

Apache Flink excels at processing unbounded and bounded data sets. Precise control of time and state enable Flink's runtime to run any kind of application on unbounded streams. Bounded streams are internally processed by algorithms and data structures that are specifically designed for fixed sized data sets, yielding excellent performance.

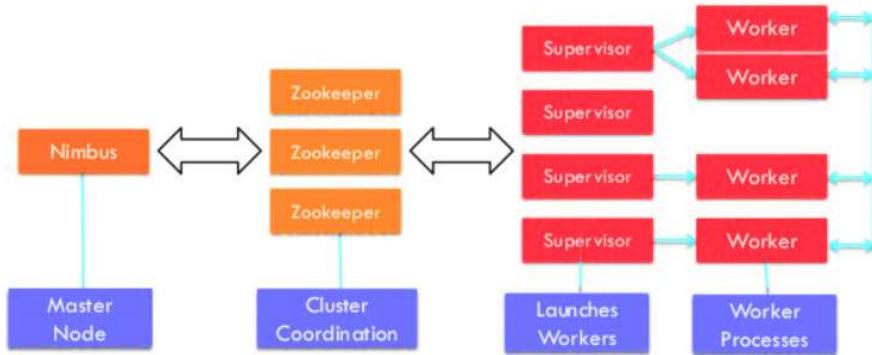


Figure 2.9: Apache Storm.

2.1.4 Streams Processing: Storm

Apache Storm [9] is a free and open source distributed realtime computation system. Storm makes it easy to reliably process unbounded streams of data, doing for realtime processing what Hadoop did for batch processing. Storm is simple, can be used with any programming language, and is easy to use.

Storm has many use cases: realtime analytics, online machine learning, continuous computation, distributed RPC, ETL, and more. Storm is fast: a benchmark clocked it at over a million tuples processed per second per node. It is scalable, fault-tolerant, guarantees your data will be processed, and is easy to set up and operate.

Storm integrates with the queueing and database technologies you already use. A Storm topology consumes streams of data and processes those streams in arbitrarily complex ways, repartitioning the streams between each stage of the computation however needed.

The Apache Storm architecture is quite similar to that of Hadoop. However there are certain differences which can be better understood by getting a closer look at its cluster in Figure 2.9: **Nodes** - There are two types of nodes in Storm cluster similar to Hadoop:

- Master node - The master node of Storm runs a daemon called 'Nimbus', which is similar to the 'Job Tracker' of Hadoop cluster. Nimbus is responsible for distributing codes, assigning tasks to machines and monitoring their performance.
- Worker node – Similar to the master node, the worker node also runs a daemon called 'Supervisor' which is able to run

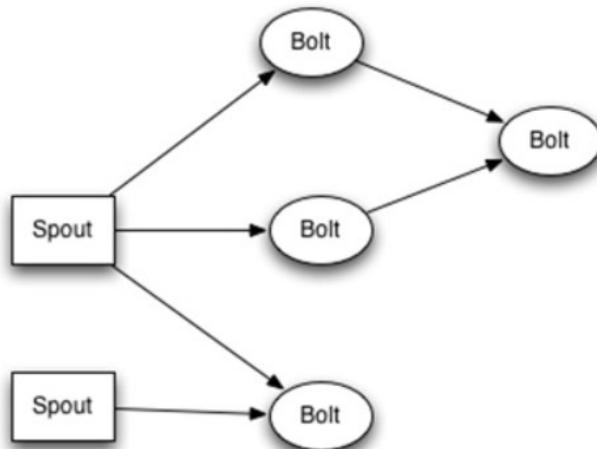


Figure 2.10: Apache Storm Components Abstraction.

one or more worker processes on its node. Each supervisor works assigned by Nimbus and starts and stops the worker processes when required. Every worker process runs a specific set of topology which consists of worker processes working around machines. Since Apache Storm does not have the abilities to manage its cluster state, it depends on **Apache Zookeeper** for this purpose. Zookeeper facilitates communication between Nimbus and Supervisors with the help of message acknowledgements, processing status, etc.

Storm Components/Abstractions (see Figure 2.10) – There are basically four components which are responsible for performing the tasks:

- Topology – Storm Topology can be described as a network made of **spouts** and **bolts**. It can be compared to the Map and Reduce jobs of Hadoop. Spouts are the data stream source tasks and Bolts are the accrual processing tasks. Every node in the network consists of processing logic's and links to demonstrate the ways in which data will pass and the processes will be executed. Each time a topology is submitted to the storm cluster, Nimbus consults the supervisor nodes about the worker nodes.
- Stream – One of the basic abstractions of the storm architecture is stream which is an unbounded pipeline of tuples. A tuple can be defined as the fundamental component in the Storm cluster containing a named list of the values or elements.

- Spout – It is the entry point or the source of streams in the topology. It is responsible for getting in touch with the actual data source, receiving data continuously, transforming those data into actual stream of tuples and finally sending them to the bolts to be processed. two types of nodes in Storm cluster.
- Bolt - Bolts keep the logic required for processing. These are responsible for emitting the streams for processing by other bolts and saving or sending the data for storage. These are capable of running functions, filtering tuples, aggregating and joining streams, linking with database, etc.

2.2 RUNTIME MANAGEMENT OF BIG DATA APPLICATIONS

Big data applications requires system scalability, fault tolerance and availability [72].

Scalability means the ability to maintain an approximate linear relationship between the size of processed data and the amount of consumed resources.

Fault tolerance is a challenge, especially when systems are complex and involve many networked nodes. By means of hardware virtualization, cloud computing services satisfies all the requested requisites, and the *elasticity* and redundancy it provides also enable big data application high availability, scalability and fault tolerance.

Another very important feature of big data applications is the *Quality of Service* or *QoS*.

QoS definition for IT applications differ by application type. Interactive applications are usually assessed according to response time or throughput, and their fulfillment depends on the intensity and variety of the incoming requests.

Big data applications might require a single batch computation on a very large dataset, thus *QoS* must consider the execution of a single run. In this domain *QoS* is often called *deadline*, or the desired duration of the computation.

We have mentioned availability, fault tolerance and availability as fundamental requirements and *QoS* as a measure of the capability to meet a user-defined deadline when processing big data. Many factors influence the duration of an application execution, surely resource allocation greatly influences the duration.

The challenges introduced by big data require resilient, flexible and self-adapting software systems [54]. Hence, *autonomic systems* and *self-adaptation* has increasingly captured the attention of researchers [88]. These systems automatically react to changes in the environment, or in their own state, and change their behaviour to satisfy functional and non-functional requirements. Meeting requirements in complex and variable execution environments is a difficult task that can be tackled at design time or at runtime. At runtime the adaptation is very

often obtained by using a well-known process called MAPE [51], a control loop composed of four phases: monitoring, analysis, planning and execution.

One of the challenges that modern software systems face is the provisioning and optimization of resources to meet a varying demand, generated by increasingly common phenomena like fluctuating workloads, unpredictable peaks of traffic and unexpected changes. Service providers cannot disregard these factors if they want to cope with the challenge of satisfying functional and non-functional requirements, usually defined in SLAs (Service Level Agreements). Hence the need arises for an automatic adjustment of system resources allocation to avoid resource saturation and unresponsiveness, users dissatisfaction and unnecessary costs. This paradigm is called elastic resource provisioning [31, 93, 74, 41].

Many approaches about elastic systems and dynamic resource allocation were proposed both in the industry and in academia. In modern technology elasticity is often enabled by cloud computing that gives to an application a theoretical infinite degree of scalability. However considering only resources is rather restrictive because of the many factors that impact application during their runtime life-cycle. In fact Dustdar et al. [31] argue that elastic computing should be designed by considering three dimensions: quality, resources and cost. Quality elasticity considers how quality is affected by a change in resource availability. Instead cost elasticity measures how resource provision is affected when a change in cost happens.

Cloud computing services provide the needed level of fault tolerance and availability required by big data applications. In the remainder of this chapter we present an overview of popular big data frameworks that can leverage cloud computing solutions and how they address elastic resource allocation to satisfy functional and non-functional application requirements. The resulting scenario represent the base for our work, where we will consider quality elasticity and not cost elasticity aspects.

This thesis shows how the application of lightweight symbolic execution techniques to deadline-based *QoS* constrained multi-DAG big data applications helps reduce the number of deadline violations and allocate resources more efficiently.

2.3 ELASTIC RESOURCE PROVISIONING

The word *elasticity* comes from physics and is defined as the ability of an object or material to resume its normal shape after being stretched or compressed. In computing, elasticity has a similar meaning and is used to characterize *autonomic systems*. Herbst et al. [41] defines elastic provisioning as reported below.

Elasticity is the capability of a system to adapt to workload changes by provisioning or de-provisioning resources automatically such that at each point in time the available resources match the current demand as closely as possible.

A system is in an *under provisioned* state if it allocates less resources than required by the current demand; it is in an *over provisioned* state if allocates more resources than required. Moreover, elasticity is determined by four attributes:

- *Autonomic Scaling*: the adaptation process used to control the system.
- *Elasticity Dimensions*: the set of scaled resources in the adaptation process.
- *Resource Scaling Units*: the minimum amount of allocable resources to each dimension.
- *Scalability Bounds*: the lower and the upper bound on the amount of resources that can be allocated to each dimension.

Additionally, two aspects must be considered in evaluating the elasticity degree of a system: speed and precision.

Speed of scaling up/down: the time it takes to switch from an underprovisioned/overprovisioned state to an optimal or overprovisioned/underprovisioned state respectively.

Precision of scaling: the absolute deviation of the current amount of allocated resources from the actual resource demand.

Scalability and efficiency are terms related to elasticity, nevertheless they differ by the following aspects:

Scalability, although is a prerequisite for elasticity, it does not consider the temporal aspects (how fast and how often) and the granularity of the adaptation actions.

Efficiency, contrary to elasticity, it takes in account all the types of resources employed to accomplish a certain amount of work, not only the resources scaled by the adaptation actions.

Elasticity reflects the (theoretical) infinite upper bound on resource scalability in cloud computing and the frictionless resource renting model. Nevertheless elasticity is not just a synonym of resource management. Elasticity is also related to trade off between cost and quality [31]. Cost elasticity describes how resources are managed in response to cost changes, while quality elasticity measures how responsive is the quality to changes in resource utilization.

A *matching function* $m(w) - r$ is a system specific function that gives the minimum quantity of resources r for any given resource type needed to meet the system's performance requirements at a certain workload intensity. A matching function is required for both up and down scaling directions.

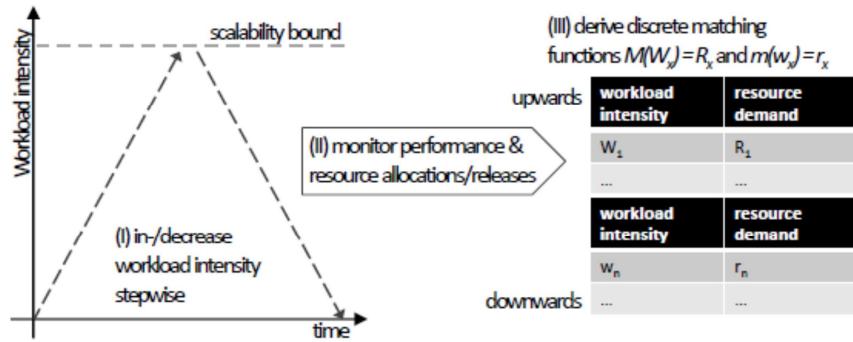


Figure 2.11: Elasticity Matching Function derivation.

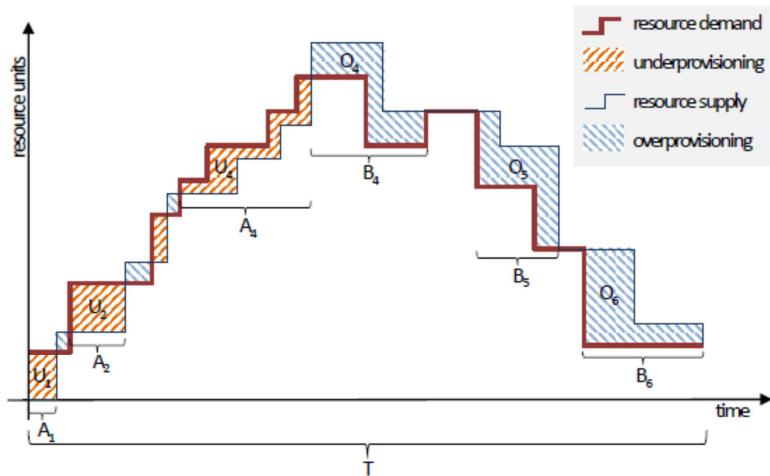


Figure 2.12: Resource Provisioning Chart.

The matching functions can be determined based on measurements, as illustrated in Figure 2.11, by increasing the workload intensity w stepwise, and measuring the resource consumption r , while tracking resource allocation changes. The process is then repeated for decreasing w . After each change in the workload intensity, the system should be given enough time to adapt its resource allocations reaching a stable state for the respective workload intensity. As a rule of thumb, at least two times the technical resource provisioning time is recommended to use as a minimum. As a result of this step, a system specific table is derived that maps workload intensity levels to resource demands, and the other way round, for both scaling directions within the scaling bounds.

An example of how *speed* and *precision* affect resource provisioning is shown in Figure 2.12.

2.4 SPARK RESOURCE PROVISIONING

It's important to remember that the cluster manager is responsible for starting executor processes and determine where and when they will run. Using Spark's embedded cluster manager might be a problem in terms of resource utilization when we want to execute different distributed applications at the same time. Using a single cluster manager for different distributed applications has the advantage of providing a global view on which applications are running and which we want to execute inside the cluster. Without a single cluster manager, we can have two main approaches in order to perform resource sharing and allocation:

- allowing every application to allocate all the resources in the cluster at the same time, this leads to an unfair situation of resource contention
- splitting the resource pool into smaller pools, one per application.

In this way we will avoid resource contention but we will have a less efficient utilization of the resources, because some of the applications might request more resources than the ones in the pool, meanwhile some others are using less resources than the allocable ones in order to execute. A more dynamic way of allocating resources will lead to a better resource utilization. Spark natively support executing on top of Apache Hadoop YARN and Apache Mesos cluster managers. Spark supports the dynamic allocation of executors, also known as elastic scaling, this feature allows to add and remove Spark executors in a dynamic way in order to match the workload. In traditional static allocation, a Spark application would allocate CPU and memory upon starting the execution, disregarding how much resources will effectively use later on. With dynamic allocation instead it is possible

to allocate as much resources as they are necessary, in order to avoid wasting them. The number of running executors is scaled up and down according to the workload, in particular idling executors are removed and when there are tasks waiting to be executed and new executors are launched. Dynamic allocation can be activate in Spark settings and should be used together with the External Shuffle Service, in this way data that have been manipulated from the executor is still available after the removal of the executor. Dynamic allocation has two different policy for scaling the executors:

- Scale Up Policy: new executor are requested when there are pending tasks, the number of executors is increased exponentially because they start slow and so the application might need a slightly higher number of them
- Scale Down Policy: idling executors are removed after a certain amount of time, this amount of time can be configured

In order for dynamic allocation to work, we must configure it, by setting the initial number of executors that are created when application starts and the minimum and maximum number of executors that can be reached when scaling down and up respectively. Dynamic allocation is available on all cluster managers currently supported by Spark, even in Standalone mode.

2.4.1 *Apache Hadoop Yarn*

Apache Hadoop YARN (Yet Another Resource Negotiator) is the next generation of Hadoop's compute platform[83]. The idea is to split the functionality of resource management and job scheduling and monitoring. This is achieved by having two different kind of daemons running, one global Resource Manager (RM) and a per-application Application Master (AM). Resource Manager (RM) and Node Manager (NM) form the data computation framework (Figure 1.4). RM is the authority that manages resources among all the applications that are running in the system, meanwhile NM is the per-machine daemon who is responsible for managing containers, monitoring and reporting. The perapplication AM has the goal of negotiating resources with RM and working with NM in order to execute and monitor tasks. Resource Manager (RM) is composed by two components: Scheduler and Applications Manager. The Scheduler is responsible for allocating resources to the various applications that are running, by taking into account constraints about capacity, queues, etc. It is a real scheduler in the sense that it does not perform monitoring or tracking of the application state. Moreover, it does not offer any guarantee that a failed application will be reexecuted after an application or hardware failure. The Scheduler performs the allocation according to the resources that are requested by an application; this is based on the abstract notion

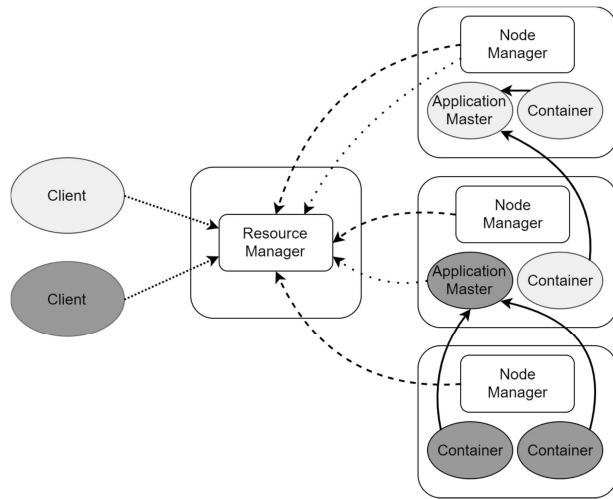


Figure 2.13: Apache Hadoop YARN Architecture.

of container which has elements as memory, CPU cores, disk and network bandwidth. There are pluggable policy that determine the repartition of resources among the different applications, for example we have the Capacity Scheduler, designed for multi-tenant clusters, and the Fair Scheduler, that shares cluster resources fairly. The Applications Manager is responsible of accepting the submission of a job, it negotiates the first container that will execute the AM and it offers a service that can be used to restart the AM in case of failure. The per-application AM has the goal of negotiating the needed containers from the Scheduler, track their status and monitor their progress. The RM keeps a global model of the cluster state and thanks to the resource requirements reported by the running applications, it makes possible to enforce a global scheduling, but it is required to have an accurate understanding of the applications' resource requirements. In response to AM requests, the RM generates containers together with tokens that grant access to resources. An extension of the protocol allows the RM to ask back resources from applications, for example when cluster resources become scarce. Application Master (AM) is the process that coordinates the execution of an application inside the cluster. It is important to remember that AM itself is run in the cluster, just like any other container. Periodically, an heartbeat is sent to the RM in order to confirm its liveliness and to update the Scheduler about its resource requests. After having modeled the application requirements, the AM encodes its preferences and constraints inside the heartbeat message. This information are stored in the form of Resource Request, containing the desired number of containers (e.g., 100 container), the resources of each container (e.g., <2 CPU, 2 GB>), the locality pref-

erences and the priority of this resource request with respect to the other ones of this application. When a container lease is received, the AM can choose to modify its execution plan in order to take into account the abundance or scarcity of resources. Node Manager (NM) is the worker daemon in YARN, its purpose is to authenticate container lease, manage dependencies, monitoring the execution of containers and offer them a set of services. After having registered with the RM, the NM sends heartbeats in order to communicate its status and receives instructions from the RM. All the containers are described by a container launch context (CLC), that keeps track of all the environment variables, the dependencies, the security tokens, but also of the payloads needed by NM services and the commands that are needed to launch the process inside the container. After having validated the authenticity of the container lease, the NM configures the container with the specified resource constraints and initializes a monitoring subsystem. In order to launch the container, dependencies are copied into local storage. NM also has the duty of killing container upon a request from RM or AM, for example when a tenant is being evicted or when an application completes. Whenever a container exits, NM needs to clean the working directory. When an application ends, all the resources held by its container on all nodes are released. NM periodically checks the state of the physical machine and informs the RM of a possible unhealthy state.

2.4.2 *Apache Mesos*

Apache Mesos is an open-source project used to manage computer clusters. The purpose of Mesos is to share cluster between different computing frameworks, such as Apache Hadoop or Message Passing Interface (MPI). The sharing increments the utilization of the cluster and prevents per-framework data replication. Mesos shares resources in a fine-grained way, allowing to achieve data locality. It presents a scheduling mechanism on two layer called resource offers. Mesos decides how many resources to offer to each of the running frameworks, meanwhile they decide how many resources to accept and which computation to execute on the granted resources. New cluster computing frameworks continue to emerge, it is clear that finding a framework that is optimal for all type of application is almost impossible. We expect that organization would like to use different frameworks inside the same cluster, picking the best one according to the kind of application that they are going to execute. Two classic solution are: i) statically partitioning the cluster and executing one framework per partition; ii) allocate a set of VMs to each of the frameworks. Unluckily these solution do not achieve high utilization and efficient data sharing. The main problem is the different allocation granularity of these solutions and the one of the existing frameworks, for example Hadoop employs

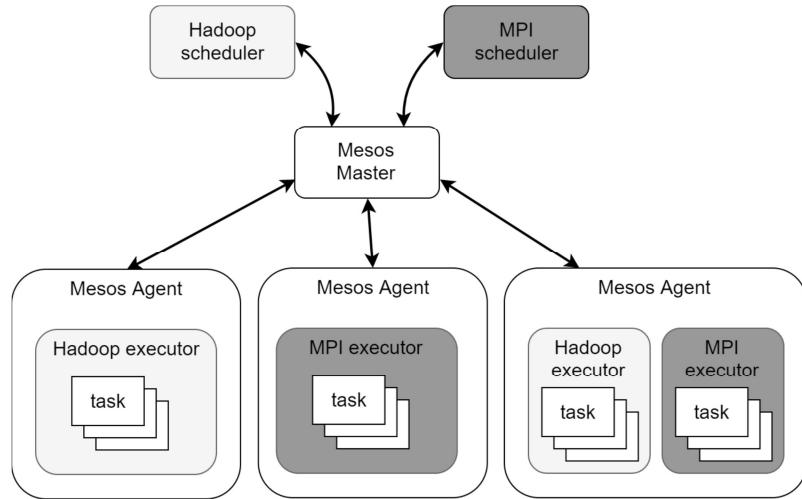


Figure 2.14: Apache Mesos Architecture.

a fine grained resource sharing model, where nodes are divided into slots and each job is composed by short tasks that match the slots. The presence of short tasks allows us to achieve high utilization, as jobs can rapidly scale when new nodes are available. But it is not possible to achieve fine grained sharing across frameworks, because they have been developed in an independent way, and thus it is difficult to efficiently share the cluster among different frameworks. Mesos delegates the control over the scheduling to the different frameworks. In this way it is possible to have the abstraction of the resource offers, that encapsulate a bundle of resources that the framework can allocate on a node in order to execute a task. Mesos decides how many resources to offer to each framework, this is based on policies, and the framework decides which resources to accept and which tasks to execute on them. Even though this approach does not lead to a globally optimum scheduling, it has been proved that it performs particularly well in practice, allowing the frameworks to obtain near perfect data locality. Mesos provides other benefits to its users, for example the possibility of running different instances of the same framework or even different versions. Mesos is composed by a master process that manages slave daemons running on each cluster node and frameworks that run tasks on these slaves, as we can see from figure 2.14 . Master implements fine-grained sharing across frameworks using resource offers. Every resource offer is a list of free resources on the different slave nodes. The master decides how many resources to offer to each framework, according to some policy such as fairness or priority.

Apache Mesos Resource Offer example. 1) Mesos Agent 1 reports free resources to the Allocation Module; 2) Allocation Module offers

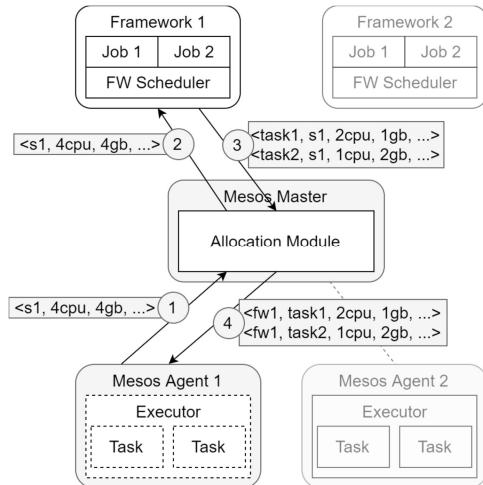


Figure 2.15: Apache Mesos Resource Offer example. 1) Mesos Agent 1 reports free resources to the Allocation Module; 2) Allocation Module offers resources to Framework 1 scheduler; 3) Framework 1 scheduler accepts resources and assign tasks; 4) Allocation Module launches tasks on the executor running in Mesos Agent 1.

resources to Framework 1 scheduler; 3) Framework 1 scheduler accepts resources and assign tasks; 4) Allocation Module launches tasks on the executor running in Mesos Agent 1.

Every framework that is running on Mesos is composed by two components: a scheduler, that registers with the master in order to obtain the resource offers, and an executor process that is launched on the slave node in order to execute framework's tasks. While the master chooses how many resources to offer, the scheduler chooses which resources to use among those offered. When an offer is accepted, the scheduler sends to the master the description of the tasks that should be executed. The resource offer process is repeated every time tasks are finished and when there are new free resources. In order to maintain a light interface, Mesos does not ask the frameworks to specify their resource requirements or constraint, instead it gives them the possibility of refusing offered resources. Mesos allows frameworks to set up a set of filters, in the form of boolean predicates, specifying the conditions on which the framework will always refuse a proposal (e.g., providing a whitelist of nodes it can run on). In Figure 1.6 we have an example of resource offer process. 1. Agent 1 reports to master that it has 4 CPUs and 4 GB of memory free. Master invokes its allocation module policy, which tells that framework 1 should be

offered all the resources. 2. Master sends a resource offer describing the resources available on agent 1 to framework 1.

2.4.3 *Spark on Yarn*

subsec:sparkOnYarn Support for running Spark on YARN was added to Spark in version 0.6.0 and has been improved in subsequent releases [71]. When running on YARN, each Spark executor is run inside a YARN container. Spark supports two different modes to run on YARN, the Yarn-cluster and Yarn-client mode. In client mode, as shown in figure 2.16, the driver program is run inside the client process. In this way, the Application Master (AM) that is run in a YARN container is used only to request resources to the Resource Manager (RM). This mode is useful for interactive applications and for debugging purposes, since you can see applications' output immediately on the client side process. If the client disconnects from the cluster, the Spark application will terminate, this is due to the fact that the driver process resides on the client. In cluster mode instead, as shown in Figure 1.10, Spark driver program is run inside the AM process managed by YARN. After initializing the application, client can disconnect from the cluster and reconnect later on. This mode makes sense when using Spark on YARN in production jobs. Running on top of YARN cluster manager has some benefits. First of all YARN allows to dynamically share the cluster resources between the different frameworks that are running together. For example we can run MapReduce jobs after running Spark jobs without the need of changing YARN configurations. Moreover, YARN supports for categorizing, isolating and prioritizing workloads and employs security policies, in this way Spark can use secure authentication between its processes.

When running on YARN, Spark executors and driver program use about 6-10% more memory with respect to the standalone execution, this is due to the fact that this extra amount of off-heap memory is allocated in order to take into account YARN overheads.

2.4.4 *Spark on Mesos*

subsec:sparkOnMesos Support for running Spark on Mesos was added to Spark in version 1.5. Spark on Mesos can be executed in two different modes: coarse-grained and fine-grained [70]. In coarse-grained mode, as shown in figure 2.18, each Spark application is submitted to Mesos master as a framework and Mesos slaves will run tasks for the Spark framework that are Spark executors. Mesos tasks are launched for each Spark executor and those Mesos tasks stay alive during the lifetime of the application unless we are using dynamic allocation or the executor is killed for various reasons. The advantage of coarse-grained mode is in a much lower task startup overhead, with

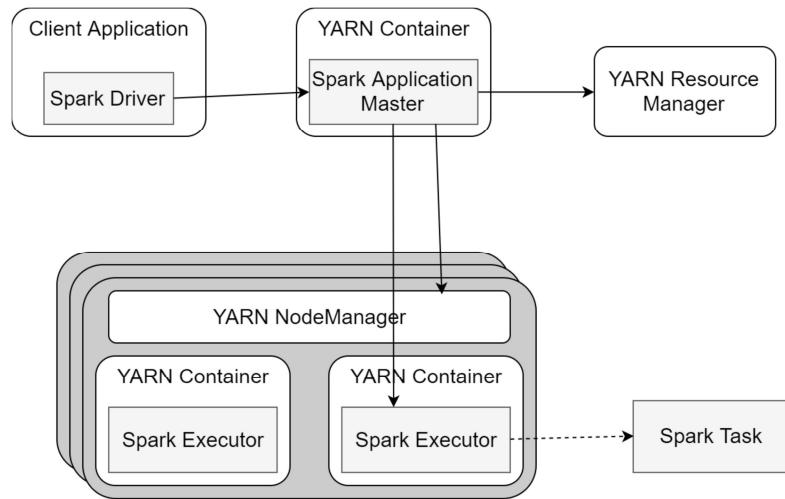


Figure 2.16: Spark on YARN Client Mode.

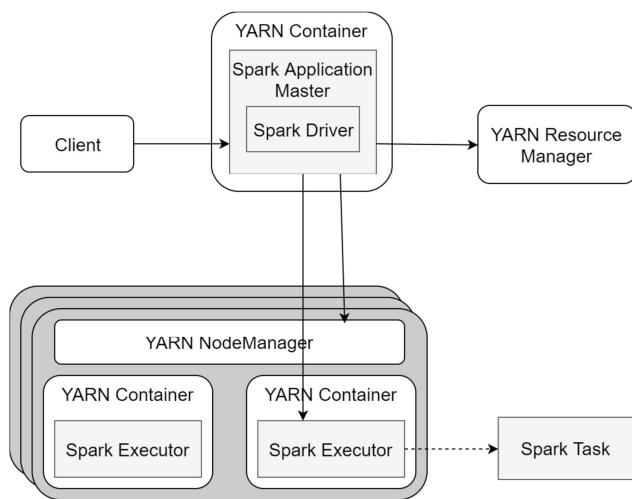


Figure 2.17: Spark on YARN Cluster Mode.

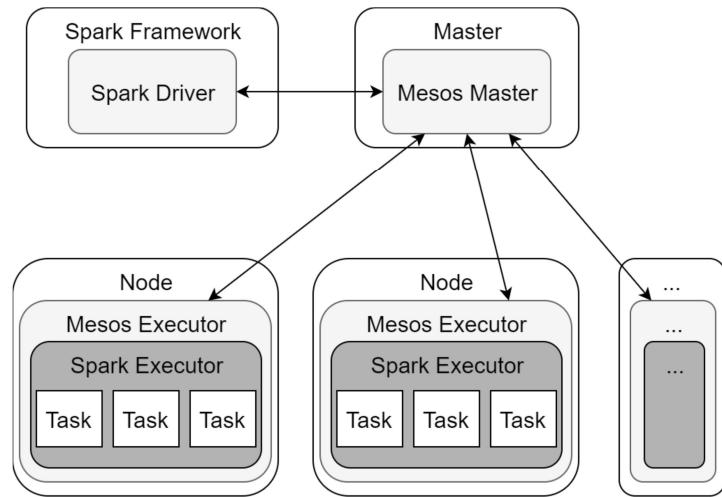


Figure 2.18: Spark on Mesos Coarse Grained Mode.

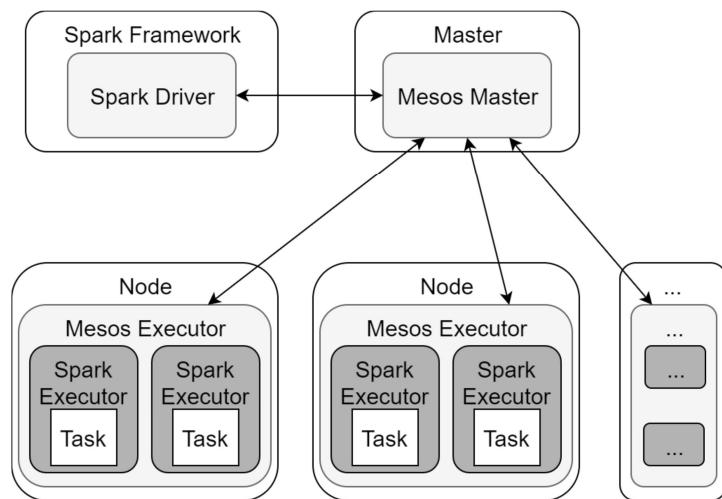


Figure 2.19: Spark on Mesos Fine Grained Mode.

respect to the other mode, and so it is good for interactive sessions. The drawback is that we are reserving Mesos resources for the complete duration of the application, unless dynamic allocation is active. Dynamic allocation allows to add and remove executors based on load: i) kill executor when they are idle, ii) add executors when tasks queue up in the scheduler. To use dynamic allocation it is required that the external shuffle service is running on each node. In fine-grained mode, shown in figure 2.19, Mesos tasks are launched for each Spark task, and those tasks die as soon as Spark tasks are done. This mode has too much overhead in case that Spark has too many tasks, for example if Spark application has 10,000 tasks, then Spark needs to be installed 10,000 times on Mesos agents. Because of this huge overhead, fine-grained mode has been deprecated since Spark version 2.0.0. This mode allows multiple instances of Spark to share cores at a very fine granularity, but it comes with an additional overhead in launching each task. Thus this mode is inappropriate for low-latency requirements like interactive queries or serving web requests, instead it is fine for batch and relatively static streaming. Similarly to what happens on YARN, it is possible to run spark in Mesos-client or Mesos-cluster mode. In client mode, the driver process is executed in the client machine that submits the job, so it is required that it stays connected to the cluster for the entire time of the application execution. In cluster mode instead, the driver program is run on a machine of the cluster.

2.5 VIRTUALIZATION AND CONTAINERIZATION

Virtualization refers to creating the virtual version of something, included hardware components, storage devices and computer networks. Virtualization is born in 1960's, as a way to logically partition the system resources offered by a mainframe computer between many different applications. From this point, the meaning of the word has been widely extended. Virtualization is a technology that allows creating multiple simulated environments or dedicated resources from a single unique physical hardware system. An hypervisor is a software that can directly connect to the hardware, with the purpose of splitting the unique physical system into separated environments, different from each other and secure, known as virtual machines (VM's). These VM's rely on the hypervisor ability to separate the hardware resources and distribute them in a proper way. The original physical machine equipped with the hypervisor is called host, meanwhile the VM's are called guests. These guests use the computation resources, such as CPU, memory and storage, as a set of resources that are easily re-allocatable. The operators can control the virtual instances of CPU, memory, storage and other resources, so that the guests can receive all the resources they need to execute their task. The words host and guest are used to distinguish the software that runs on the physical

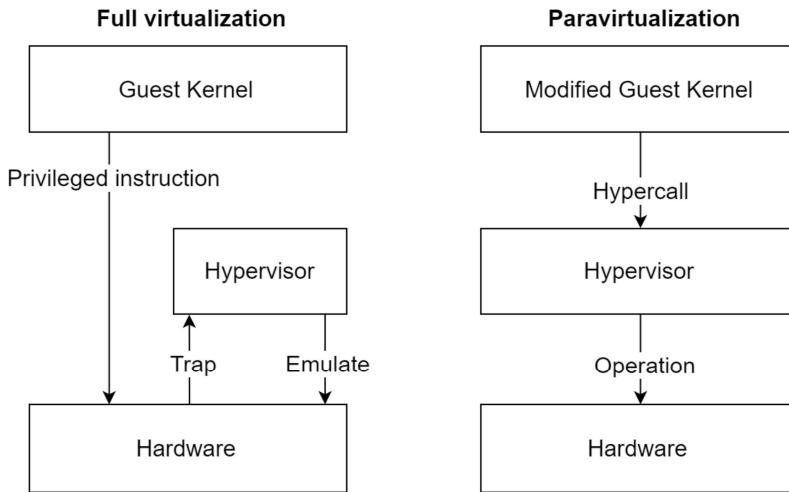


Figure 2.20: Full virtualization and paravirtualization.

hardware from the software that is running on the virtual machines. Hardware virtualization or platform virtualization refer to the creation of a VM that acts like a real computer with an OS. The software that is run in this VM is separated from the underlying hardware. This allows us to run particular configurations, for example we run a computer with a Windows OS that hosts a VM with Linux as guest OS. There are at least two different hardware virtualization types:

- full virtualization: it completely simulates the hardware in order to allow the software, typically a guest OS, to be run without the need of modifications
- paravirtualization: the hardware environment is not simulated, anyway guest programs are run in isolated domains, as if they were run in completely separated systems. Guest programs need to be modified in a specified way in order to run in this kind of environment.

In figure 2.20, we can see the differences between the two kind of virtualization. In paravirtualization, the VM presents a different interface compared to the one of a physical machine. This requires modification in the guest OS in order to allow its execution inside the VM. The hypervisor exposes a set of APIs that the guest OS must use to execute privileged instructions. Calls to these particular functions are often defined as Hypercalls. In full virtualization instead, VM have the same interface as the physical ones. Ideally, the guest OS would not be able to determine if it is being run on a physical or virtual machine. The great advantage of full virtualization is that we do not need to

modify the OS. In this way the hypervisor can adopt a trap system to execute privileged instructions. We can improve the efficiency of the virtualization by using hardware assisted virtualization, in particular we can decide to use CPU's that provide efficient support for virtualizing on hardware, but also other kind of hardware components that can improve the performance of the guest environments. Hardware virtualization can be seen as a trend of the enterprise IT that includes autonomic computing, a scenario in which the different environments are able to manage themselves based on the detected level of activity, and utility computing, where the processing power is seen as a utility that users pay only when needed. The purpose of virtualization is to centralize the administrative tasks, offering scalability and good resource utilization. With virtualization, different OS can run in parallel on a single CPU. This parallelism reduces overhead cost in a way different from multitasking, where different programs are executed in parallel on the same OS. Thanks to virtualization, an enterprise IT can better handle updates and rapid changes in OS and applications, with little impact on users. Virtualization allows organizations to have better efficiency and availability of resources and applications. It is important to remember that hardware emulation is a complete different thing from hardware virtualization, in particular with emulation we have a piece of hardware that imitates another piece of hardware. In virtualization instead a hypervisor, which is a piece of software, mimics a piece of hardware or even an entire computer. Moreover, a hypervisor is not an emulator, even though both are software programs that mimic hardware, their domain of use is different. Containerization is a OS-level virtualization technique that allows deploying and executing distributed applications without the need of launching an entire VM for each of the applications (figure 2.21). These multiple isolated systems are called containers. They are executed on top of a single host controller and access a single kernel. Since containers share the same OS kernel of the host, they can be a lot more efficient than a VM, that instead needs a separate instance of the OS. Containers own all the different components that are needed in order to execute the desired software, such as files, environment variables and libraries. The host OS controls the access of the container to the physical resources, such as CPU and memory, in order to prevent a single container from occupying the entire resources offered by the host. The main advantages of containerization come from efficiency in terms of memory, CPU and storage, when compared to traditional hardware virtualization. Since containers do not have the same overhead of VM, in particular we do not need different instances of the OS and is possible to support more containers on the same infrastructure. Containerization offers better performances since there is a single OS that takes care of all the hardware calls. A special advantage of the containers is the fact that they can be created much

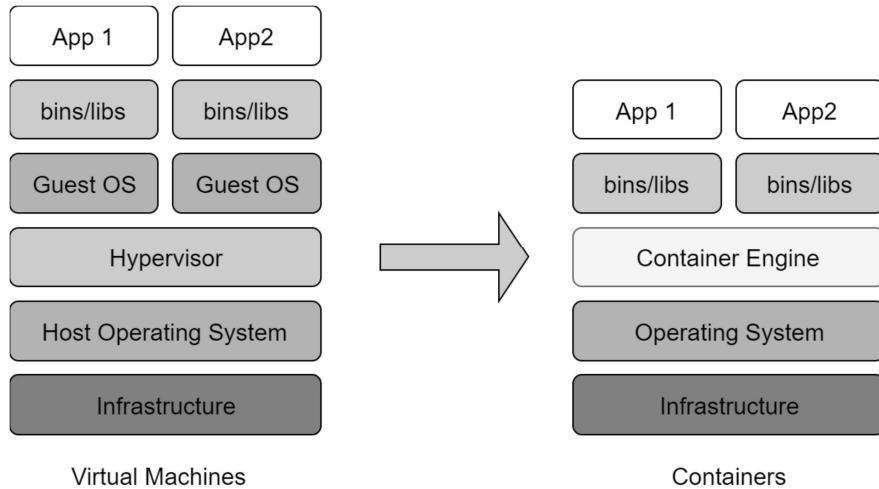


Figure 2.21: The difference in architecture between virtual machines and containers.

faster than instances based on an hypervisor. This allows to have a more agile environment and allows the creation of new approaches to virtualization, such as microservices and continuous integration and delivery of services. Potential disadvantages of containerization are the absence of isolation from the host OS. Since containers share the same host OS, a potential security threat might easily gain access to the entire system. This does not happen when using virtualization based on a hypervisor, since in this case the only compromised component would be the VM. A way to circumvent this problem, is the creation of containers inside an OS run from a VM. This prevents the security breach at container level from letting the attacker gain access to the OS of the physical host. Another little disadvantage of containerization is that containers must execute the same OS as the base OS, meanwhile instances based on an hypervisor are allowed to execute different OSs. Because of this, a container that is running on a Linux host, can neither execute an instance of Windows OS nor a Windows designed application. Containerization has gained more and more relevance thanks to the wide spread of the open source software Docker, that has developed enhanced portability to the containers, allowing them to be moved from different systems that share the same kind of host OS without the need to change even a single line of code. In particular, with Docker container there are no environment variables that must be set on the guest OS or library dependencies that need to be managed.

2.6 SYMBOLIC EXECUTION

In this section we present an overview of *symbolic execution*, taken primarily from a research work of Baldoni et al [14].

In computer science, the term symbolic execution refers to a software program analysis technique used to determine which data inputs cause the execution of each part of a program. It was introduced in the mid '70s [52, 20, 53, 43] mainly to test whether a software program could violate certain properties, e.g. that no divisions by zero are performed, no null pointers get dereferenced, no access to protected data can happen by unauthorized users, etc. In general, it's not possible to decide every possible program property by means of automated checks, for example we cannot predict the target of an indirect jump. In practice, approximate and heuristic-based analyses are used in many cases, even in the field of mission-critical and security applications.

Software testing is performed to check that certain program properties hold for any possible usage scenario. A viable approach would be to test the program using a wide range of different, possibly random inputs. As the problem may occur only for very specific input values, we need to automate the exploration of the domain of the possible input data.

With symbolic execution many possible execution paths are explored in parallel, without necessarily requiring concrete inputs. The idea is to replace the fully specified input data with symbols, that are their abstract representation, devolving to constraint solvers the construction of actual instances that would cause property violations.

The symbolic execution interpreter walks through all the steps of the program, associating symbolic values to inputs rather than obtaining their actual values, building expressions in terms of those symbols and program variables, and constraints in terms of the symbols corresponding to the possible outcomes of each conditional branch.

When a program is run with a specific set of input data (a concrete execution), a single control flow path is explored. Hence, concrete executions can only under-approximate the analysis of the property of interest. With symbolic execution, multiple paths that a program could take under different inputs can be simultaneously explored. This means that a sound analyses can be done, giving stronger guarantees about the checked property [14]. When a program runs with *symbolic* – rather than concrete – input values, the execution is performed by a *symbolic execution engine*, which builds and updates a structure to hold (i) a first-order Boolean *formula* describing the conditions satisfied by all the traversed branches along that path, and (ii) a *symbolic memory store* mapping variables to symbolic expressions or values, for each path traversed by the control flow. Execution of a branch updates the formula, while assignments update the symbolic store. Finally, a *model checker*, commonly based on a *satisfiability modulo theories* (SMT)

```

1. void foo(int a, int b) {
2.     int x = 1, y = 0;
3.     if (a != 0) {
4.         y = 3+x;
5.         if (b == 0)
6.             x = 2*(a+b);
7.     }
8.     assert(x-y != 0);
9. }
```

Figure 2.22: Example: which values of *a* and *b* make the `assert` fail?

solver [17], is used to verify if the property is violated somewhere along each explored path and if the path itself is concretely feasible, i.e., if any assignment of concrete values to the program’s symbolic arguments exists that satisfies its formula [14].

Symbolic execution techniques have been emphasized to a wide audience following the DARPA announcement in 2013 at the Cyber Grand Challenge, a two-year competition pursuing the creation of automatic systems for vulnerability detection, exploitation, and patching in near real-time [76]. More important, symbolic execution based engines have been running 24/7 in the testing process of many Microsoft applications since 2008, discovering nearly 30% of all the flaws discovered by file fuzz testing during the development of Windows 7, which other program analyses and blackbox testing techniques missed [38].

EXAMPLE With reference to the C code in Figure 2.22, let’s say we want to discover which of the 2^{32} possible 4-byte inputs make the `assert` at line 8 of function `foo` fail. If we address the problem by running concretely the function `foo` on randomly generated inputs, we will unlikely pick up exactly the assert-failing inputs. Symbolic execution go beyond this limitation by reasoning on *classes of inputs*, rather than single input values, thanks to the evaluation of the code using symbols for its inputs, instead of concrete values,

Going further into details, a symbol α_i is associated to each value that cannot be resolved by a static analysis of the code, e.g. an actual parameter of a function or data read from an input stream. At any time, the symbolic execution engine maintains a state $(\text{stmt}, \sigma, \pi)$ where:

- `stmt` is the next statement to evaluate. At this time, we assume that `stmt` can be an assignment, a conditional branch, or a jump.
- σ is a *symbolic store* that associates program variables with concrete values of expressions or symbolic values α_i .
- π identifies the *path constraints*, i.e., it is a formula expressing a set of assumptions on the symbols α_i as a result of branches

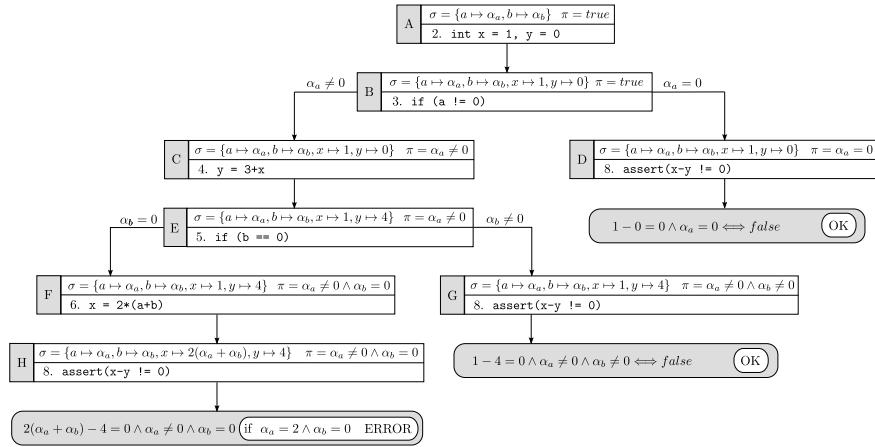


Figure 2.23: Symbolic execution tree of function foo given in Figure 2.22.

Each execution state, labeled with an upper case letter, shows the statement to be executed, the symbolic store σ , and the path constraints π . Leaves are evaluated against the condition in the assert statement. Image courtesy of Association for Computing Machinery [14].

taken in the execution to reach stmt. At the start of the analysis, $\pi = \text{true}$.

Depending on stmt, the symbolic engine changes the state as follows:

- The evaluation of an assignment $x = e$ updates the symbolic store σ by associating x with a new symbolic expression e_s . We express this association with $x \mapsto e_s$, where e_s is obtained by evaluating e in the context of the current execution state and can be any expression involving unary or binary operators over symbols and concrete values.
- The evaluation of a conditional branch `if e then strue else sfalse` affects the path constraints π . The symbolic execution is forked by creating two execution states with path constraints π_{true} and π_{false} , respectively, which correspond to the two branches: $\pi_{\text{true}} = \pi \wedge e_s$ and $\pi_{\text{false}} = \pi \wedge \neg e_s$, where e_s is a symbolic expression obtained by evaluating e . Symbolic execution independently proceeds on both states.
- The evaluation of a jump goto s updates the execution state by advancing the symbolic execution to statement s.

Figure 2.23 shows a symbolic execution of function foo, that can be adequately represented as a tree. In the initial state (execution state A) the path constraints are true and input arguments a and b are associated with symbolic values. After local variables initialization x and y at line 2, the symbolic store is updated by associating x and y with concrete values 1 and 0, respectively (execution state B). A conditional branch is met in line 3 and the execution is forked: according to

which branch is taken, a different statement is evaluated and different assumptions are made on symbol α_a (execution states C and D). In the branch corresponding to $\alpha_a \neq 0$, variable y is assigned to $x + 3$, obtaining $y \mapsto 4$ in state E because $x \mapsto 1$ in state C. Arithmetic expression evaluation, generally, change only the symbolic values. When the assert at line 8 is reached by fully expanding every execution state on all branches, we can check which input values for parameters a and b can make the assert fail. By analyzing execution states {D, G, H}, we can conclude that only H can make $x-y = 0$ true. The path constraints for H at this point implicitly define the set of inputs that are unsafe for foo. In particular, any input values such that:

$$2(\alpha_a + \alpha_b) - 4 = 0 \wedge \alpha_a \neq 0 \wedge \alpha_b = 0$$

will make assert fail. An instance of unsafe input parameters can be eventually determined by invoking an *SMT solver* [17] to solve the path constraints, which in this example would yield $a = 2$ and $b = 0$.

2.6.0.1 Challenges of Symbolic Execution

The example shown in Section 2.6 represents a case where symbolic execution can derive *all* the inputs that make the assert fail, by exploring all the possible execution states. With regards to the underpinning theory, exhaustive symbolic execution represents a *sound* and *complete* methodology for any decidable analysis [14]. Soundness (no false negatives) means that all possible unsafe inputs are guaranteed to be found, while completeness (no false positives) means that input values deemed unsafe are actually unsafe. Exhaustive symbolic execution cannot be easily scalable beyond small-sized applications. In many practical cases, a trade-off between soundness and performance approach is used.

Symbolic execution applied to real world applications faces many challenges, deriving from the increased complexity of those applications compared to the simplicity of our example. These challenges are pointed out by the following questions and observations:

- *Pointers, arrays, or other complex objects*: how are they handled by the symbolic engine? Code manipulating pointers and data structures may identify memory addresses that are described by symbolic expressions, in addition to symbolic stored data.
- *Interactions across the software stack*: how are they handled by the engine? System code/library calls can have side effects, e.g. file creations, calls back to user code or callback functions that could affect the execution at a later stage of the computation and must be taken in account. Evaluating all possible interaction outcomes may be unfeasible.

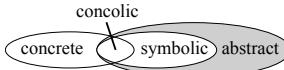


Figure 2.24: Concrete and abstract execution machine models. Image courtesy of Association for Computing Machinery [14].

- *State space explosion*: how does symbolic execution deal with path explosion? The presence of iterative constructs in the language (iteration loops) can dramatically increase the number of execution states, up to the point where it becomes unfeasible for a symbolic execution engine to exhaustively explore all the possible states within a reasonable amount of time.
- *Constraint solving*: in practice, what can a constraint solver do? SMT solvers can solve up to complex combinations of constraints over hundreds of variables. However, constructs such as non-linear arithmetic pose a major obstacle to efficiency.
- *Binary code*: what issues can arise when symbolically executing binary code? While our example of Section 2.6 is written in C, in many cases programs are only available as binary code. The availability of the source code of an application usually simplifies symbolic execution, as it can exploit high-level properties (e.g., object morphology) that can be inferred statically by analyzing the source code.

The above questions and assumptions are addressed in different ways and lead to different choices according to the specific contexts where symbolic execution is used. Typically these choices affect soundness or completeness, however this does not represent a serious limitation in many real cases where a partial exploration of the domain of possible execution states may provide enough information to reach the goal (e.g., determining which input values provoke an application failure) within specified time limits.

2.6.1 Symbolic Execution Engines

This section presents some principles for the design of symbolic executors and introduces the concept of *concolic* execution.

MIXING SYMBOLIC AND CONCRETE EXECUTION As shown in the previous example, modeling all runs of concrete executions of real-world software programs is desirable but very often impossible in practice.

Symbolic execution, as we have described it so far, cannot explore feasible executions that would result in path constraints that cannot be dealt with [25]. Loss of soundness originates from many sources, e.g. untraceable external code, complex constraints involving non-linear

arithmetic or transcendental functions. As constraint solving is a major performance barrier for an engine, we can assess the solvability in an absolute sense but also in terms of efficiency. Considering that practical programs are typically not self-contained, it is very challenging to statically analyze the whole software stack, especially considering the evaluation of every possible side effects.

One way to overcome this issue in practice is to mix concrete and symbolic execution, aka *concolic execution*, where the term concolic is a portmanteau of the words “concrete” and “symbolic” (Figure 2.24). Some applications of this general principle are briefly discussed in the remainder of this section.

DYNAMIC SYMBOLIC EXECUTION A common approach to concolic execution, known as *dynamic symbolic execution* (DSE) or *dynamic test generation* [36], is to have concrete execution drive the symbolic one. A concrete store σ_c is maintained, in addition to the symbolic and the path constraints stores, by the execution engine. It chooses an initial arbitrary input and executes the program both concretely and symbolically, by simultaneously updating the two stores and the path constraints. The symbolic execution walks through the same branch taken by the concrete execution, and the constraints extracted from the branch condition are added to the current set of path constraints. The symbolic engine does not need to invoke the constraint solver to decide if a branch condition is satisfiable, as this is already tested by the concrete execution. Since not all the paths might be explored by concrete execution, the path conditions associated to a branch can be negated and passed to the SMT solver to determine a satisfying assignment for the new constraints, i.e., to generate a new input. This strategy can be repeated until the desired path coverage is reached.

Example (extracted from [14] by R.Baldoni et al.). Consider the C function in Figure 2.22 and suppose to choose $a = 1$ and $b = 1$ as input parameters. Under these conditions, the concrete execution takes path $A \rightsquigarrow B \rightsquigarrow C \rightsquigarrow E \rightsquigarrow G$ in the symbolic tree of Figure 2.23. Besides the symbolic stores shown in Figure 2.23, the concrete stores maintained in the traversed states are the following:

- $\sigma_c = \{a \mapsto 1, b \mapsto 1\}$ in state A ;
- $\sigma_c = \{a \mapsto 1, b \mapsto 1, x \mapsto 1, y \mapsto 0\}$ in states B and C ;
- $\sigma_c = \{a \mapsto 1, b \mapsto 1, x \mapsto 1, y \mapsto 4\}$ in states E and G .

After checking that the assert conditions at line 8 succeed, we can generate a new control flow path by negating the last path constraint, i.e., $\alpha_b \neq 0$. The solver at this point would generate a new input that satisfies the constraints $\alpha_a \neq 0 \wedge \alpha_b = 0$ (for instance $a = 1$ and $b = 0$) and the execution would continue in a similar way along the path $A \rightsquigarrow B \rightsquigarrow C \rightsquigarrow E \rightsquigarrow F$.

DSE uses concrete inputs to drive the symbolic execution toward a specific path, but needs to select a branch to negate whenever a new path has to be explored. Also, each concrete execution may add new branches that will have to be visited. Since the set of unexplored branches across all performed concrete executions can be very large, it's fundamental to adopt effective search heuristics. For example, DART [36] chooses the next branch to negate using a depth-first strategy. Additional strategies for picking the next branch to negate have been presented in literature, like the *generational search* of SAGE [37] that systematically yet partially explores the state space, maximizing the number of new tests generated while also avoiding redundancies in the search. This is achieved by negating constraints following a specific order and by limiting the backtracking of the search algorithm. Since the state space is only partially explored, the initial input plays a crucial role in the effectiveness of the overall approach. The importance of the first input is similar to what happens in traditional *black-box fuzzing*; hence, symbolic engines such as SAGE are often referred to as *white-box fuzzers*.

The symbolic information maintained during a concrete run can be exploited by the engine to obtain new inputs and explore new paths. The next example shows how DSE can handle invocations to external code that is not symbolically tracked by the concolic engine.

Example (extracted from [14] by R.Baldoni et al.). Consider function `foo` in Figure 2.25a and suppose that `bar` is not symbolically tracked by the concolic engine (e.g., it could be provided by a third-party component, written in a different language, or analyzed following a black-box approach). Assuming that $x = 1$ and $y = 2$ are randomly chosen as the initial input parameters, the concolic engine executes `bar` (which returns $a = 0$) and skips the branch that would trigger the error statement. At the same time, the symbolic execution tracks the path constraint $\alpha_y \geq 0$ inside function `foo`. Notice that branch conditions in function `bar` are not known to the engine. To explore the alternative path, the engine negates the path constraint of the branch in `foo`, generating inputs, such as $x = 1$ and $y = -4$, that actually drive the concrete execution to the alternative path. With this approach, the engine can explore both paths in `foo` even if `bar` is not symbolically tracked.

A variant of the previous code is shown in Figure 2.25b, where function `qux` – differently from `foo` – takes a single input parameter but checks the result of `bar` in the branch condition. Although the engine can track the path constraint in the branch condition tested inside `qux`, there is no guarantee that an input able to drive the execution toward the alternative path is generated: the relationship between `a` and `x` is not known to the concolic engine, as `bar` is not symbolically tracked. In this case, the engine could re-run the code

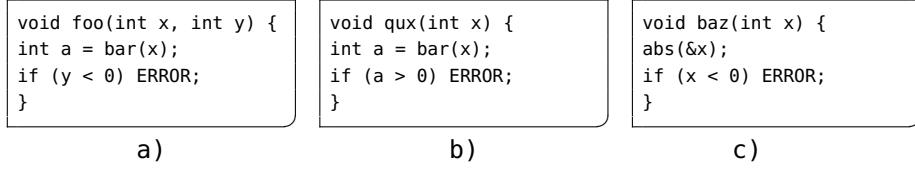


Figure 2.25: Concolic execution: (a) testing of function `foo` even when `bar` cannot be symbolically tracked by an engine, (b) example of false negative, and (c) example of a path divergence, where `abs` drops the sign of the integer at `&x`. Image courtesy of Association for Computing Machinery [14].

using a different random input, but in the end it could fail to explore one interesting path in `qux`.

A related issue is presented by Figure 2.25c. We observe a *path divergence* when inputs generated for a predicted path lead execution to a different path. In general, this can be due to symbol propagation not being tracked, resulting in inaccurate path constraints, or to imprecision in modeling certain (e.g., bitwise, floating-point) operations in the engine. In the example, function `baz` invokes the external function `abs`, which performs a side effect on `x` by assigning it with its absolute value. Choosing $x = 1$ as the initial concrete value, the concrete execution does not trigger the error statement, but the concolic engine tracks the path constraint $\alpha_x \geq 0$ due to the branch in `baz`, trying to generate a new input by negating it. However the new input, e.g., $x = -1$, does not trigger the error statement due to the (untracked) side effects of `abs`. Interestingly, the engine has no way of detecting that no input can actually trigger the error.

Two notable downsides of dynamic symbolic execution are false negatives (i.e., missed paths) and path divergences as depicted by the example. DSE trades soundness for performance and implementation effort: false negatives can happen, because some program executions – and therefore possible erroneous behaviors – may be missed, leading to a *complete*, but *under-approximate* form of program analysis [14].

SELECTIVE SYMBOLIC EXECUTION S²E [28] takes a different approach by mixing symbolic and concrete execution considering a partial exploration (i.e. limited to some of the components) of a software stack, dropping the others. *Selective symbolic execution* accurately interleaves concrete and symbolic execution, while maintaining the meaningfulness of the overall exploration.

Suppose a function A calls a function B and the execution mode changes at the call site. Two scenarios arise: (1) *From concrete to symbolic and back*: the arguments of B are made symbolic and B is explored symbolically in full. B is also executed concretely and its concrete result is returned to A. After that, A resumes concretely. (2) *From symbolic to concrete and back*: the arguments of B are concretized, B is executed

concretely, and execution resumes symbolically in A. This may impact both soundness and completeness of the analysis: (i) *Completeness*: to make sure that symbolic execution skips any paths that would not be realizable due to the performed concretization (possibly leading to false positives), S²E collects path constraints that keep track of how arguments are concretized, what side effects are made by B, and what return value it produces. (ii) *Soundness*: concretization may cause missed branches after A is resumed (possibly leading to false negatives). To remedy this, the collected constraints are marked as *soft*: whenever a branch after returning to A is made inoperative by a soft constraint, the execution backtracks and a different choice of arguments for B is attempted. To guide re-concretization of B's arguments, S²E also collects the branch conditions during the concrete execution of B, and chooses the concrete values so that they enable a different concrete execution path in B.

2.6.2 Path Selection

Since enumerating all paths of a program can be prohibitively expensive, in many software engineering activities related to testing and debugging the search is prioritized by looking at the most promising paths first. Among several strategies for selecting the next path to be explored, we now briefly overview some of the most effective ones. We remark that path selection heuristics are often tailored to help the symbolic engine achieve specific goals (e.g., overflow detection). Finding a universally optimal strategy remains an open problem.

Depth-first search (DFS), which expands a path as much as possible before backtracking to the deepest unexplored branch, and *breadth-first search* (BFS), which expands all paths in parallel, are the most common strategies. DFS is often adopted when memory usage is at a premium, but is hampered by paths containing loops and recursive calls. Hence, in spite of the higher memory pressure and of the long time required to complete the exploration of specific paths, some tools resort to BFS, which allows the engine to quickly explore diverse paths detecting interesting behaviors early. Another popular strategy is *random path selection*, that has been refined in several variants. For instance, KLEE [24] assigns probabilities to paths based on their length and on the branch arity: it favors paths that have been explored fewer times, preventing starvation caused by loops and other path explosion factors.

Fitness functions have been largely used in the context of search-based test generation [59]. A fitness function measures how close an explored path is to achieve the target test coverage. Several works, e.g., [91, 25], have applied this idea in the context of symbolic execution. As an example, [91] introduces *fitnex*, a strategy for flipping branches in concolic execution that prioritizes paths likely *closer* to

take a specific branch. In more detail, given a target branch with an associated condition of the form $|a - c| == 0$, the closeness of a path is computed as $|a - c|$ by leveraging the concrete values of variables a and c in that path. Similar fitness values can be computed for other kinds of branch conditions. The path with the lowest fitness value for a branch is selected by the symbolic engine. Paths that have not reached the branch yet get the worst-case fitness value.

2.6.3 Symbolic Backward Execution

Symbolic backward execution (SBE) [26, 29] is a variant of symbolic execution in which the exploration proceeds from a target point to an entry point of a program. The analysis is thus performed in the reverse direction than in canonical (forward) symbolic execution. The main purpose of this approach is typically to identify a test input instance that can trigger the execution of a specific line of code (e.g., an assert or throw statement). This can be very useful for a developer when performing debugging or regression testing over a program. As the exploration starts from the target, path constraints are collected along the branches met during the traversal. Multiple paths can be explored at a time by an SBE engine and, akin to forward symbolic execution, paths are periodically checked for feasibility. When a path condition is proved unsatisfiable, the engine discards the path and backtracks.

[55] discusses a variant of SBE dubbed *call-chain backward symbolic execution* (CCBSE). The technique starts by determining a valid path in the function where the target line is located. When a path is found, the engine moves to one of the callers of the function that contains the target point and tries to reconstruct a valid path from the entry point of the caller to the target point. The process is recursively repeated until a valid path from the main function of the program has been reconstructed. The main difference with respect to the traditional SBE is that, although CCBSE follows the call-chain backwards from the target point, inside each function the exploration is done as in traditional symbolic execution.

A crucial requirement for the reversed exploration in SBE, as well as in CCBSE, is the availability of the inter-procedural control flow graph which provides a whole-program control flow and makes it possible to determine the call sites for the functions that are involved in the exploration. Unfortunately, constructing such a graph can be quite challenging in practice. Moreover, a function may have many possible call sites, making the exploration performed by a SBE still very expensive. On the other hand, some practical advantages can arise when the constraints are collected in the reverse direction.

THIS chapter illustrates the previous works done on xSpark, that is the base for the developments made and contribution given by this thesis. xSpark is a Spark extension that offers optimized and elastic provisioning of resources in order to meet execution deadlines. This is obtained by using nested control loops. A centralized control loop, implemented on the master node, controls the execution of the different stages of an application; at the same time multiple local loops, one per executor, focus on task execution. In Figure 3.1 we can see a high level representation of xSpark execution flow. A preliminary Profiling Phase, that is executed once per application, is obtained by executing the application once to obtain information about its runtime execution. The log generated by the application are used to generate its Profiling Data. In the Execution Phase, we control the application by means of xSpark's control loops. The centralized control loop is represented as a Heuristic Based Planner, which exploits the Profiling Data. The profiling data is used to understand the amount of work that is needed to execute the application. This component uses the provided profiling data to determine the amount of resources to assign to executors for each of the application stages, in order to complete the execution within the given deadline. The local loops instead are represented as Control Theory Controllers. They perform fine-grained tuning of the resources assigned to each of the executors using a control theory based controller. This component is used to counteract the possible imprecision of the estimated needed resources, which may be caused by different factors, such as the available memory, etc. Usually Spark applications are run multiple times, as they are reusable and longlasting assets. xSpark exploits an initial profiling phase to create an enriched DAG describing the entire application execution flow, by collecting information about all the stages of the application. For each stage, xSpark annotates the DAG with the execution time (stage duration), the number of task processed, the number of input records read, the number of output record written and the nominal rate, defined as the number of records that a single core processes during one second of execution.

In Figure 3.2 we can see a fragment, relative to stage number zero, of the profiling data of a Louvain application executed in xSpark. The duration field contains the serialized total duration of the tasks in milliseconds. At runtime, the annotated DAG allows the xSpark scheduler function to know how much work was already completed and how much remains to be done. This means that xSpark can only

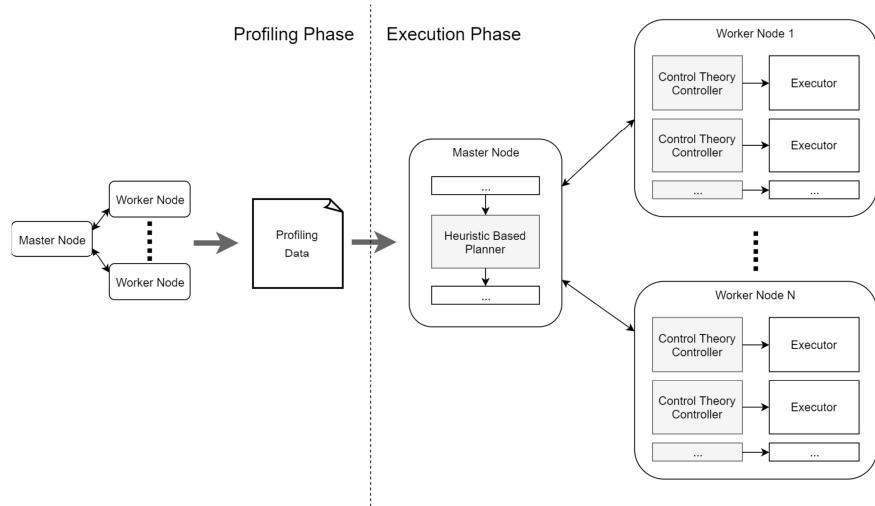


Figure 3.1: High level setup and execution flow of xSpark.

```
{
  "0": {
    "actual_records_read": 20000000.0,
    "actual_records_write": 20000000.0,
    "actualtotalduration": 186000.0,
    "bytesread": 337777780.0,
    "byteswrite": 0.0,
    "cachedRDDs": [
      5
    ],
    "duration": 11000.0,
    "genstage": false,
    "id": 0.0,
    "io_factor": 1.0,
    "monocoreduration": 381388.0,
    "monocoretotalduration": 10127715.0,
    "name": "mapPartitions at VertexRDD.scala:356",
    "nominalrate": 52440.03482018312,
    "nominalrate_bytes": 885653.9272342077,
    "numtask": 500,
    "parentsIds": [],
    "recordsread": 20000000.0,
    "recordswrite": 0.0,
    "shufflebytesread": 0.0,
    "shufflebyteswrite": 114175342.0,
    "shufflerecordsread": 0.0,
    "shufflerecordswrite": 20000000.0,
    "skipped": false,
    "t_record_ta_executor": 0.0190694,
    "totalduration": 207000.0,
    "weight": 28.777443181222274
  },
}
```

Figure 3.2: Example of profiling data from a Louvain application.

optimize the allocation of the resources if all the executions of the same application use the same DAG. This might not always be the case, for example when the code contains branches or loops, because these might need to be resolved in different ways at runtime. If the application DAG at runtime differs from the one obtained during the profiling phase, xSpark is not able to estimate the amount of remaining work. The centralized control loop is activated before the execution of every stage, it uses an heuristic, explained in Section 3.2, in order to assign a deadline to the stage, calculate the amount of CPU cores that are needed to satisfy it, and assign cores to the allocated executors. The per-stage deadline takes into account the amount of work already completed, the consumed time, and the overall deadline. All the computations done by the heuristic are based on the information stored in the DAG and obtained during the profiling phase. Unfortunately many factors can influence the actual performance and invalidate the prediction, such as the amount of records that have been filtered-out, the available memory, the number of used nodes, the storage layer dimension, and so on. It is important to remember that Spark mostly uses in-memory data, but there are operations like `textFile`, `saveAsTextFile` and `saveAsSequenceFile` that impose restricting constraints on the storage layer. If not correctly dimensioned, the storage layer might become a bottleneck, causing the throughput to degrade and thus making the provisioning predicted by xSpark incorrect. Local control loops, explained in Section 3.3, counteract this imprecision by dynamically modifying the amount of CPU cores assigned to the executors during the execution of a stage. This can lead the executor to use more or less resources than the ones previously assigned by the centralized control loop. The local loop controls the progress of a specific executor with respect to the tasks it has assigned. A control theory algorithm determines the amount of CPU cores that must be allocated to the executor for the next control period, typically one second, and assigns them. xSpark uses Docker in order to dynamically allocate CPU cores and memory, as explained in Section 2.1.1. Memory allocation simply sets an upper bound to the memory that each docker container (executor) can use. CPU cores instead are allocated in a more sophisticated way. Using Linux cgroups, Docker can support CPU shares, reservations and quotas. In particular, xSpark uses CPU quotas. CPU shares are not able to limit the number of CPU cores used by a container in a deterministic way, in particular it is not independent of other processes running on the same machine. CPU reservation instead does not have the fine granularity we are looking for, indeed the allocation would be limited to entire cores. By using CPU quotas instead, we have a reliable and tunable mechanism that provides also fine granularity allocation, in particular it allows xSpark to allocate fractions of cores to the containers (executors), with a precision up to 0.05 cores.

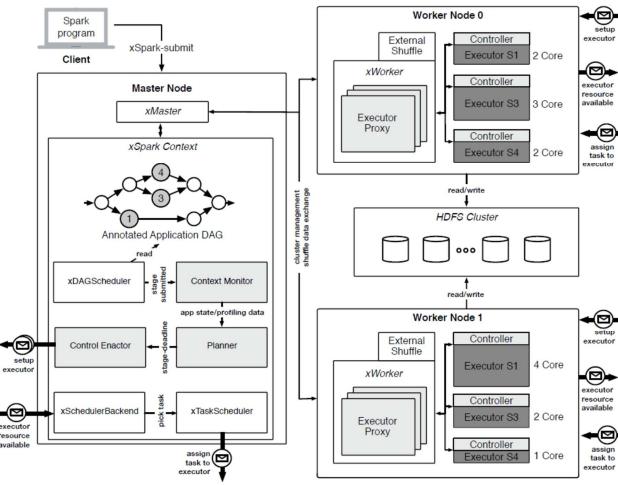


Figure 3.3: Architecture of xSpark. New components are represented in light grey boxes, meanwhile those that start with an x are the modified ones. In dark grey are represented the containerized components (the executors).

3.1 ARCHITECTURE

To achieve the objectives of xSpark, the architecture and processing model of Spark have been modified. In Figure 2.2 we can see how xSpark architecture differs from Spark. The principal architectural change introduced by xSpark is an increased focus on stages. Instead of considering entire applications, xSpark reasons on per-stage deadlines. xSpark instantiates an executor per stage per worker, instead of a single executor per worker for all the stages that will be executed. This way the resources that are allocated to a single executor only impact the performance of the stages that are associated with it, this leads to a fine grained control over the different stage, and thus on the entire application. When multiple stages are run in parallel, multiple executors can be running on the same worker node. When a stage is submitted for execution, one executor per worker is created and bound to that stage. This way computation and data are equally spread across the entire cluster. Thanks to containers, xSpark can isolate the execution of the different executors that are running on the same worker node and achieve quick, fine-grained resource provisioning. On average, containers can be modified in less than one second, allowing a more precise way of allocating cores. Users submit applications and their deadlines to the master node using the submit command. This creates an *xSparkContext* on the master node, containing an *xMaster* object that is used to manage the cluster and that knows all the

resources that are available on each worker node. A *xSparkContext* is composed by six components:

- *xDagScheduler* schedules the stages according to the application's DAG. The submitted stage is enriched with the information obtained during the profiling phase
- *ContextMonitor* monitors the progress of the application, taking into account stage scheduling and completion. It also stores information about the performance of the execution, that will be later used to calculate the deadlines and the resources needed by upcoming stages
- *Planner* is heuristics based and is used to calculate the deadlines and resources associated with a stage
- *ControlEnactor* determines when a stage is ready to be executed, meaning there are enough resources (cores) and sufficient executors become available. It also has the duty of initializing the different executors.
- *xSchedulerBackend* controls the stage execution. In particular it launches new tasks by taking into account the resources availability and registers their completion
- *xTaskScheduler* also controls stage execution, with the goal of allocating tasks to the available cores to optimize data locality. In general, the closer a data partition required by a task is to the task's executor, the better.

xSpark also modified the worker nodes. Each node contains a *xWorker* that connects to a *xMaster*, generates local controllers for the executors, and controls the evolution of its executor by dynamically scaling their resources. *xWorker* creates an *Executor Proxy* for each of the executors that are running on its node. These proxies are placed between the executors and the *xSchedulerBackend*, and are used to monitor the execution progress of the stage assigned to the executors. It is important to remember that each executor focuses on a single stage at a time. The heuristic calculates how many tasks must be executed by each of the executors of the same stage, the strategy is to have an equal number of tasks assigned to each of the executors. This way the deadline assigned to each of the executors coincides with the deadline of their stage. This allows the different executors of the same stage to not be synchronized. Native Spark instead requires that executors with free resources spontaneously require new tasks to execute from the master node.

Every *xWorker* uses an External Shuffle Service. Native Spark moves data across the cluster in different ways. If the data is stored on executor's memory, then the executor itself manages the data exchange.

If the data is stored on an external storage system (e.g., HDFS cluster), they can be retrieved by using different communication protocols. If the data is stored in the internal storage of a worker node, then the data is managed by the External Shuffle Service. Notice that this is not the default, but xSpark adopts this technique to be able to assign zero CPU cores to an executor, without loosing the ability to read data, since it is effectively not performed by the executor but by the external service.

3.2 HEURISTIC

xSpark uses a heuristic to compute per-stage deadlines and to estimate how many cores must be allocated for a stage to successfully fulfill the deadline. In order to do this, at submission time the user is asked to specify three parameters: i) the application deadline, ii) the cluster size, and iii) the number of cores per worker node. Before executing the application, xSpark performs a feasibility check given the available resources.

When a stage is submitted for execution, its deadline is computed

$$\text{deadline}(\text{sk}) = \frac{\alpha \cdot \text{ApplicationDeadline} - \text{SpentTime}}{\text{weight}(\text{sk})}$$

where SpentTime is the time already spent for execution and α a value between 0 and 1 that xSpark uses to be more conservative with respect to the provided ApplicationDeadline. The weight is computed

$$\begin{cases} w1(\text{sk}) = \#(\text{RemainingStages} + 1) \\ w2(\text{sk}) = \frac{\sum_{i=k}^{k+w1} \text{duration}(s_i)}{\text{duration}(\text{sk})} \\ \text{weight}(\text{sk}) = \beta \cdot w1(\text{sk}) + (1 - \beta) \cdot w2(\text{sk}) \end{cases}$$

where $w1$ is the number of stages still to be scheduled (s included) and $w2$ is the rate between the duration of s and the duration of the remaining stages (s included). xSpark then proceeds to estimate how many cores are needed to execute the stage:

$$\text{estimatedCores}(\text{sk}) = \lceil \frac{\text{inputRecords}(\text{sk})}{\text{deadline}(\text{sk}) \cdot \text{nominalRate}(\text{sk})} \rceil$$

where inputRecords is the number of records that will be processed by sk and nominalRate is the number of records processed by a single core per second in stage sk.

Since xSpark controls the resource allocation of a stage before and during the execution, the maximum amount of allocable cores needs to be greater than the estimated one, in order to be able to accelerate when progressing slower than expected

$$\text{maxAllocableCores}(\text{sk}) = \text{overscale} \cdot \text{estimatedCores}(\text{sk})$$

The final step is to determine the initial number of cores that should be assigned to the different executors, xSpark distributes the cores equally amongst the available workers by creating one executor per stage per worker. In this way, it is guaranteed that executor performances will be equal, and that xSpark can compute the same deadline for all the executors. The initial number of cores per executor is computed as

$$\text{initCorePerExec}(\text{sk}) = \lceil \frac{\text{maxAllocableCores}(\text{sk})}{\text{overscale} \cdot \text{cq} \cdot \text{numExecutors}} \rceil \cdot \text{cq}$$

where `numExecutors` is the number of executors and `cq` is the core quantum, a constant that defines the quantization applied to resource allocation, the smaller this value is, the more precise the allocation.

3.3 CONTROLLER

Each containerized executor has an associated local controller, whose goal is to fulfill the per-stage deadline taking into account external disturbances by dynamically allocating CPU cores. The controllers use control theory, with no heuristic involved.

The centralized control loop determines the desired stage duration, the maximum and the initial number of cores that should be assigned to the executors and the number of tasks that must be processed. Local controllers adjust the number of allocated cores, according to the work that has already been accomplished.

Executors that are dedicated to different stages are implicitly independent, and thus their controllers are also independent. The executors that are running in parallel on the same stage must complete the same amount of work (number of tasks) in the same desired time. This means that local controllers are independent and do not need to communicate amongst themselves. Moreover, the heuristic is relegated outside the local controller, so that it cannot compromise the controller's stability. In the controller, the progress set point is chosen based on the desired completion time. Its value is received from the centralized control loop. In Figure 3.4 we see the prescribed completion percentage, in particular t_{co} is the desired completion time and $\alpha \in (0, 1]$ is a configuration parameter used to determine how much earlier we are willing to complete the execution with respect to requested deadline. In order to track the set point ramp, we need to use a Proportional plus Integral (PI) controller. As a result, the discrete-time controller in state space form reads:

$$\begin{cases} x_C(k) = x_C(k-1) + (1 - \alpha)(a_{\%}^o(k-1) - a_{\%}(k-1)) \\ c(k) = Kx_C(k) + K(a_{\%}^o - a_{\%}(k)) \end{cases}$$

where ($a_{\%}^o$ is the prescribed progress percentage at each k control step and $a_{\%}$ is the accomplished completion percent at each k control step).

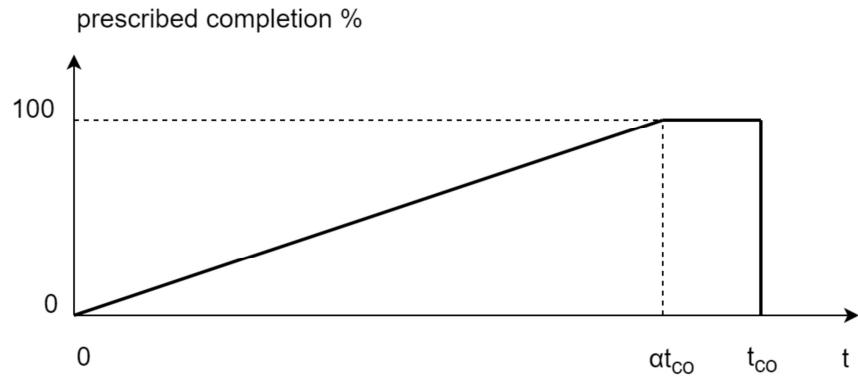


Figure 3.4: Set point generation for an executor controller.

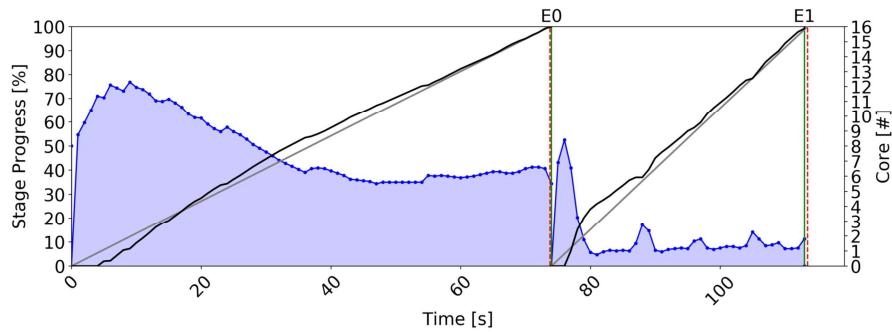


Figure 3.5: CPU cores allocated to the application aggregate-by-key running on a xSpark worker node. The blue line represents the allocated cores over the time, gray and black line are desired and actual progress rate respectively. Green line represents the obtained stage ending, while dashed red line is the desired one.

Notice that it is possible that the controller computes a negative value for $c(k)$, the CPU cores that need to be allocated. To fix this problem $c(k)$ must be clamped between a minimum c_{\min} and a maximum c_{\max} . To maintain consistency, we need to recompute the state $x_C(k)$ as

$$x_C(k) = \frac{c(k)}{K} - a_{\%}^o(k) + a_{\%}(k)$$

In Figure 3.5 the cores allocated to an application executor running on a worker node are shown. The running application is *aggregateby-key* and is composed by two stages. The black line represents the progress of the stage, meanwhile the gray one is the desired progress rate. The goal of the controller is to reduce the error, i.e., the distance between the two lines. We want the black line to follow as much as possible the gray line.

PROBLEM STATEMENT AND SOLUTION OVERVIEW

BIG data applications are widely used in the industry and research fields. Specialized distributed frameworks are used to execute these applications on clusters of computers, very often made by cloud computing resources, like virtual machines and virtual storage. Apache Spark is probably the most popular big data processing framework. Spark organizes computations in directed acyclic graphs (DAGs) or *Parallel Execution Plans (PEPs)* and its declarative API offers the capability to make transformations to datasets and return results to the client program, usually a standard Java, Python or Scala application.

Spark starts executing a program by identifying jobs, delimited by the presence of actions in the code, and stages (within jobs), delimited by operations that require data to be shuffled (i.e., moved among executors), thus breaking locality. Indeed, Spark distinguishes between narrow and wide transformations specifically for this purpose; the former do not shuffle data (e.g., map, filter, etc.), while the latter do (e.g., reduceByKey, etc.). Spark identifies all the operations to be executed, up to the first action, and materializes them in a *PEP* represented as a directed acyclic graph. The *PEP* is not the control-flow graph of the job's code. It does not contain branches and loops since they were already resolved during the execution of the driver program. The *PEP* defines the execution order among stages, and defines the extent to which stages can be executed in parallel. For each job, Spark computes a *PEP* to maximize the parallelism while executing an application. In fact, a stage is, by definition, executed in parallel, and also different stages can be executed at the same time. For this reason, Spark materializes *PEPs* as directed acyclic graphs of stages while the complete *PEP* of an application is simply the sequence of the *PEPs* of its jobs. A Spark application is usually composed by several jobs executed using a LIFO queue, therefore the *PEP* of a Spark application is composed by the sequence of the *PEPs* generated by the application jobs (an action corresponds to a job). In fact, Spark does not "compile" the code of the driver program to generate the application *PEP*, instead it incrementally generates it as soon as an action is reached.

4.1 PROBLEM STATEMENT

The literature contains several works exploring adaptation capabilities, formal guarantees or response time estimation for Spark applications

based on their *PEP*-based structure [48, 16, 15], assuming that the *PEP* of the application does not change with respect to different data input or parameters. However, the *PEP* uniqueness assumption does not hold if conditional branches or loops are present in the control flow of the client program.

This is particularly relevant in Spark because of the possibility of evaluating partial results through actions. In fact, these values can be used in the code as part of conditional expressions that can create branches in the control flow graph. In this cases the conditional branches govern the final structure of the *PEP* and also the operations that form a stage while loops influence the number of repetition of either transformations, stages or actions.

Many of the research studies [35, 77, 48, 56] use the execution graph to establish the amount of work to be done, the degree of parallelism, duration and other specific elements of the application. In all these case studies the assumption that the graph does not change is always taken for granted, since many loops and conditional branches are embedded in the code (eg Filter, map). As mentioned earlier in this document, this condition is not verified when the code contains explicit loop and conditional statements.

As an example, let's take a simple application made by the concatenation of two jobs: the first gets some records from an input file and filters them according to certain criteria, then the second one sorts the records and returns the first x records. We could think of executing the second job only if the number of the filtered records c , (i.e. the output of the first job), is greater than zero. The execution graph would then be composed by two jobs if $c \geq 0$ and contain only the first job if $c = 0$. This example shows how partial results (c) returned to the driver program by Spark actions can be evaluated in conditional expressions that decide the outcome of conditional statements like loops or branches, thus materializing different execution graphs.

To bypass this problem, xSpark and many other approaches [77, 48] leverage an initial profiling phase to retrieve the execution graph and collect performance data. Considering again the simple example used above, running a profiling phase on it would return the execution graph generated by the input data fed to the application. However, there is no way to adapt the execution to the actual situation, in case the actual data drives the execution to walk through a different execution path than the one used to profile the application. In the end, the quality of the results is impacted. Consider that, even if a conservative approach were used by retrieving the worst-case execution graph (i.e., the graph corresponding to the two jobs in the previous example), this would result in an over-allocation of resources and/or over-estimation of the execution time. Conversely, if we adopted the best case (one job), we would forecast too few resources and a too short execution time.

```

1  from pyspark import SparkContext
2  def run(numIterations, threshold):
3      sc = SparkContext('local','example')
4      x = sc.textFile(...).map(...).groupBy(...)
5          .map(...).aggregate(...)
6      y = sc.textFile(...).map(...).groupBy(...)
7          .map(...).aggregate(...)
8      if x > threshold and y > threshold:
9          for i in range(numIterations):
10             z = sc.parallelize(...).map(...).sort(...).take(10)
11     if x > y:
12         w = sc.parallelize(...).map(...).filter(...).count()

```

Figure 4.1: Example Spark application with conditional branches and loops.

User parameters or local can also affect the execution graph, so they must be considered in a sound analysis. For example, Figure 4.1 shows the code of application that has two input parameters *numIterations* and *threshold*. Its execution graph depends on both the user parameters and the input dataset. The first two *aggregate* jobs¹ are always executed (line 4 and 6) and the results are assigned to variables *x* and *y*, respectively. Line 8 checks if both variables are greater than *threshold*. If it is the case, a *take* job (line 10) is repeated *numIterations* times (*for* loop). Finally, if *x > y* (line 11) *count* (line 12) is executed.

Four possible execution graphs (Figure 4.2) correspond to the code in Figure 4.1. They are: i) the two cascaded *aggregates* (case where the two conditional statements are both false) ii) the two cascaded *aggregates* and *take* repeated *numIterations* times (case where the first conditional statement is true and the second is false) iii) the two cascaded *aggregates* and *count* (case where the second conditional statement is true and the first is false), and iv) the concatenation of the two *aggregates*, *take* repeated *numIterations* times, and *count* (if both conditional statements are true) [12].

4.2 SOLUTION OVERVIEW

This section contains a functional level description of the proposed solution, of its components and their interactions to make the solution work.

Our solution is based on xSpark, a modified version of Apache Spark, developed at Politecnico di Milano [16, 15], that has demonstrated the capability to execute deadline-constrained single-*PEP* applications by using resources more efficiently than Spark would do in running the same applications. xSpark is able to use less resources than native Spark and can complete executions with less than 1% error in terms of set deadlines.

¹ In Spark *aggregateByKey* is a transformation while *aggregate* is an action.

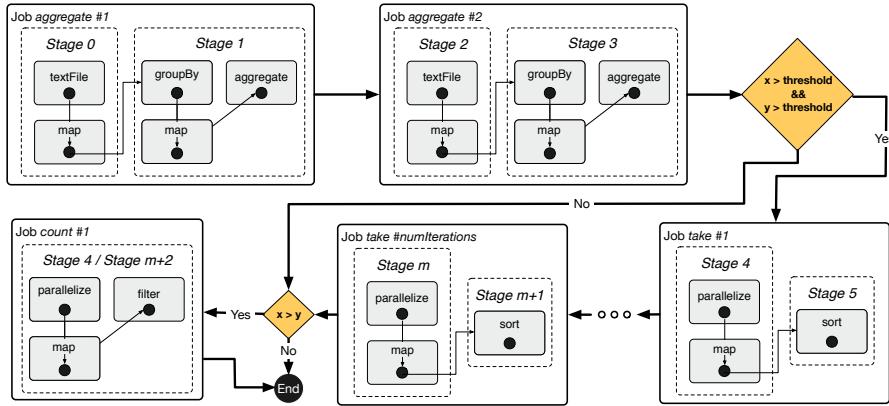


Figure 4.2: The four *PEPs* of the application of Figure 4.1.

This thesis presents $xSpark_{SEEPPEP}$, a toolchain providing the capability to manage the efficient execution of deadline-based QoS constrained multi-*PEP* Spark applications.

$xSpark_{SEEPPEP}$ is the result of the integration of *SEEPPEP* [12], a tool exploiting symbolic execution techniques to generate the path condition associated to each possible *PEPs* produced by different inputs and parameters, with a modified version of *xSpark*.

SEEPPEP is a four-fold stage-set composed by the: i) application of a lightweight symbolic execution of the Spark application driver program to derive a representative set of the control-flow paths execution conditions in the program, ii) use of these execution conditions with a search-based test generation algorithm, to compute sample input datasets that make each path to execute, iii) execution of the target application with these datasets as input, to profile the *PEP* generated by each path, and synthesize the *PEP** accordingly, iv) generation of an artifact called *GuardEvaluator* that returns the feasible *PEPs* given a partial set of concrete values of the symbolic variables.

Moreover, $xSpark_{SEEPPEP}$ generates a launcher with a synthesized dataset for each *PEP* and an artifact to retrieve the feasible *PEPs* given a set of symbolic variables resolved to a value. Finally, we integrated this approach with *xSpark*, an extension of Spark that can control the duration of Spark applications according to specified deadlines through dynamic resource allocation.

We take advantage of the information in the *PEP** extracted with *SEEPPEP* to extend *xSpark* with the capability of tuning its adaptation strategy according to the worst-case behaviour of the application. At runtime, our extended version of *xSpark* uses the *GuardEvaluator* to refine the control policy by recomputing the worst-case estimation every time the current worst-case refers to a program path for which the execution condition stored in the *PEP** becomes unsatisfiable.

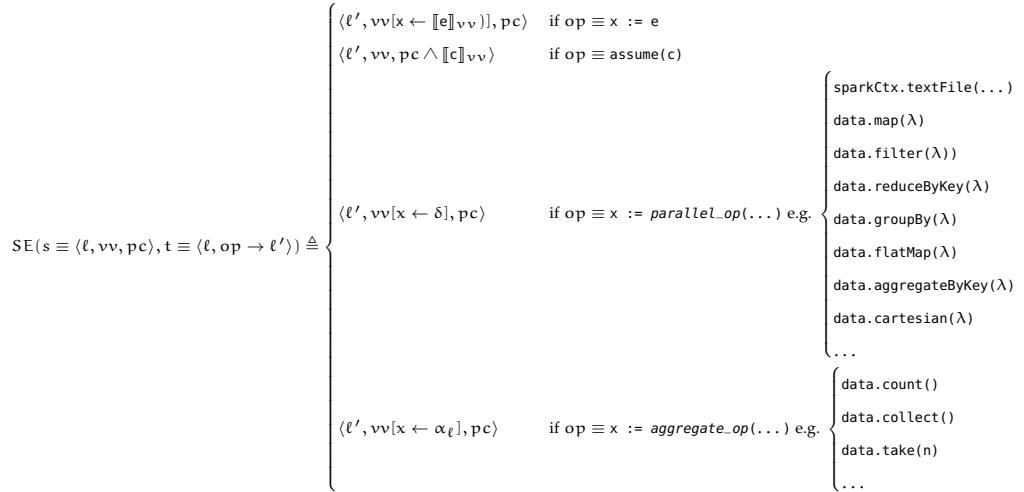


Figure 4.3: Symbolic execution algorithm of SEEPEP.

The evaluation of shows that SEEPEP is able to effectively extract all the *PEPs* generated by Spark applications and that xSpark reduces the number of deadline violation thanks to the presented integration.

4.2.1 $x\text{Spark}_{\text{SEEPEP}}$

This section describes how *SEEPEP* integrates within xSpark: the resulting tool-chain is called $x\text{Spark}_{\text{SEEPEP}}$. *SEEPEP* produces the path conditions associated with the different *PEPs* of the application, a set of test cases for each *PEP* for profiling, and a *GuardEvaluator* to allow xSpark to select the most appropriate *PEP* at runtime. At each execution step, *GuardEvaluator* always returns the *PEPs* whose associated path conditions still hold true.

Figure 4.4 shows the main elements of the tool-chain and exemplifies it on profiling and controlling the example application of Figure 4.2 (*app*, hereafter). As first step, *SEEPEP* generates the 4 (n, in general) launchers, which activate the four (n) *PEPs* of the program, and an application specific *GuardEvaluator*.

The new toolchain associates each *PEP* to its path condition, it uses the generated launchers to obtain the profiling data of each *PEP* and synthesize the *PEP** for the application that is then stored in the master node in component *PEP* Store*. Finally the generated *GuardEvaluator*, which implements a common interface to be dynamically instantiated and used without a static import in the source code of xSpark, is also stored in the master node (component *GuardEvaluator Store*).

After this phase, *app* can be executed and controlled by xSpark. Since the application parameters can be part of a path condition as symbols, we modified component *SparkSubmit* to store their values.

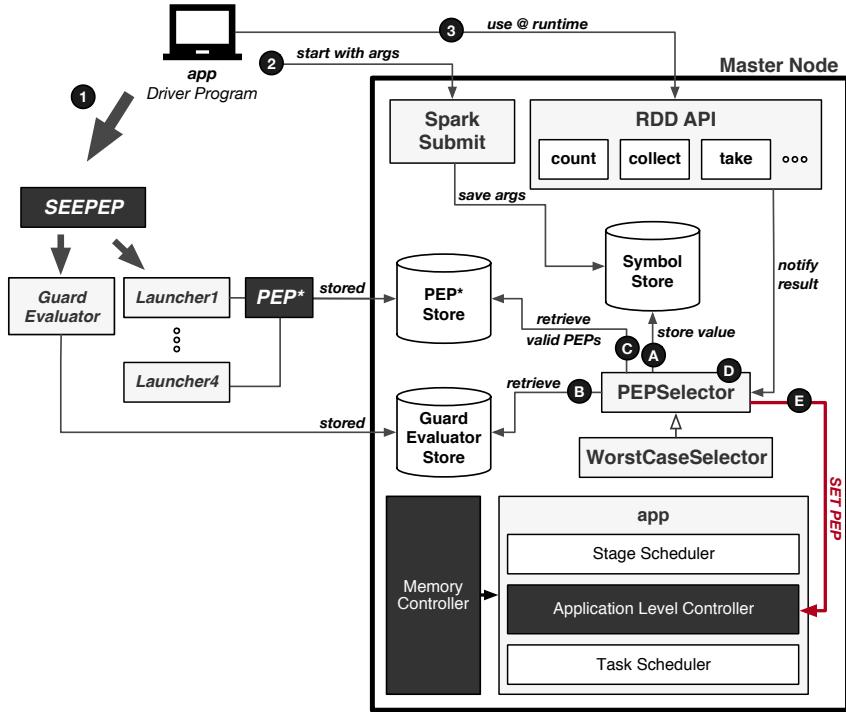


Figure 4.4: xSpark SEEPEP

A symbol is identified by application name, action name, code line in the driver program, and a counter to take multiple executions into account (i.e., loops or recursive functions).

At runtime, everytime an action is executed, that is, a result is computed and returned to the driver program, our modified version of the *RDD API* notifies component *PEPSelector* that a new result is available. This component is in charge of selecting the *PEP* and its profiling data then used by *Application Level Controller* to compute the local deadlines for the next stages and thus to provision resources. *PEPSelector* saves computed results into component *Symbol Store*, retrieves an instance of the dedicated *GuardEvaluator*, and feeds it with all the symbols resolved by the aforementioned results. *GuardEvaluator* [12] returns the list of *PEP* whose path conditions still hold.

This means, for example, that at the beginning of the execution of function *run* of *app*, four *PEPs* are valid since neither *x* nor *y* have been resolved to a value. The job at line 4 produces the value of *x* and if the value is less than or equal to threshold, the if statement of line 8 is not evaluated. Therefore, even if the value of *y* is still unknown, *GuardEvaluator* only returns two *PEPs*, that is, the only two *PEPs* whose path conditions still hold: it excludes all the path conditions that depends on the expression *x > threshold*. This way, since the *PEP* is updated constantly, xSpark becomes aware of what has been actually done, and can use this information to refine resource provisioning.

Note that *PEPSelector* receives all the valid *PEPs* and computes the next *PEP* to use. This selection can be customized by the user. Currently, we always select the worst-case *PEP*, that is, the *PEP* with the greatest number of remaining stages to be conservative and minimize deadline violations. If one wanted to optimize different performance indicators (e.g., deadlines are not strict and used resources must be minimized), the selection could privilege a *PEP* that corresponds to an average case instead of the worst one.

IMPLEMENTATION

IN this chapter we show the implementation details of $xSpark_{SEEP}$, which consists of the modifications to existing xSpark component, new components added to xSpark, *SEEPEP* concrete application *launchers* and the *Python tool* embedded in $xSpark_{SEEP}$ that was used to launch the experiments that generated the results for the evaluation of the solution.

5.1 OVERVIEW

Figure 5.1 shows a simplified overview of the components of the solution. New components are highlighted with a yellow dotted-pattern background, modified components are highlighted with a grey background.

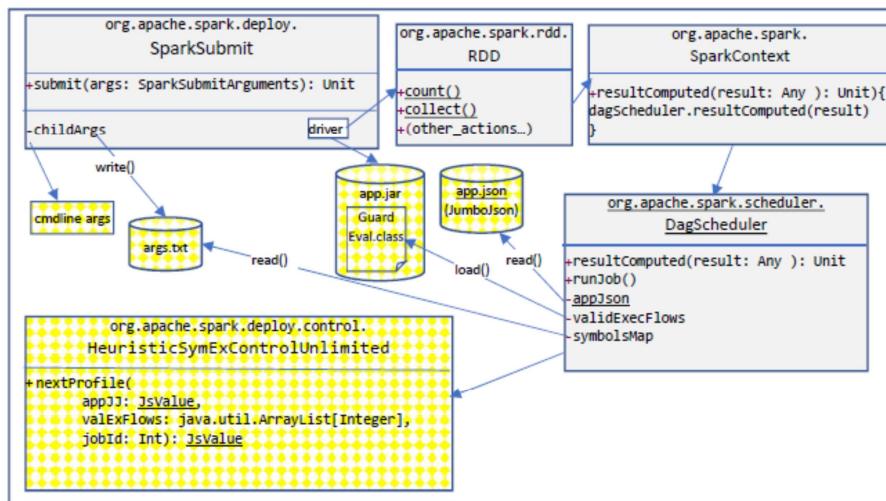


Figure 5.1: Simplified solution components overview.

5.1.1 Background: Current xSpark Heuristic

xSpark uses a heuristic to compute per-stage deadlines and to estimate how many cores must be allocated for a stage to successfully fulfill the deadline. See Section 3.2 for an exhaustive description of its functionality.

5.1.2 Current xSpark Scheduling Limitation

At runtime, an annotated DAG or *Parallel Execution Plan*, (from here onwards *PEP*) allows the xSpark scheduler function to comprehend how much work has already been completed and how much work still needs to be done. This means that xSpark can only optimize the allocation of the resources if the execution of all jobs of the application use the same *PEP*. This might not always be the case, for example when the code contains branches or loops, because these might need to be resolved in different ways at runtime. This is a severe limitation of the xSpark capability to manage real-world applications.

5.2 IMPLEMENTATION SCOPE AND OBJECTIVE

The present work, by addressing the xSpark limitation explained above, aims at extending the scope of applicability of xSpark, enhancing it with the capability to manage the case of applications that can potentially generate, at runtime, a different *PEP* at each execution. The code in this kind of applications includes conditional branches or iterative loops whose outcomes can only be resolved at runtime because they depend on user input values or results from previous computations that cannot be predicted or folded to constant values by the compiler.

5.2.1 Symbolic Execution

Executing a program symbolically means to simultaneously explore multiple paths that a program could take under different inputs. The key idea is to allow a program to take on symbolic – rather than concrete – input values. Execution is performed by a symbolic execution engine, which maintains for each explored control flow path: (i) a first-order Boolean formula that describes the conditions satisfied by the branches taken along that path, and (ii) a symbolic memory store that maps variables to symbolic expressions or values.

When a conditional branch is met, both sides of the branch are executed. Branch execution updates the formula, while assignments update the symbolic store. A symbolic execution tree is generated with an execution state associated with each node, containing the statement to be executed, the symbolic store, and the path conditions (a formula that expresses a set of assumptions on the symbols). The leaves of the tree identify the end of the computations, and tracing back from each leaf up to the root of the tree allows us to reconstruct, in reverse order, all the possible execution paths of the program. Our implementation exploits the functionality provided by *GuardEvaluator*, an extension to Spark applications explained in Section 5.6, which is

generated by an external implementation of a lightheadsymbolic execution algorithm.

5.2.2 $xSpark_{\text{SEEPEP}}$ vs. $xSpark$

The first important difference between $xSpark_{\text{SEEPEP}}$ and $xSpark$ is in the profiling of the applications. $xSpark$ requires the generation of a single *PEP* profile per application, while $xSpark_{\text{SEEPEP}}$ requires a family of *PEP* profiles, one for each possible execution path, each of them associated to a unique set of "Path Conditions". Profile information is collected in special JSON files, called "JSON profiles". Listing 5.1 shows an example of a JSON profile.

Listing 5.1: Example of JSON profile.

```

1 {
2     "o": {
3         "RDDIds": {
4             "o": {
5                 "callsite": "textFile at PromoCalls.java:34",
6                 "name": "hdfs://10.0.0.4:9000//user/ubuntu/
7                     last24HoursLocalCalls.txt"
8             },
9             "i": {
10                 "callsite": "textFile at PromoCalls.java:34",
11                 "name": "hdfs://10.0.0.4:9000//user/ubuntu/
12                     last24HoursLocalCalls.txt"
13             }
14         },
15         "actual_records_read": 597900000.0,
16         "actual_records_write": 597900000.0,
17         "actualtotalduration": 43000.0,
18         "bytesread": 55618499408.0,
19         "byteswrite": 0.0,
20         "cachedRDDs": [],
21         "duration": 17000.0,
22         "genstage": false,
23         "id": 0.0,
24         "io_factor": 1.0,
25         "jobs": {
26             "o": {
27                 "id-symb": "count_PromoCalls.java:34_o",
28                 "stages": [
29                     0
30                 ]
31             },
32             "i": {
33                 "id-symb": "count_PromoCalls.java:42_o",
34                 "stages": [
35                     1
36                 ]
37             }
38         }
39     }
40 }
```

```

35     },
36     "2": {
37         "id-symb": "count_PromoCalls.java:45_o",
38         "stages": [
39             2
40         ]
41     },
42 },
43 "monocoreduration": 705582.0,
44 "monocoretotalduration": 2170942.0,
45 "name": "count at PromoCalls.java:34",
46 "nominalrate": 847385.562556868,
47 "nominalrate_bytes": 78826414.8008311,
48 "numtask": 500,
49 "parentsIds": [],
50 "recordsread": 597900000.0,
51 "recordswrite": 0.0,
52 "shufflebytesread": 0.0,
53 "shufflebyteswrite": 0.0,
54 "shufflerecordsread": 0.0,
55 "shufflerecordswrite": 0.0,
56 "skipped": false,
57 "t_record_ta_executor": 0.0011801003512293025,
58 "totalduration": 43000.0,
59 "weight": 3.0384051747351832
60 },
61 "1": {
62     "RDDIds": {
63         "2": {
64             "callsite": "textFile at PromoCalls.java:35",
65             "name": "hdfs://10.0.0.4:9000//user/ubuntu/
66                 last24HoursLocalCalls.txt"
67         },
68         "3": {
69             "callsite": "textFile at PromoCalls.java:35",
70             "name": "hdfs://10.0.0.4:9000//user/ubuntu/
71                 last24HoursLocalCalls.txt"
72         },
73         "4": {
74             "callsite": "filter at PromoCalls.java:36",
75             "name": "MapPartitionsRDD"
76         }
77     },
78     "actual_records_read": 597900000.0,
79     "actual_records_write": 597900000.0,
80     "bytesread": 55618499408.0,
81     "byteswrite": 0.0,
82     "cachedRDDs": [],
83     "duration": 13000.0,
84     "genstage": false,
85     "id": 1.0,
86     "io_factor": 1.0,

```

```

85     "monocoreduration": 737425.0,
86     "name": "count at PromoCalls.java:42",
87     "nominalrate": 810794.3180662441,
88     "nominalrate_bytes": 75422584.54486898,
89     "numtask": 500,
90     "parentsIds": [],
91     "recordsread": 597900000.0,
92     "recordswrite": 0.0,
93     "shufflebytesread": 0.0,
94     "shufflebyteswrite": 0.0,
95     "shufflerecordsread": 0.0,
96     "shufflerecordswrite": 0.0,
97     "skipped": false,
98     "t_record_ta_executor": 0.001233358421140659,
99     "weight": 1.9935654473336273
100   },
101   "2": {
102     "RDDIds": {
103       "5": {
104         "callsite": "textFile at PromoCalls.java:44",
105         "name": "hdfs://10.0.0.4:9000//user/ubuntu/
106           last24HoursAbroadCalls.txt"
107       },
108       "6": {
109         "callsite": "textFile at PromoCalls.java:44",
110         "name": "hdfs://10.0.0.4:9000//user/ubuntu/
111           last24HoursAbroadCalls.txt"
112       },
113       "7": {
114         "callsite": "filter at PromoCalls.java:45",
115         "name": "MapPartitionsRDD"
116       }
117     },
118     "actual_records_read": 597900000.0,
119     "actual_records_write": 597900000.0,
120     "bytesread": 55618499408.0,
121     "byteswrite": 0.0,
122     "cachedRDDs": [],
123     "duration": 13000.0,
124     "genstage": false,
125     "id": 2.0,
126     "io_factor": 1.0,
127     "monocoreduration": 727935.0,
128     "name": "count at PromoCalls.java:45",
129     "nominalrate": 821364.5449112902,
130     "nominalrate_bytes": 76405859.60010166,
131     "numtask": 500,
132     "parentsIds": [],
133     "recordsread": 597900000.0,
134     "recordswrite": 0.0,
135     "shufflebytesread": 0.0,
136     "shufflebyteswrite": 0.0,

```

```

135     "shufflerecordsread": 0.0,
136     "shufflerecordswrite": 0.0,
137     "skipped": false,
138     "t_record_ta_executor": 0.001217486201705971,
139     "weight": 1.0
140   }
141 }
```

The profiling information for a *xSpark_{SEEP}* application is obtained by combining the JSON profiles, obtained by driving the application with different sets of input data so to drive the execution of all the possible execution paths, into a JSON file that we will call with a jargon “JumboJSON”, as shown in Figure 5.2. Furthermore, each single json profile is enhanced with information about the jobs composing the application, as shown in Figure 5.3. Inside *xSpark_{SEEP}* is kept a symbolic memory store that maps symbolic values (or symbols) to actual values. To this structure, initially empty, a new entry is added every time a symbolic value gets assigned a concrete value. Each entry is a key-value pair containing the symbol as key and the assigned value as value. The convention adopted for naming the symbols is the following:

- **Commandline arguments:** prefix “arg_” followed by an integer reflecting the position of the argument on the commandline. For example: “arg_o”, “arg_1” etc...
- **Program variables:** Spark action name followed by “_”, followed by program name followed by “：“, followed by the program line number where the action is called, followed by “_”, followed by an integer representing the number of times the action in the same line of code is being repeated. For example: “count_PromoCalls.java:34_2”

Table 5.1: Example of Symbolic Memory Store contents.

Entry#	Key	Value
0	arg_o	100
1	arg_1	200
2	arg_2	300
3	count_PromoCalls.java:42_o	2350
4	count_PromoCalls.java:45_o	1920
5	count_PromoCalls.java:45_1	3800

Table 5.1 shows an example of symbolic memory store contents during the execution of the application, run with three commandline arguments having value “100”, “200”, “300” and two Spark actions already executed, of which the second was executed twice. The application is

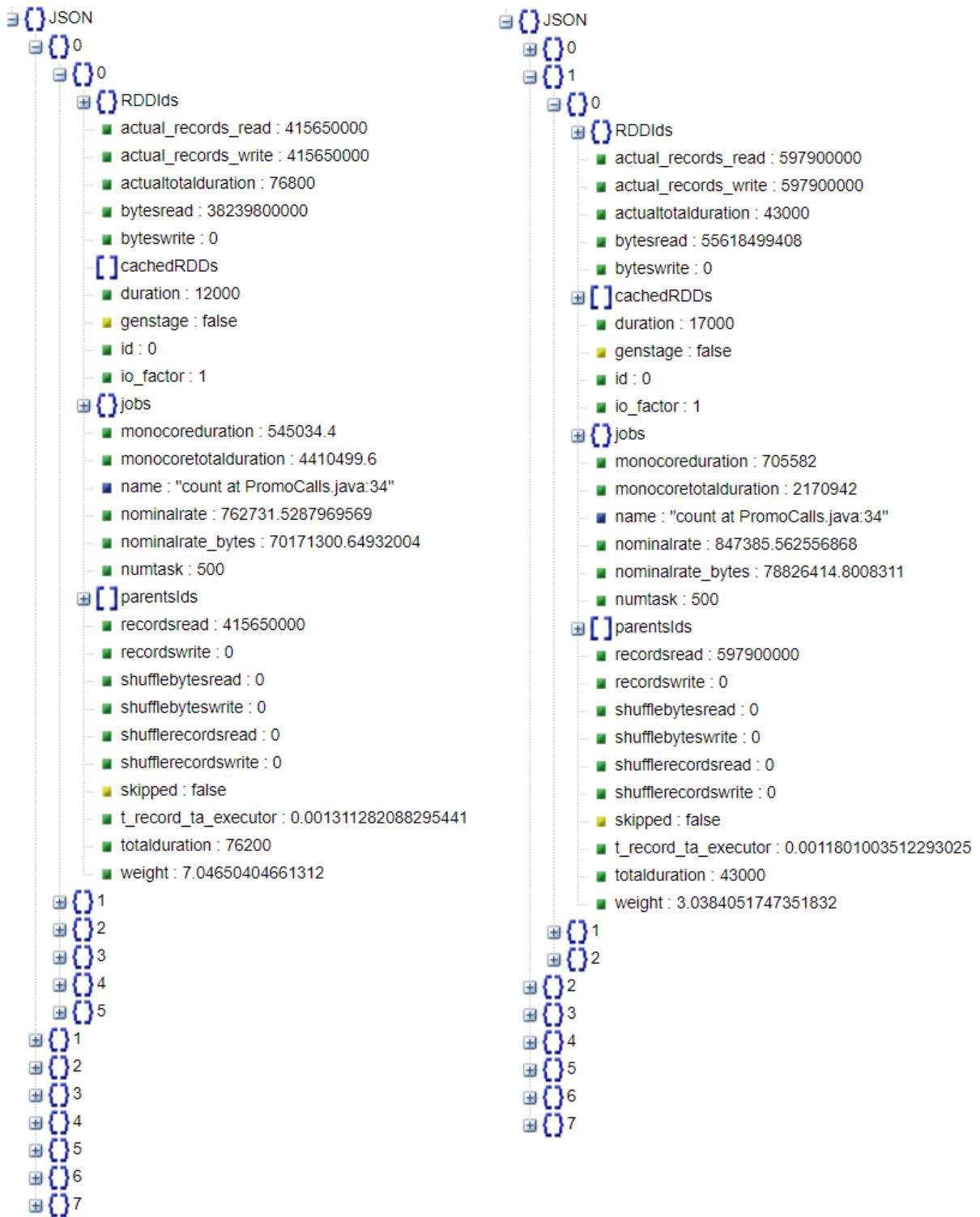


Figure 5.2: Structure of profile JumboJSON. Each profile is identified by its profile id's "0", "1", ..., "7" and, in general, is characterized by its own number of stages and key values. Left and right images show the same JumboJSON profile where profile "0" (left) and profile "1" (right) are expanded, to show their difference in terms of number of stages (6 stages for profile 0 and 3 for profile 1) and totalduration (76,200 ms for profile 0 ms and 43,000 ms for profile 1).

```
{  
    "0": {  
        "0": {  
            "jobs": {  
                "0": {  
                    "id-symb": "count_PromoCalls.java:34_0",  
                    "stages": [  
                        0  
                    ]  
                },  
                "1": {  
                    "id-symb": "count_PromoCalls.java:42_0",  
                    "stages": [  
                        1  
                    ]  
                },  
                "2": {  
                    "id-symb": "count_PromoCalls.java:45_0",  
                    "stages": [  
                        2  
                    ]  
                },  
                "3": {  
                    "id-symb": "collect_PromoCalls.java:51_0",  
                    "stages": [  
                        3  
                    ]  
                },  
                "4": {  
                    "id-symb": "collect_PromoCalls.java:52_0",  
                    "stages": [  
                        4  
                    ]  
                },  
                "5": {  
                    "id-symb": "collect_PromoCalls.java:68_0",  
                    "stages": [  
                        5  
                    ]  
                },  
                ...  
            }  
        },  
        ...  
    }  
}
```

Figure 5.3: Information about jobs in json DAG profile.

also required to implement the *GuardEvaluator* functionality (explained in Section 5.6) by providing a *Java* class implementing a method called `evaluateGuards` that receives in input a Map of a symbolic memory store (as the one described above) and returns a list of the profile id's whose *PEPs* are still executable (i.e. they contain execution paths whose Path Conditions are satisfiable).

5.2.3 A new Heuristic

xSpark_{SEEPEP} implements `HeuristicSymExControlUnlimited`, a new heuristic that extends `HeuristicControlUnlimited`. The class diagram with the relationships between the heuristic classes is shown in Figure 5.4. A new method, `nextProfile`, is implemented by the new heuristic. It takes the following input parameters:

- 1) a json containing the application profiles obtained by the concrete execution of every possible execution path of the application;
- 2) a list of the profile id's that are still satisfiable;
- 3) the id of the job being submitted;

and returns the json of the application profile to be used during the execution of the next job. Keeping in mind that the scheduler uses the `HeuristicSymExControlUnlimited` to estimate how many cores are needed to execute the stage, given the application deadline and the parameters in the application profile, we expect the new heuristic to choose the profile so as not to jeopardize the controller's ability to meet the deadline. This can be achieved by choosing a profile that will lead to not underestimate the cores needed to execute the remaining stages. All the following parameters seem to be good proxies for estimating the remaining computing effort:

- 1) Number of remaining stages to be executed
- 2) Sum of duration of remaining stages to be executed
- 3) Weighted combination of 1 and 2 above

The above parameters can be calculated using the data inside the application json profile. Current implementation of the heuristic uses proxy #1 (number of remaining stages to be executed). It calculates the value for each of the satisfiable profiles, and then selects the profile associated to the maximum value of the proxy. This way, we supply the “worst case” profile to the heuristic, so that the stage deadline is not overestimated and consequently the number of cores to be assigned for the next stage execution is not underestimated.

5.3 APPLICATION PARAMETERS

As explained in Chapter 4, the application parameters can be part of a path condition as they could have been associated to a symbol by the

symbolic executor part of *SEEPEP*. Hence, we have introduced in xSpark a mechanism to intercept and store these application parameters in the *xSpark_{SEEPEP} Symbol Store*, as shown in Figure 4.4. Listing 5.2 shows part of the code of method *submit* of xSpark class *SparkSubmit*, that was modified in order to read the values of the application's runtime arguments passed via the Spark *submit* command and write them as separate lines to textfile *args.txt*. This file is a component of the *PEP* Store*. Records from this file are read by xSpark at a later stage, when lazily executing the application by means of job scheduling.

Listing 5.2: Changes to SparkSubmit method "submit".

```

1 private def submit(args: SparkSubmitArguments): Unit = {
2   val (childArgs, childClasspath, sysProps, childMainClass) =
3     prepareSubmitEnvironment(args)
4   val argsFile = sys.env.getOrElse("SPARK_HOME", ".") + "/conf/
5     args.txt"
6   val bw = new BufferedWriter(new FileWriter(argsFile))
7   if (childArgs.size > 0) {
8     bw.write(childArgs(0) + "\n")
9     bw.write(args.primaryResource + "\n")
10    for ( i <- 1 to childArgs.size - 1 ) {
11      bw.write(childArgs(i)+"\n")
12    }
13  }
14  bw.close()
15  ...

```

5.4 APPLICATION PROFILING

As shown in Figure 4.4, for each set of input parameters identified by *SEEPEP* a *Launcher* is generated. A *Launcher* is a Java class which contains the command to run the application with a specific set of arguments, which is in a 1:1 relationship with the application parameters of the corresponding *PEP*. An example of Launcher class is shown in Listing 5.3.

A *Profiling* (see Figure 3.1) of the application is done by running it with each *Launcher*'s set of arguments. Each profiling run generates the corresponding *PEP* in a specialized JSON file. At the end of this process, all the generated *PEPs* are packaged into another JSON file, in jargon called *JumboJSON*, to form the *PEP**. All the generated JSON files are then stored along with the *PEP** in the *PEP* Store* on the *Spark Master* server.

Listing 5.3: Example of Launcher Code .

```

1 package it.polimi.deepse.dagsymb.launchers;/*
2  * This file was automatically generated by EvoSuite
3  * Wed May 16 13:17:45 GMT 2018

```

```

4  */
5
6 import it.polimi.deepse.dagsymb.examples.PromoCalls;
7 import it.polimi.deepse.dagsymb.examples.UserCallDB;
8
9 public class Launcher0 {
10
11 //Test case number: 0
12 /*
13 * 1 covered goal:
14 * Goal 1. com.xspark.varyingdag.examples.calls.PromoCalls.
15 * run_driver(IJJI)V: path condition EvoSuiteWrapper_0_0 (id
16 * = 0)
17 */
18
19 public static void main(String[] args) {
20     int threshold = 2772;
21     long minLocalLongCalls = 2772;
22     long minAbroadLongCalls = 1397;
23     int pastMonths = 0;
24     int last24HLocalCallsLength = 0;
25     int last24HLocalCallsSize = 0;
26     int last24HAbroadCallsLength = 3361;
27     int last24HAbroadCallsSize = 2794; // 1397 * 2
28     int MonthCallsLength = 2990;
29     int MonthCallsSize = 3000;
30     int num_partitions = 500;
31     PromoCalls promoCalls0 = new PromoCalls();
32     boolean genData = false;
33     String appName = "";
34     if (args[12] != null && args[12].startsWith("-g")) genData
35         = true;
36     if (args[13] != null && !args[13].startsWith("-"))
37         appName
38         = args[12];
39     //UserCallDB.addCallsToLast24HoursAbroadCalls(3361, 1397);
40     //UserCallDB.addCallsToLast24HoursAbroadCalls(3361, 1397);
41     //promoCalls0.run(2772, 2772, 1397, 0);
42     promoCalls0.run(threshold, minLocalLongCalls,
43                     minAbroadLongCalls, pastMonths,
44                     last24HLocalCallsLength, last24HLocalCallsSize,
45                     last24HAbroadCallsLength, last24HAbroadCallsSize,
46                     MonthCallsLength, MonthCallsSize, num_partitions,
47                     genData, appName);
48 }
49 }

```

5.5 PEP*

The previous version of xSpark required a single *PEP* to be present in the *PEP* Store*, so we had to modify the xSpark class *DAGScheduler* to take in account that the data in the *JumboJSON file* now contains

all the *PEPs*. To maintain backwards compatibility, a new variable was introduced, *JumboJson*, to hold the contents of the *PEP** Store. The modified code is in charge of checking the *heuristic type* in the *SparkContext* instance *sc*, to understand if it should expect the contents of the *PEP** Store to be a single *PEP* or a whole set of *PEPs* representing a *PEP**. The heuristic type is initialized with the value of the key *spark.control.heuristic* specified in the xSpark configuration file *spark-defaults.conf*.

Listing 5.4: Changes to class DAGScheduler.scala - reading PEPs.

```

1  val jsonFile = sys.env.getOrElse("SPARK_HOME", ".") +
2    "/conf/" + sc.appName + ".json"
3
4  val appJumboJson = if (Files.exists(Paths.get(jsonFile))) {
5    io.Source.fromFile(jsonFile).mkString.parseJson
6  } else null
7
8  var appJson = if (appJumboJson != null && heuristicType > 2)
9    heuristic.nextProfile(appJumboJson)
10   else appJumboJson

```

5.6 GUARDEVALUATOR

With the term *GuardEvaluator* we collectively refer to the interface class *IGuardEvaluator* and its implementation class defining the *evaluateGuards* method, that is in charge of returning the list of valid profiles (*PEPs*) when it is called with the *HashMap* of the known symbols and their values. The code of the *IGuardEvaluator* interface is shown in Listing 5.5, while Listing 5.6 shows an example of an implementation class and its method *evaluateGuards* for a specific application.

Listing 5.5: Interface class *IGuardEvaluator*.

```

1 package it.polimi.deepse.dagsymb.examples;
2 import java.util.List;
3 import java.util.Map;
4
5 public interface IGuardEvaluator {
6
7   public List<Integer> evaluateGuards(Map<String, Object>
8                                         knownValues);
9 }

```

Listing 5.6: Class *GuardEvaluatorPromoCallsFile*, implementing the *IGuardEvaluator* interface.

```

1 package it.polimi.deepse.dagsymb.examples;
2 import java.util.ArrayList;
3 import java.util.List;
4 import java.util.Map;
5
6 public class GuardEvaluatorPromoCallsFile implements
7     IGuardEvaluator {
8
9
10    @Override
11    public List<Integer> evaluateGuards(Map<String, Object>
12                                         knownValues) {
13        satisfiableGuards = new ArrayList<>();
14
15        extractValues(knownValues);
16
17        evaluateActualGuards();
18
19        return satisfiableGuards;
20    }
21
22    private void evaluateActualGuards() {
23        //path condition evaluation
24        if (
25            ( !arg0_known || arg0 > 100 ) &&
26            ( !arg3_known || arg3 >= 0 ) &&
27            ( !arg3_known || arg3 <= 2 ) &&
28            ( !count_PromoCalls_java_42_0_known || !arg1_known || count_PromoCalls_java_42_0 - arg1 <= 0 ) &&
29            ( !arg2_known || !count_PromoCalls_java_45_0_known || count_PromoCalls_java_45_0 - arg2 > 0 ) &&
30            true) {
31            satisfiableGuards.add(0);
32        }
33
34        if (
35            ( !arg0_known || arg0 > 100 ) &&
36            ( !arg3_known || arg3 >= 0 ) &&
37            ( !arg3_known || arg3 <= 2 ) &&
38            ( !count_PromoCalls_java_42_0_known || !arg1_known || count_PromoCalls_java_42_0 - arg1 <= 0 ) &&
39            ( !arg2_known || !count_PromoCalls_java_45_0_known || count_PromoCalls_java_45_0 - arg2 <= 0) &&
40            true) {
41            satisfiableGuards.add(1);
42        }
43
44        if (
45            ( !arg0_known || arg0 > 100 ) &&
46            ( !arg3_known || arg3 >= 0 ) &&
47            ( !arg3_known || arg3 <= 2 ) &&

```

```

47      ( !count_PromoCalls_java_42_0_known || !arg1_known ||  

48          count_PromoCalls_java_42_0 - arg1 > 0 ) &&  

49      ( !arg3_known || 1 > arg3 ) &&  

50      ( !arg2_known || !count_PromoCalls_java_45_0_known ||  

51          count_PromoCalls_java_45_0 - arg2 > 0 ) &&  

52  true) {  

53      satisfiableGuards.add(2);  

54  }  

55  

56  if (  

57      ( !arg0_known || arg0 > 100 ) &&  

58      ( !arg3_known || arg3 >= 0 ) &&  

59      ( !arg3_known || arg3 <= 2 ) &&  

60      ( !count_PromoCalls_java_42_0_known || !arg1_known ||  

61          count_PromoCalls_java_42_0 - arg1 > 0 ) &&  

62      ( !arg3_known || 1 > arg3 ) &&  

63      ( !arg2_known || !count_PromoCalls_java_45_0_known ||  

64          count_PromoCalls_java_45_0 - arg2 <= 0 ) &&  

65  true) {  

66      satisfiableGuards.add(3);  

67  }  

68  

69  if (  

70      ( !arg0_known || arg0 > 100 ) &&  

71      ( !arg3_known || arg3 >= 0 ) &&  

72      ( !arg3_known || arg3 <= 2 ) &&  

73      ( !count_PromoCalls_java_42_0_known || !arg1_known ||  

74          count_PromoCalls_java_42_0 - arg1 > 0 ) &&  

75      ( !arg3_known || 1 <= arg3 ) &&  

76      ( !arg3_known || 2 > arg3 ) &&  

77      ( !arg2_known || !count_PromoCalls_java_45_0_known ||  

78          count_PromoCalls_java_45_0 - arg2 > 0 ) &&  

79  true) {  

80      satisfiableGuards.add(4);  

81  }  

82  

83  if (  

84      ( !arg0_known || arg0 > 100 ) &&  

85      ( !arg3_known || arg3 >= 0 ) &&  

86      ( !arg3_known || arg3 <= 2 ) &&  

87      ( !count_PromoCalls_java_42_0_known || !arg1_known ||  

88          count_PromoCalls_java_42_0 - arg1 > 0 ) &&  

89      ( !arg3_known || 1 <= arg3 ) &&  

90      ( !arg3_known || 2 > arg3 ) &&  

91      ( !arg2_known || !count_PromoCalls_java_45_0_known ||  

92          count_PromoCalls_java_45_0 - arg2 <= 0 ) &&  

93  true) {  

94      satisfiableGuards.add(5);  

95  }  

96  

97  if (  

98      ( !arg0_known || arg0 > 100 ) &&

```

```

91      ( !arg3_known || arg3 >= 0 ) &&
92      ( !arg3_known || arg3 <= 2 ) &&
93      ( !count_PromoCalls_java_42_0_known || !arg1_known || 
94          count_PromoCalls_java_42_0 - arg1 > 0 ) &&
95      ( !arg3_known || 1 <= arg3 ) &&
96      ( !arg3_known || 2 <= arg3 ) &&
97      ( !arg2_known || !count_PromoCalls_java_45_0_known || 
98          count_PromoCalls_java_45_0 - arg2 > 0 ) &&
99      true) {
100         satisfiableGuards.add(6);
101     }
102
103     if (
104         ( !arg0_known || arg0 > 100 ) &&
105         ( !arg3_known || arg3 >= 0 ) &&
106         ( !arg3_known || arg3 <= 2 ) &&
107         ( !count_PromoCalls_java_42_0_known || !arg1_known || 
108             count_PromoCalls_java_42_0 - arg1 > 0 ) &&
109             ( !arg3_known || 1 <= arg3 ) &&
110             ( !arg3_known || 2 <= arg3 ) &&
111             ( !arg2_known || !count_PromoCalls_java_45_0_known || 
112                 count_PromoCalls_java_45_0 - arg2 <= 0) &&
113             true) {
114             satisfiableGuards.add(7);
115         }
116     }
117
118     private boolean arg0_known; private Integer arg0;
119     private boolean count_PromoCalls_java_42_0_known; private
120         Long count_PromoCalls_java_42_0;
121     private boolean arg2_known; private Long arg2;
122     private boolean arg1_known; private Long arg1;
123     private boolean count_PromoCalls_java_45_0_known; private
124         Long count_PromoCalls_java_45_0;
125     private boolean arg3_known; private Integer arg3;
126
127     private void extractValues(Map<String, Object> knownValues) {
128         arg0_known = (knownValues.get("argo") != null);
129         arg0 = arg0_known ? Integer.parseInt((String) knownValues
130             .get("argo")) : null;
131
132         count_PromoCalls_java_42_0_known = (knownValues.get("count_PromoCalls.java:42_0") != null);
133         count_PromoCalls_java_42_0 =
134             count_PromoCalls_java_42_0_known ? (Long) knownValues
135                 .get("count_PromoCalls.java:42_0") : null;
136
137         arg2_known = (knownValues.get("arg2") != null);
138         arg2 = arg2_known ? Long.parseLong((String) knownValues.
139             get("arg2")) : null;
140
141     }

```

The application is in charge of providing the *GuardEvaluator* as a java class, implementing the interface `IGuardEvaluator`, and packaged inside the jar of the application. This class is loaded dynamically at runtime by the new code added for this purpose to the xSpark class `DAGScheduler` and shown in Listing 5.7.

Listing 5.7: Changes to class DAGScheduler.scala - Loading GuardEvaluator.

```
1 var guardEvalObj:Any = null
2 var guardEvalMethod: java.lang.reflect.Method = null
3 if (heuristicType > 2) {
4 /*
5 * DB - DagSymb enhancements
6 * The following variables are needed to load the
7 GuardEvaluator class from the application jar
8 */
9 val jarfile = new File(appJar)
10 val classLoader = new URLClassLoader(Array(jarfile.toURI.toURL))
11 val guardEvalClass = classLoader.loadClass(guardEvalClassname)
12 val guardEvalConstructor = guardEvalClass.getConstructor()
13 guardEvalObj = guardEvalConstructor.newInstance()
14 val methods = guardEvalClass.getDeclaredMethods
15 for (m <- methods) {
16   m.getName match {
17     case "evaluateGuards" => guardEvalMethod = m
18     case _ =>
19   }
20 } else {
21   guardEvalObj = new core/src/main/scala/org/apache/spark.
22   scheduler.GuardEvaluator
23 }
```

5.7 SYMBOL STORE

In order to take advantage of the symbolic execution, we need to maintain an updated *Symbol Store* containing all the symbols that can be part of a path condition and their associated determinations (assigned values). We added the code into the xSpark class `DAGScheduler` to abstractely represent the *Symbol Store* as a `HashMap[String, Any]`. Each entry of this `HashMap` stores a *known symbol* name and its value. By *known symbol* we mean a symbol that has been associated to a value during the concrete execution of program code. Given this defintion, at the very beginning of the computation the only known symbols are the runtime arguments passed to the application. We added to the xSpark class `DAGScheduler` the code to read the arguments and their values and create the corresponding *Symbol Store* entries. The code is shown in Listing 5.8, where we can notice that the first two arguments loaded when var `iter` is set to negative values (`-2` and `-1`) are respectively the `GuardEvaluator` class name and the application jar name, that are not symbols. They are not kept in the *Symbol Store*, instead they are used to initialize the variables `guardEvalClassname` and `appJar` which will be needed in a later step of the execution to identify the location and dynamically load the class implementing the `GuardEvaluator` function.

Listing 5.8: Changes to class `DAGScheduler.scala` - Initializing Symbol Store.

```

1 var symbolsMap = new java.util.HashMap[String, Any]()
2 var symbolName: String = ""
3 var guardEvalClassname: String = ""
4 var appJar: String = ""
5 val argsFile = sys.env.getOrElse("SPARK_HOME", ".") +
6   "/conf/args.txt"
7 var iter: Int = -2
8 if (Files.exists(Paths.get(argsFile))) {
9   for (line <- Source.fromFile(argsFile).getLines) {
10     iter match {
11       case -2 => guardEvalClassname = line
12       case -1 => appJar = line.split(":")(1)
13       case _   => symbolsMap.put("arg" + iter, line)
14     }
15     iter += 1
16   }
17 ...

```

5.8 HEURISTIC

The heuristic used by xSpark is determined by the value of configuration parameter `spark.control.heuristic` and is an implementation of the class `HeuristicBase`, whose class diagram is shown in Figure 5.4.

In Listing 5.9, we can see that the heuristic *HeuristicControl* is used by default, but other heuristics can be selected. *HeuristicFixed* and *HeuristicControlUnlimited* were already available in xSpark, while we implemented a new heuristic **HeuristicSymExControlUnlimited** to exploit *Symbolic Execution*.

Listing 5.9: Changes to class ControlEventListener.scala - selecting the heuristic.

```

1 val heuristicType = conf.getInt("spark.control.heuristic", 0)
2 val heuristic: HeuristicBase =
3   if (heuristicType == 1 &&
4       conf.contains("spark.control.stagecores") &&
5       conf.contains("spark.control.stagedeadlines") &&
6       conf.contains("spark.control.stage"))
7     new HeuristicFixed(conf)
8   else if (heuristicType == 2)
9     new HeuristicControlUnlimited(conf)
10  else if (heuristicType == 3)
11    new HeuristicSymExControlUnlimited(conf)
12  else
13    new HeuristicControl(conf)

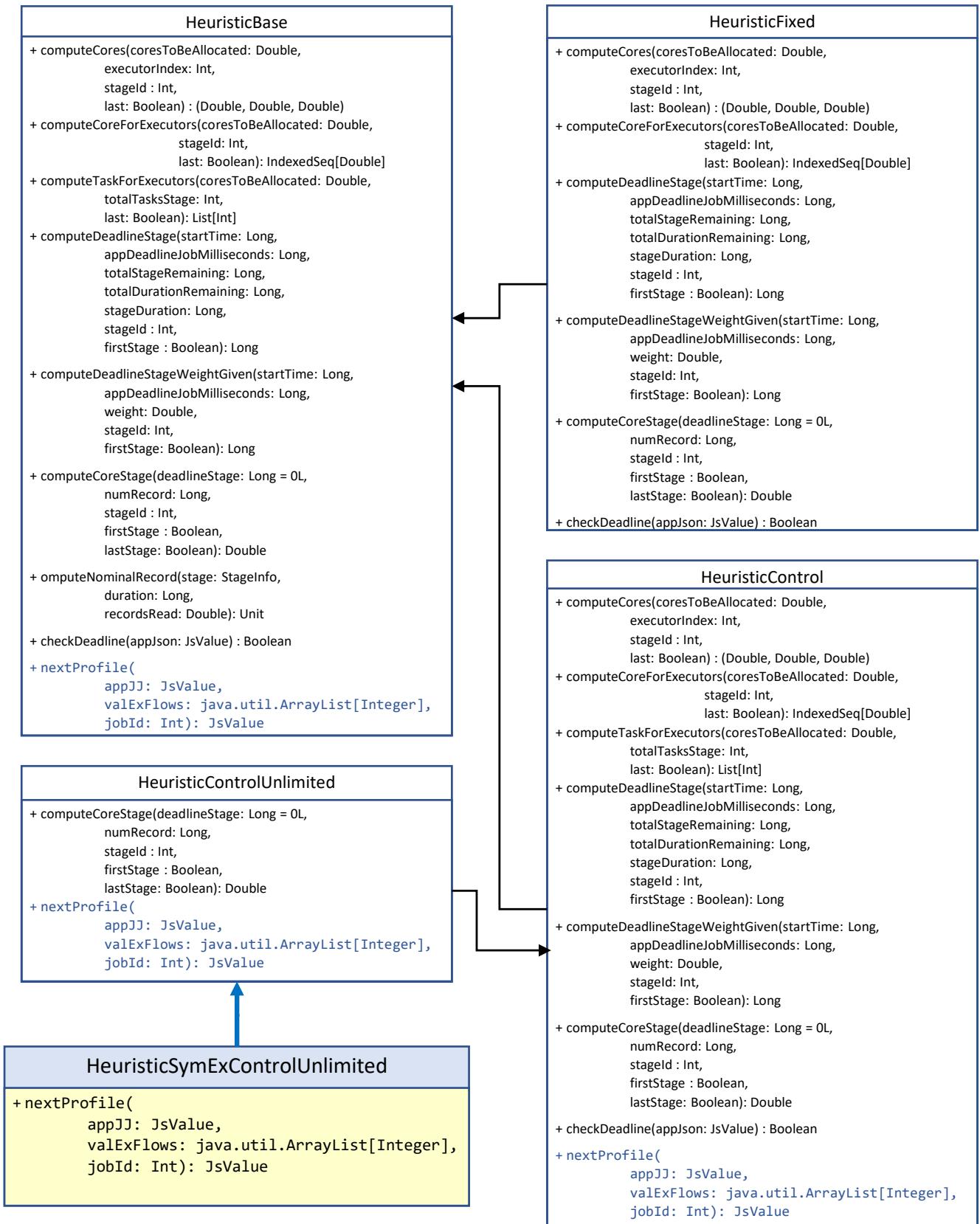
```

HeuristicSymExControlUnlimited extends *HeuristicControlUnlimited* by adding the implementation of a new method, *nextProfile*, taking parameters *appJson*, the JumboJSON containing the PEP* representation and *valExFlows*, a list containing the id's of the valid application profiles, (i.e. the list of the *PEPs* whose path conditions still hold true), and returns the *PEP* of the profile to be used during the executing of the next scheduled job.

Listing 5.10: Class HeuristicSymExControlUnlimited.scala implementation.

```

1 class HeuristicSymExControlUnlimited(conf: SparkConf)
2   extends HeuristicControlUnlimited(conf) {
3   override def nextProfile(appJJ: JsValue,
4     valExFlows: java.util.ArrayList[Integer] = null,
5     jobId: Int = 0): JsValue = {
6     var setP = appJJ.asJsObject.fields
7     val stageId = if (valExFlows != null)
8       setP(valExFlows.get(0).toString()).asJsObject.fields("o")
9       .asJsObject.fields("jobs")
10      .asJsObject.fields(jobId.toString())
11      .asJsObject.fields("stages")
12      .convertTo[List[Int]].sortWith((x, y) => x < y).apply(0)
13     if (valExFlows != null)
14       setP = setP.filter(
15         {case (k,v) => valExFlows.exists(x => x == k.toInt)})
16     var wCaseProfId = setP.toList.map({ case (k, profile) => {
17       val numStage = profile.asJsObject.fields.filter(
18         {case (k, stage) => {
19           !stage.asJsObject.fields("skipped")
```

**Figure 5.4:** Class diagram of Heuristic related classes.

```

20         .convertTo[Boolean]}))
21     .size())
22   (k, numStage)})).reduce({ (x, y) => {
23     if (x._2 > y._2) x else y
24   }})
25   ._1;
26   setP(wCaseProfId)
27 }
}

```

Class `HeuristicSymExControlUnlimited` implementation code is shown in Listing 5.10. The *PEP* selection performed by method `nextProfile` is made by choosing the “*worst case*” among the valid *PEPs*, that is the *PEP* with the maximum number of stages still to be executed, as we want to be conservative and minimize the deadline violations. If we wanted to optimize another performance indicator, like minimum resource utilization in absence of strict deadline commitment, we could choose the profile with an average number of remaining stages to be executed.

STAGE DEADLINE ADJUSTMENT As it was explained in Section 3.2, when a stage is submitted for execution, the heuristic calculates its stage deadline `deadline(sk)` based on the `ApplicationDeadline`, the `SpentTime`, which is the sum of the durations of the stages already executed, and the stage weight `weigth(sk)`. Our solution implies that we select potentially a new *PEP* at each new job execution boundary. This means that we have to face a situation where `heuristic` calculates the new stage deadline using the value of `SpentTime` that is calculated from the *PEP* used to execute the previous stages and the value of `weigth(sk)` that is calculated from the newly selected *PEP*. This is incorrect, since the new *PEP*, in general, has its own number of stages and stage durations. We fixed it by changing the way `SpentTime` is calculated by the following changes to the `xSpark` code: i) `case class SparkStageWeightSubmitted` in class `SparkListener.scala` shown in Listing 5.11, where optional formal parameter `executedstagesduration` is added to the `case` class declaration ii) class method `submitMissingTasks` of class `DAGScheduler.scala` shown in Listing 5.12, where the variable `executedstagesduration` is initialized with the cumulated duration of the executed stages from the selected *PEP* and passed as a parameter to the `SparkStageWeightSubmitted` and iii) method `onStageWeightSubmitted` in class `ControlEventListener.scala` shown in Listing 5.13, where `stageSubmitted.executedstagesduration` is used to calculate the value of the variable `totaldurationremaining`.

Listing 5.11: Changes to companion object `case class SparkStageWeightSubmitted` of class `SparkListener.scala`

```

1 @DeveloperApi
2 case class SparkStageWeightSubmitted

```

```

3 (stageInfo: StageInfo, properties: Properties = null, weight:
4   Long, duration: Long,
5   totalduration: Long, parentsIds: List[Int],
6   nominalrate: Double, genstage: Boolean, stageIds: List[String],
7   executedstagesduration: Long = 0L)
    extends SparkListenerEvent

```

Listing 5.12: Changes to method submitMissingTasks of class DAGScheduler.scala

```

1 private[spark] class DAGScheduler(
2   ...
3   private def submitMissingTasks(stage: Stage, jobId: Int) {
4     ...
5     if (tasks.nonEmpty) {
6       ...
7       if (appJson != null) {
8         val stageJson = appJson.asJsObject.fields(stage.id.
9           toString)
10        val submittedStageId = stage.id
11        var stageId = stage.id
12        // Adding resilience in case of profile does not match no
13        . stages
14        val highestStageIdInProfile = appJson.asJsObject.fields.
15          keys.size - 1
16        if (stage.id > highestStageIdInProfile) {
17          stageId = highestStageIdInProfile
18          logInfo(s"Submitted Stage ID not contained in appJSON
19            profile. Submitted Stage ID: $submittedStageId,
20            " +
21            s"Highest Stage ID in appJSON profile:
22            $highestStageIdInProfile" )
23        }
24        val stageJson = appJson.asJsObject.fields(stageId.
25          toString)
26        val totalduration = appJson.asJsObject.fields("o").
27          asJsObject.fields("totalduration").convertTo[Long]
28        val duration = stageJson.asJsObject.fields("duration").
29          convertTo[Long]
30        val weight = stageJson.asJsObject.fields("weight").
31          convertTo[Long]
32        val stageJsonIds = appJson.asJsObject.fields.keys.toList.
33          filter(id =>
34            appJson.asJsObject.fields(id).asJsObject.fields(".
35            nominalrate").convertTo[Double] != 0.0)
36        val executedstagesduration = appJson.asJsObject.fields.
37          filter(stage =>
38            stage._1.toInt <
39              stageId)
40          .foldLeft(0L){ (acc, elem) =>
41            acc + elem._2.asJsObject

```

```
26     .fields("duration")
27     convertTo[Long] }
28   listenerBus.post(SparkStageWeightSubmitted(stage.
29     latestInfo, properties,
30     weight,
31     duration,
32     totalduration,
33     stageJson.asJsObject.fields("parentsIds").convertTo[
34       List[Int]],
35     stageJson.asJsObject.fields("nominalrate").convertTo[
36       Double],
37     stageJson.asJsObject.fields("genstage").convertTo[
38       Boolean],
39     stageJsonIds))
40   stageJsonIds,
41   executedstagesduration))
42 }
43 else {
44   logError('NO JSON FOR APP: ' + jsonFile)
45   ...
46 }
```

Listing 5.13: Changes to method `onStageWeightSubmitted` of class `ControlEventListener.scala`

```
1 class ControlEventListener(conf: SparkConf) extends
2     JobProgressListener(conf) with Logging {
3     var totaldurationremaining = -1L
4     var totalStageRemaining = -1L
5     var previous_profile_totalStages = 0L
6     var previous_profile_totalduration = 0L
7     ...
8     override def onStageWeightSubmitted
9     (stageSubmitted: SparkStageWeightSubmitted): Unit =
10    synchronized {
11        val stage = stageSubmitted.stageInfo
12        val genstage = if (firstStageId != -1) 1 else 0
13        stageIdToParentsIds(stage.stageId) = stageSubmitted.
14            parentsIds
15
16        if (totaldurationremaining == -1L) {
17            totaldurationremaining = stageSubmitted.totalduration
18            previous_profile_totalduration = totaldurationremaining
19        } else {
20            totaldurationremaining = stageSubmitted.totalduration -
21                stageSubmitted.executedstagesduration
22        }
23        if (totalStageRemaining == -1L) {
24            totalStageRemaining = stageSubmitted.stageIds.size - 1 +
25                genstage
26            previous_profile_totalStages = totalStageRemaining
27        } else {
28    }
```

```

23     totalStageRemaining += stageSubmitted.stageIds.size - 1 +
24         genstage - previous_profile_totalStages
25     stageIdToDuration(stage.stageId) = stageSubmitted.duration
26 ...
27     previous_profile_totalduration = stageSubmitted.totalduration
28     previous_profile_totalStages = stageSubmitted.stageIds.size -
29         1 + genstage
30 ...
31 }

```

STAGE CARDINALITY MISMATCH ADJUSTMENT As it will be presented in Chapter 6 dedicated to the evaluation of the solution, we have to be able to run the applications under xSpark configured to use heuristic `HeuristicControlUnlimited` and providing a single *PEP* corresponding to the best case (i.e. lowest number of stages) and the worst case (i.e. highest number of stages). These test settings are necessary to compare the performance of xSpark with and without *Symbolic Execution* capabilities, however this situation leads to the problem that we call *Stage Cardinality Mismatch*. It happens when the number of stages of the application does not match the number of stages in the *PEP*. We can have two cases: the *PEP* contains more stages than the application or, vice-versa, the *PEP* contains less stages than the application. In the former case there are no problems, and xSpark terminates correctly the execution of the application, while in the latter case xSpark tries to retrieve information from the *PEP* about a non-existent stage and it terminates the execution of the application in error. To avoid this unwanted situation, we have changed the code in method `submitMissingTasks` of class `DAGScheduler.scala` shown in Listing 5.12, where we added the variable `highestStageIdInProfile` that is initialized with the cardinality of the stages in the *PEP* and its value assigned to the variable `stageId` if the `id` of the requested stage is higher than `highestStageIdInProfile`. This creates the effect called *stage id extension* depicted in Figure 5.6

5.9 SYMBOLS

As mentioned in Section 5.7, we have to update the Symbol Store everytime a variable associated to a symbol is evaluated by the concrete execution of the application. We adopted the convention to identify a symbol by the string `arg_n` if it refers to a runtime application argument, where `n` is the position of the argument (e.g. `arg_0`), or by a string obtained by concatenating its *CallSite* and *IterationNumber* separated by an underscore character `_`, where *Call Site* is a string obtained by concatenating *SparkActionName*, *ApplicationClassName*, *SourceLanguageName:SourceLineNumber* separated by an underscore

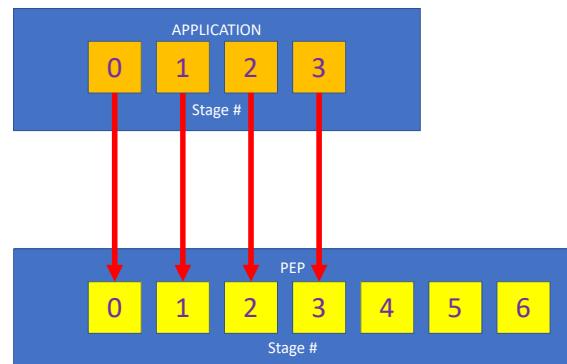


Figure 5.5: Stage Cardinality Mismatch - Acceptable situation.

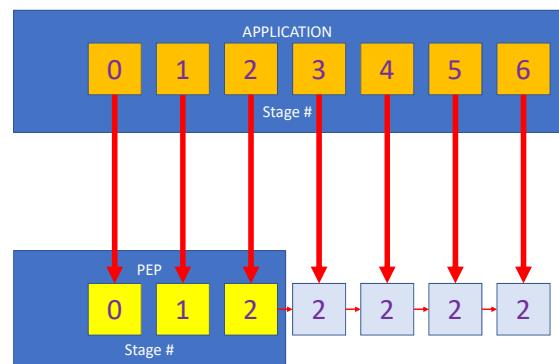


Figure 5.6: Stage Cardinality Mismatch - stage id extension.

character “_”, and *IterationNumber* is an integer starting from 0 and incremented everytime a call is originated by the same *CallSite*. i.e. the same line of code is re-executed (e.g. due to iterative loops). Examples of formal symbol names are shown in Figure 5.7.

```
arg_0
count_PromoCalls_java:45_2
```

Figure 5.7: Example of symbol formal names.

As stated above, a symbol, if not associated to an application runtime argument, is identified by its *CallSite*. This means that to identify the symbols we have to intercept calls coming from their *CallSite*. Since the values associated to a symbol can only be changed by xSpark *actions* which delimit jobs, we added code to the method `runJob` of class `DAGScheduler` to extract the *CallSite*, generate a symbol and push it in the *Symbol Store* everytime the method `runJob` is called. The code is shown in Listing 5.14.

Listing 5.14: Changes to class `DAGScheduler.scala` - method `runJob`.

```
1 var symbolMap = HashMap[String, Int]()
2 var symbolName: String = ""
3 def runJob[T, U](
4   rdd: RDD[T],
5   func: (TaskContext, Iterator[T]) => U,
6   partitions: Seq[Int],
7   callSite: CallSite,
8   resultHandler: (Int, U) => Unit,
9   properties: Properties): Unit = {
10   val actionCallSite = callSite.shortForm.replace(" at ", "_")
11   symbolName = actionCallSite + "_" + symbolMap.getOrDefault(
12     actionCallSite, 0).toString()
13   symbolsMap.put(symbolName, null)
14   symbolMap(actionCallSite) += 1
15 }
```

The variable `actionCallSite` is initialized using the value of the `callSite` parameter. An auxilliary structure, the `HashMap[String, Int]` `symbolMap` keeps track of every `actionCallSite` and counts how many times each of them has called the method `runJob`. The value of the count determines the suffix of each symbol. Symbols are initially assigned a `null` value, and stay in the symbol store waiting to be assigned the result of the *action* originated from the `actionCallSite`.

This task is performed by the new method `resultComputed`, that was added to the xSpark class `DAGScheduler`. `resultComputed` is called by the homonymous method, added to the class `SparkContext`, that passes to it the computed result of the action received from an *action* method in the class `RDD`. Lastly, we have indeed modified the methods

that execute the *actions* in the RDD xSpark class by inserting a call to the method `resultComputed` of the `SparkContext` instance `sc` and passing to it the computed result of the action.

Updating the value of *symbols* in the *Symbol Store* is not the only task fulfilled by the `resultComputed` method of class `DAGScheduler`. It also executes the method stored in the variable `guardEvalMethod` with the map of the known *symbols* getting in return the variable `new_validExecFlows` containing the list of valid profiles. It is also in charge of selecting the profile `appJson` to be used to run the next job. It fulfills this task by calling the method `nextProfile` of the `HeuristicSymExControlUnlimited` instance `heuristic` and passing to it the list of valid profiles. In returns, it gets the json profile to be used in the next job to be run, stored in the variable `appJson`.

The new and modified methods are shown in Listings 5.15, 5.16, 5.17.

Listing 5.15: Changes to class `DAGScheduler.scala` - new method `resultComputed`.

```

1 private[scheduler] def numTotalJobs: Int = nextJobId.get()
2
3 def resultComputed(result: Any ): Unit = {
4     if (heuristicType > 2) {
5         symbolsMap(symbolName) = result
6         val resultType = ClassTag(result.getClass)
7         var new_validExecFlows = guardEvalMethod.invoke(guardEvalObj,
8             symbolsMap).asInstanceOf[java.util.ArrayList[Integer]]
9         if (new_validExecFlows.size() > 0)
10            validExecFlows = new_validExecFlows
11        else
12            println("Warning! GuardEvaluator returned an empty set of
13                profile ids")
13        val highestJobId: Int =
14            if (validExecFlows != null) {
15                appJumboJson.asJsObject.fields(validExecFlows.get(0).
16                    toString())
16                .asJsObject.fields("o").asJsObject.fields("jobs")
17                .asJsObject.fields.keys.max.toInt}
18            else 0
19        appJson = if (numTotalJobs <= highestJobId) {
20            heuristic.nextProfile(appJumboJson,
21                validExecFlows, nextJobId.get())
21        else appJson

```

Listing 5.16: Changes to class `SparkContext.scala` - new method `resultComputed`.

```

1 def resultComputed(result: Any ): Unit = {
2     dagScheduler.resultComputed(result)
3 }

```

Listing 5.17: Changes to class RDD.scala - modified methods count, collect and reduce.

```

1 def count(): Long = {
2   val res = sc.runJob(this, Utils.getIteratorSize _).sum
3   sc.resultComputed(res)
4   res
5 }
6
7 def collect(): Array[T] = withScope {
8   val results = sc.runJob(this, (iter:Iterator[T]) => iter.toArray)
9   val res = Array.concat(results: _*)
10  sc.resultComputed(res)
11  res
12 }
13
14 def reduce(f: (T, T) => T): T = withScope {
15   val cleanF = sc.clean(f)
16   val reducePartition: Iterator[T] => Option[T] = iter => {
17     if (iter.hasNext) {
18       Some(iter.reduceLeft(cleanF))
19     } else {
20       None
21     }
22   }
23   var jobResult: Option[T] = None
24   val mergeResult = (index: Int, taskResult: Option[T]) => {
25     if (taskResult.isDefined) {
26       jobResult = jobResult match {
27         case Some(value) => Some(f(value, taskResult.get))
28         case None => taskResult
29       }
30     }
31   }
32   sc.runJob(this, reducePartition, mergeResult)
33   jobResult.getorElse(throw new UnsupportedOperationException(
34     "empty collection"))
35   val res = jobResult.getorElse(throw new
36     UnsupportedOperationException("empty collection"))
37   sc.resultComputed(res)
38   res
39 }

```

5.10 INTERCEPTING ACTIONS RESULTS

Jobs in xSpark are delimited by *actions* (e.g. *count()*, *collect()* etc...) and are composed by a number of stages, that reflect the *transformations* operated on data. Data are abstracted by Resilient Distributed Datasets (*RDD*'s classes, which contains specialized methods implementing the actions that are in charge of calling the method *runJob* of the *SparkContext* instance *sc*. As we have seen in Section 5.9, the result of

a Spark *action* is associated to a symbol through its *CallSite*. Hence, we have changed the code of the methods implementing the *actions* in the xSpark class RDD to store the result of the action in the local variable *res*, which is passed as a parameter in the call to the *resultComputed* method of the active SparkContext instance *sc*. The process to update the *xSpark_{SEEP} Symbol Store* with the result of the action, which is associated to a unique symbol, is explained in Section 5.9.

5.11 PYTHON TOOL

xSpark-dagsymb is the name of the Python tool, which is part of the *xSpark_{SEEP}* solution that exploits symbolic execution techniques to safely run multi-dag applications in xSpark ¹. It combines two distinct functionalities, application profiling and application execution, which are part of *xSpark_{SEEP}* application lifecycle, in one integrated tool.

The tool is composed by ten principal modules: **xSpark_dagsymb.py**, **launch.py**, **run.py**, **log.py**, **plot.py**, **metrics.py**, **configure.py**, **processing.py**, **average_runs.py**, **process_on_server.py**, in addition to the configuration files **credentials.json**, **setup.json**, **control.json**.

5.11.1 Core Functionality

The **launch.py** module manages the startup of spot request instances on *Amazon EC2* or virtual machines on *Microsoft Azure* and waits until the instances are created and are reachable from the network via their public ip's. Subsequently the **run.py** module receives as input the instances on which to configure the cluster (*HDFS* or *Spark*), configures and runs the applications to be executed and waits for the applications to complete. The module **log.py** downloads and saves the logs created by the applications run. The **plot.py** and **metrics.py** modules respectively generate graphs and calculate metrics. The **process_on_server.py** module can be called to remotely execute the log analysis, graphs generation and metrics calculation on the xSpark master server, and download the results to the client. This option is very useful to speed-up the processing especially in case of sizeable logfiles.

5.11.2 Cloud Environment Configuration

The Cloud environment must be properly initialized in order to allow **xSpark_dagsymb** to access and modify resources in the cloud. We will present the procedure to setup the cloud environment on Microsoft Azure [60];

¹ <https://github.com/gioenn/xSpark-dagsymb.git>

Azure Follow the instructions to create an identity called `service principal`² and assign to it all the required permissions:

- 1) Check that your account has the **required permissions**³ to create an identity.
 - 2) Create an **Azure Active Directory application**
 - 3) Get the **Application ID** and an **Authentication Key**⁴. The **Application ID** and **Authentication Key** values replace respectively the <AZ-APP-ID> and the <AZ-SECRET> values in the credentials.json file described in the next paragraph.

5.11.3 Tool Configuration

The **configure.py** module contains the Config class (shown in Listing 5.18) used to instantiate configuration objects that are initialized with default values. The **credentials.json** file, shown in Figure 5.8, contains *Amazon EC2* and/or *Microsoft Azure* credential information. The **setup.json** contains Cloud environment and *Amazon EC2* and/or *Microsoft Azure* image parameters. The **control.json** file contains xS-spark controller configuration parameters. Information in the **credentials.json**, **setup.json** and **control.json** files are used to customize the configuration object used by other modules during the application execution.

Listing 5.18: Fragment of configuration class Config.

```
1 class Config(object):
2     """
3         Configuration class for xSpark-dagsymb
4     """
5
6     class Heuristic(Enum):
7         CONTROL = 0
8         FIXED = 1
9         CONTROL_UNLIMITED = 2
10        SYMEX_CONTROL_UNLIMITED = 3
11
12        REGION = ""           #####Region of AWS to use#####
13        PROVIDER = ""         #####Provider to be used#####
14        AZ_KEY_NAME = ""       ##### name of Azure private key #####
15        AZ_PUB_KEY_PATH = ""   ##### path of Azure public key #####
16        AZ_PRV_KEY_PATH = ""   ##### path of Azure private key #####
17
18
19
20
21 https://docs.microsoft.com/en-us/azure/azure-resource-manager/
22 resource-group-create-service-principal-portal
23 https://docs.microsoft.com/en-us/azure/azure-resource-manager/
24 resource-group-create-service-principal-portal?view=azure-cli-latest#
25 required-permissions
26 https://docs.microsoft.com/en-us/azure/azure-resource-manager/
27 resource-group-create-service-principal-portal?view=azure-cli-latest#
28 get-application-id-and-authentication-key
```

```

16     AWS_ACCESS_ID = ""      """ Azure access id """
17     AWS_SECRET_KEY = ""      """ Azure secret key """
18     AZ_APPLICATION_ID = ""  """ Azure application id """
19     AZ_SECRET = ""          """ Azure secret """
20     AZ_SUBSCRIPTION_ID = "" """ Azure subscription id """
21     AZ_TENANT_ID = ""       """ Azure tenant id """
22
23     ...
24
25     def __init__(self):
26         self.config_credentials("credentials.json")
27         self.config_setup("setup.json")
28         self.config_control("control.json")
29         self.update_config_parms(self)
30
31     ...
32
33     def config_experiment(self, filepath, cfg):
34
35     def config_setup(self, filepath):
36
37
38 config_instance = Config()

```

AWS and/or MS-Azure Credentials: Open the *credentials_template.json* file and add the credentials for **xSpark_dagsymb** (see instructions below to retrieve missing credentials):

```

1 {
2     "AzTenantId": "", 
3     "AzSubscriptionId": "", 
4     "AzApplicationId": "", 
5     "AzSecret": "", 
6     "AzPubKeyPath": "", 
7     "AzPrvKeyPath": "", 
8     "AwsAccessId": "", 
9     "AwsSecretId": "", 
10    "KeyPairPath": ""
11 }

```

Figure 5.8: Credential template file.

Save the file as *credentials.json*.

How to retrieve your Azure credentials (using the Azure Command Line Interface):

Install the [Azure CLI](#)⁵. Launch the following command from a console terminal:

```
$ az login
```

Note, we have launched a browser for you to login. For old experience with device code, use "az login --use-device-code"

a browser authentication windows is open to allow you to login to the Azure portal. If login is successful, you should get an output similar to the following:

You have logged in. Now let us find all the subscriptions to which you have access...

```
[  
  {  
    "cloudName": "AzureCloud",  
    "id": "< AZ-SUBSCRIPTION-ID >",  
    "isDefault": true,  
    "name": "Microsoft Azure Sponsorship xx",  
    "state": "Enabled",  
    "tenantId": "< AZ-TENANT-ID >",  
    "user": {  
      "name": "*your_username*",  
      "type": "user"  
    }  
  }  
]
```

<AZ-SUBSCRIPTION-ID> and <AZ-TENANT-ID> parameters should be replaced by the actual values from the *credentials.json* file.

Launch the following command from a console terminal to create the private and public RSA cryptography keys:

```
$ ssh-keygen -t rsa
```

Save the generated files in your favorite folder and replace the values <AZ-PUB-KEY-PATH> and <AZ-PRV-KEY-PATH> in the *credentials.json* file respectively with the fully qualified file name of the public and the private keys.

Setup the xSpark and the Virtual Machine Cloud environment: edit the *setup.json* file to set the values to your need. An example using Microsoft Azure VM Cloud Service is shown in Figure 5.9.

⁵ <https://docs.microsoft.com/it-it/cli/azure/install-azure-cli?view=azure-cli-latest>

```

1 {
2     "Provider": "AZURE",
3     "VM": {
4         "Core": 16,
5         "Memory": "100g"
6     },
7     "ProcessOnServer": false,
8     "InstallPython3": false,
9     "Azure": {
10         "ResourceGroup": "xspark-davide",
11         "SecurityGroup": "cspark-securitygroup2",
12         "StorageAccount": {
13             "Sku": "standard_lrs",
14             "Kind": "storage",
15             "Name": "xsparkstoragedavide"
16         },
17         "Subnet": "default",
18         "NodeSize": "Standard_D14_v2_Promo",
19         "Network": "cspark-vnet2",
20         "Location": "centralus",
21         "NodeImage": {
22             "BlobContainer": "vhds",
23             "StorageAccount": "xsparkstoragedavide",
24             "Name": "vm3-os.vhd"
25         }
26     },
27     "Spark": {
28         "ExternalShuffle": "true",
29         "Home": "/opt/spark/",
30         "LocalityWaitRack": 0,
31         "CpuTask": 1,
32         "LocalityWaitProcess": 1,
33         "LocalityWait": 0,
34         "LocalityWaitNode": 0
35     },
36     "xSpark": {
37         "Home": "/usr/Local/spark/"
38     },
39     "SparkSeq": {
40         "Home": "/opt/spark-seq/"
41     }
42 }

```

Figure 5.9: Credential template file.

5.11.4 Application Profiling: PEP and PEP* generation

Profiling is the first logical phase of the *xSpark_{SEEPEP}* application life-cycle. *xSpark_{SEEPEP}* allows to obtain the PEP* of an application. This can be obtained by launching the profiling command **profile_symex** followed by the list of experiment files initialized with the input data set identified by the *SEEPEP launchers* and specifying the number n of iterations the single experiment should be repeated. In profiling mode, experiments are run using the “vanilla” Spark version. Then the **processing.py** module is called to analyze the logs and create the “application profile”, that is a JSON file containing the annotated DAG (*PEP*) of the executed stages plus additional information intended to be used by the controller in the execution phase. For each experiment the **average_runs.py** module is called to create a JSON profile called *<appname>.json* containing the average values of the n *PEPs* obtained by iterating the profiling n times. When all the experiments have produced their averaged *PEP*, a JSON file also called *JumboJSON* containing the *PEP** is created by collecting the averaged *PEPs*. Finally, the *PEP** file and all the *PEPs* are uploaded to the *xSpark* configuration directory of the *xSpark* master server. For example, the command

```
$ python3 xSpark\_dagsymb profile\_symex -r 5
exp\_Louvain\_1.json exp\_Louvain\_2.json
exp\_Louvain\_3.json
```

runs the profiling of application Louvain using launchers data sets 1, 2 and 3 and generating PEP files Louvain_1.json, Louvain_2.json, Louvain_3.json and PEP* file Louvain.json.

5.11.5 Application Execution

Applications are executed using *xSpark*, and require the application profile *<appname>.json* to be present in the *xSpark* configuration directory. The name of the application and the experiment parameters are inserted into a JSON format “experiment files” and passed as commandline arguments to the **submit_symex** command. As an example, an experiment files for CallsExample and one for Louvain are shown here below:

```
CallsExample experiment file example:
{
    "Deadline": 91200,
    "AppName": "CallsExample-6",
    "MetaProfileName": "CallsExample",
    "AppJar": "dagsymb/target/dagsymb-1.0-jar-with-dependencies.jar",
    "AppClass": "it.polimi.deepse.dagsymb.launchers.Launcher",
    "GuardEvaluatorClass": "it.polimi.deepse.dagsymb.examples
```

```

        .GuardEvaluatorPromoCallsFile",
    "NumPartitions": 500,
    "AppConf": {
        "0": {"Name": "threshold", "Value" : 10060000},
        "1": {"Name": "minLocalLongCalls", "Value" : 10060000},
        "2": {"Name": "minAbroadLongCalls", "Value" : 10060000},
        "3": {"Name": "pastMonths", "Value" : 2},
        "4": {"Name": "last24HLocalCallsLength", "Value" : 16100000},
        "5": {"Name": "last24HLocalCallsSize", "Value" : 16140000},
        "6": {"Name": "last24HAbroadCallsLength", "Value" : 20030000},
        "7": {"Name": "last24HAbroadCallsSize", "Value" : 20030000},
        "8": {"Name": "MonthCallsLength", "Value" : 29900000},
        "9": {"Name": "MonthCallsSize", "Value" : 30000000}
    }
}

Louvain experiment file example:
{
    "Deadline": 287040,
    "AppName": "Louvain-2",
    "MetaProfileName": "Louvain",
    "AppRepo": "https://github.com/gioenn/louvain-modularity-spark.git",
    "AppBranch": "master",
    "AppJar": "louvain/target/louvainspark-1.0-SNAPSHOT-jar-with-
                dependencies.jar",
    "AppClass": "it.polimi.dagsymb.launchers.LouvainLauncher",
    "GuardEvaluatorClass": "it.polimi.dagsymb.GuardEvaluatorLouvain",
    "NumPartitions": 500,
    "DataMultiplier" : 1,
    "DeadlineMultiplier" : 1,
    "AppConf": {
        "0": {"Name": "inputFile", "Value" : "/louvain_2.txt"},
        "1": {"Name": "outputFile", "Value" : "/louvain_2-out"},
        "2": {"Name": "parallelism", "Value" : 500},
        "3": {"Name": "minimumCompressionProgress",
               "Value" : -2147483648},
        "4": {"Name": "progressCounter", "Value" : 2},
        "5": {"Name": "delimiter", "Value" : ","},
        "6": {"Name": "size", "Value" : 9999970},
        "7": {"Name": "id1", "Value" : 0},
        "8": {"Name": "id2", "Value" : 1}
    }
}

```

Many files are downloaded from the xSpark Master and Slaves containing application log information and system resources utilization, that are processed to create resource utilization reports and statistics and a series of charts about the application progress, the worker

progress, including application and stages deadlines, actual execution times and cpu utilization. One of the features that was added by this work is the creation of Matlab interactive and zoomable charts that allow an accurate evaluation of the details of the data shown on the charts. Examples of Matlab charts, contained in the Matlb executable file `charts.m` are shown in Figure 5.10, Figure 5.11,Figure 5.12,Figure 5.13

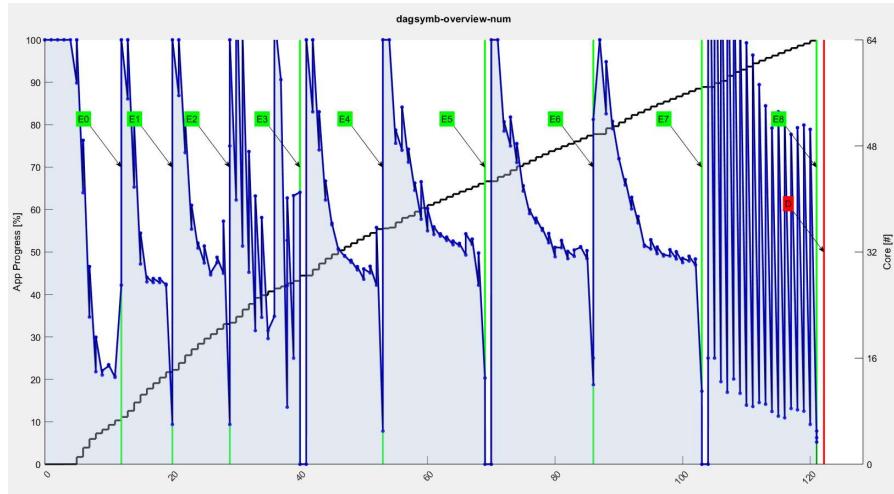


Figure 5.10: Example - Application Matlab Chart.

5.11.6 Download & Requirements

```
$ git clone https://github.com/gioenn/xSpark-dagsymb.git
$ cd xSpark-dagsymb
$ pip3 install -r requirements.txt"
```

5.11.7 xSpark-dagsymb commands

xSpark-dagsymb run command syntax:

```
$ cd xSpark-dagsymb
$ python3 xSpark_dagsymb.py *command [*args]*
```

command [*args] syntax:

```
[setup | reboot | terminate | log | profiling | time_analysis | check |
profile | submit | launch_exp] [*args*]
```

where ***args*** is a set of command-specific arguments list or options.

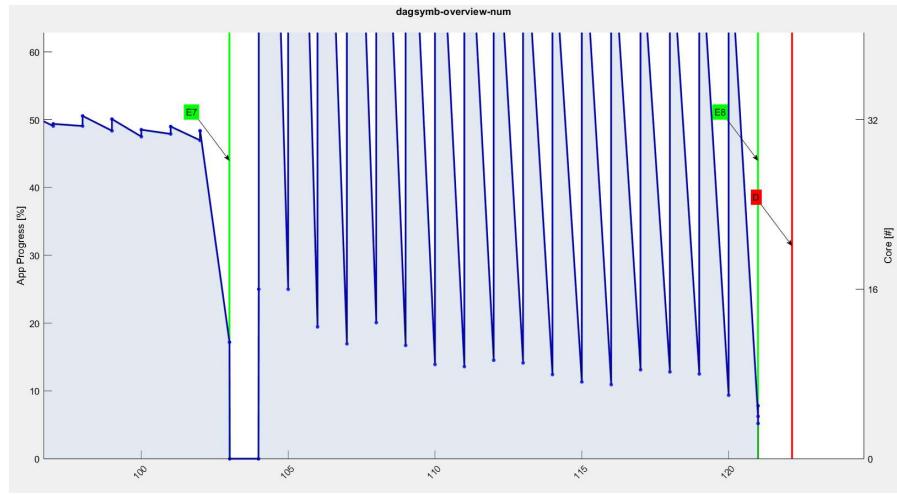


Figure 5.11: Example - Zoomed Details Application Matlab Chart.

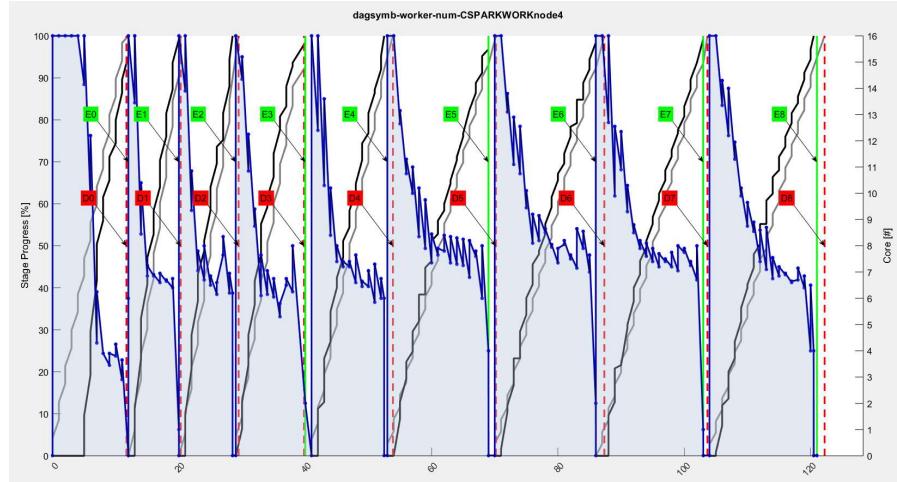


Figure 5.12: Example - Worker Matlab Chart.

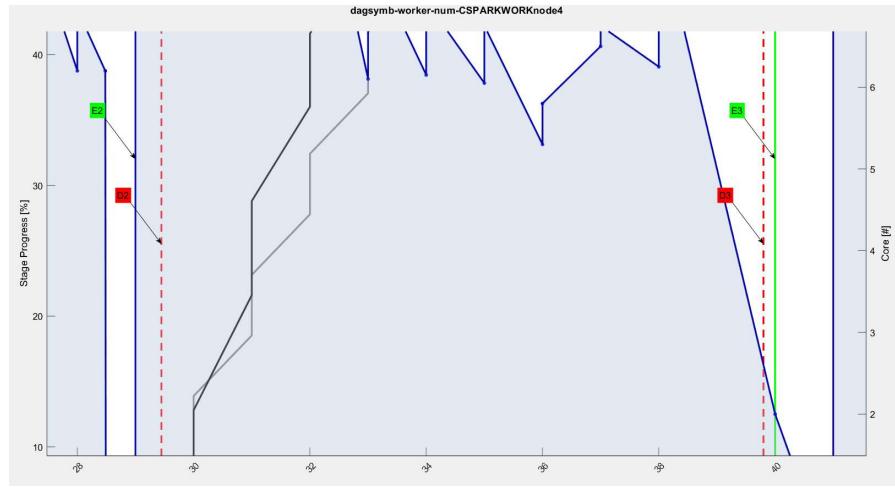


Figure 5.13: Example - Zoomed Details Worker Matlab Chart.

setup command syntax:

```
setup [hdfs | spark | all | generic]
  {[-n | --num-instances] *numinstances*}
  {[-y | --assume-yes]}
```

where *numinstances* is the number of nodes to add to the specified cluster (default is 5), -y or *--assume-yes* option sets default affirmative answer to interactive confirmation requests.

reboot command syntax:

```
reboot [hdfs | spark | all | generic]
```

reboots all nodes in the specified cluster.

terminate command syntax:

```
terminate [hdfs | spark | all | generic]
```

deletes (destroys) all nodes and their connected resources in the specified cluster.

check command syntax:

```
check [hdfs | spark | all | generic]
```

checks the status of all nodes in the specified cluster.

***profile* command syntax:**

```
profile [*exp_file_paths*] {[[-r | --num-runs] *numruns*]
  {[[-R | --reuse-dataset]}}
  {[[-q | --spark-seq]]}
```

where **exp_file_paths** is a non-empty space separated list of experiment file paths, **numruns** is the number of times to repeat the profiling for each experiment file (default is 1), *-R* or *--reuse-dataset* option instructs xSpark to reuse (not to delete) application data in hdfs master node, *-q* or *--spark-seq* option instructs xSpark to use Spark data sequencing home directory.

***submit* command syntax:**

```
submit [*exp_file_paths*] {[[-r | --num-runs] *numruns*]
  {[[-R | --reuse-dataset]]}
```

where **exp_file_paths** is a non-empty space separated list of experiment file paths, **numruns** is an integer specifying the number of times to repeat the profiling for each experiment file (default is 1), *-R* or *--reuse-dataset* option instructs xSpark to reuse (not to delete) application data in hdfs master node.

***log_profiling* command syntax:**

```
log_profiling {[[-L | --local]}}
```

where *-L* or *--local* option instructs xSpark use default local output folders.

***time_analysis* command syntax:**

```
time_analysis {[[-i | --input-dir] *dir*]}
```

where *dir* is the directory where the log files are located.

5.11.8 Example: Profile and Test CallsExample

- 1) Create the credential.json file as instructed above.
- 2) Configure the setup.json file as instructed above.
- 3) Configure the control.json file as instructed above.
- 4) Create and initialize a hdfs cluster with 5 nodes:

```
$ python3 xSpark\_dagsymb.py setup hdfs -n 5
```

- 5) Create and initialize a spark cluster with 5 nodes:

```
$ python3 xSpark\_dagsymb.py setup spark -n 5
```

- 6) Create *CallsExample_1.json* *CallsExample_2.json* file with the following contents:

```
CallsExample_1.json:
{
    "Deadline": 150000,
    "AppName": "CallsExample-1",
    "MetaProfileName": "CallsExample",
    "AppJar": "dagsymb/target/dagsymb-1.0-jar-with-dependencies.jar",
    "AppClass": "it.polimi.deepse.dagsymb.launchers.Launcher",
    "GuardEvaluatorClass": "it.polimi.deepse.dagsymb.examples.
                                GuardEvaluatorPromoCallsFile",
    "NumPartitions": 500,
    "DataMultiplier" : 1,
    "DeadlineMultiplier" : 1,
    "AppConf": {
        "0": {"Name": "threshold", "Value" : 1994},
        "1": {"Name": "minLocalLongCalls", "Value" : 1994},
        "2": {"Name": "minAbroadLongCalls", "Value" : 1994},
        "3": {"Name": "pastMonths", "Value" : 0},
        "4": {"Name": "last24HLocalCallsLength", "Value" : 1993},
        "5": {"Name": "last24HLocalCallsSize", "Value" : 1993},
        "6": {"Name": "last24HAbroadCallsLength", "Value" : 1993},
        "7": {"Name": "last24HAbroadCallsSize", "Value" : 1993},
        "8": {"Name": "MonthCallsLength", "Value" : 2990},
        "9": {"Name": "MonthCallsSize", "Value" : 3000}
    }
}

CallsExample_2.json:
{
    "Deadline": 150000,
    "AppName": "CallsExample-2",
    "MetaProfileName": "CallsExample",
    "AppJar": "dagsymb/target/dagsymb-1.0-jar-with-dependencies.jar",
    "AppClass": "it.polimi.deepse.dagsymb.launchers.Launcher",
    "GuardEvaluatorClass": "it.polimi.deepse.dagsymb.examples.
                                GuardEvaluatorPromoCallsFile",
    "NumPartitions": 500,
    "DataMultiplier" : 1,
    "DeadlineMultiplier" : 1,
    "AppConf": {
        "0": {"Name": "threshold", "Value" : 2600},
        "1": {"Name": "minLocalLongCalls", "Value" : 2600},
        "2": {"Name": "minAbroadLongCalls", "Value" : 2600},
    }
}
```

```
    "3": {"Name": "pastMonths", "Value" : 0},
    "4": {"Name": "last24HLocalCallsLength", "Value" : 2635},
    "5": {"Name": "last24HLocalCallsSize", "Value" : 2635},
    "6": {"Name": "last24HAbroadCallsLength", "Value" : 2635},
    "7": {"Name": "last24HAbroadCallsSize", "Value" : 2635},
    "8": {"Name": "MonthCallsLength", "Value" : 2990},
    "9": {"Name": "MonthCallsSize", "Value" : 3000}
  }
}
```

- 7) Run the Profiling with 5 iterations:

```
$ python3 xSpark\_dagsymb.py profile -r 5 CallsExample\_1.json CallsExample\_1
```

Run the Application Test with 5 iterations:

```
$ python3 xSpark\_dagsymb.py submit -r 5 CallsExample\_1.json CallsExample\_1
```

6

EVALUATION

IN this chapter we describe the experiments we have performed to evaluate the feasibility and quality of the solution and to address the research questions identified in Section 1.3. The evaluation took place by using *xSpark_{SEEPEP}*, which is the integration of *SEEPEP* with xSpark, to control the parallel execution of the Spark applications.

6.1 TEST ENVIRONMENT

The test environment was built on Virtual Machines (from here onward VMs) of type *Standard_D14_v2* [86] provided by Microsoft Azure [61], each of them equipped with 5 CPUs, 112 GiB of memory, 800 GiB of local SSD memory and 6000 Mbps of network bandwidth. This type of VM is optimized for memory usage, with a high memory/core ratio. The os and software packages installed on these VMs are: Canonical Ubuntu Server 14.04.5-LTS [82], Oracle Java 8 [49], Apache Hadoop 2.7.2 [6], Apache Spark 2.0.2 [7] and xSpark. All VM software is stored in a 200 GiB virtual hard drive maintained in the persistent layer of an Azure Blob Storage [24]. The VMs are organized in clusters composed of 5 VMs running HDFS (storing the input datasets) and 5 hosting Apache Spark and xSpark (5 for old xSpark and 5 for xSpark + *SEEPEP*). The Spark cluster runs Spark application either under Spark with the default configuration or under xSpark + *SEEPEP* or xSpark with the configuration parameters shown in Figure 6.1 or in Figure 6.2 respectively.

6.2 TESTED APPLICATIONS

We performed the experiments with two applications: PromoCalls and Louvain. PromoCalls is an example application that was developed at Politecnico di Milano in the Deib Labs¹. It resembles a batch application from a telecommunications company that calculates promotional discounts based on the number of daily domestic and foreign calls (calls longer than a parametric threshold) made by customers. If a customer makes more than min_l local long calls or more than min_a abroad long calls (or both) in a day, she may receive discounts on the calls made in that day, in the last m months, or in the current month. Only some or all of the discounts may be applied to the customer depending on the possible combinations of the trigger conditions.

¹ <https://github.com/seepep/promocalls>

```
{  
    "Alpha": 1.0,  
    "Beta": 0.3,  
    "OverScale": 2,  
    "K": 50,  
    "Ti": 12000,  
    "TSample": 500,  
    "Heuristic": "SYMEX_CONTROL_UNLIMITED",  
    "CoreQuantum": 0.05,  
    "CoreMin": 0,  
    "CpuPeriod": 100000  
}
```

Figure 6.1: control.json $xSpark_{SEEP}$ configuration parameters.

```
{  
    "Alpha": 1.0,  
    "Beta": 0.3,  
    "OverScale": 2,  
    "K": 50,  
    "Ti": 12000,  
    "TSample": 500,  
    "Heuristic": "CONTROL_UNLIMITED",  
    "CoreQuantum": 0.05,  
    "CoreMin": 0,  
    "CpuPeriod": 100000  
}
```

Figure 6.2: control.json $xSpark$ configuration parameters.

PromoCalls uses Spark to efficiently analyze the data of all calls and calculate the applicable discounts. PromoCalls was used as a reference application during the development of SEEPEP and for a preliminary assessment of the accuracy of the technique used.

Instead, to evaluate our approach to a real world application, we selected Louvain, a Spark implementation of the Louvain algorithm [19] that we downloaded from a highly ranked GitHub repository². Louvain uses *GraphX*, a Spark library specialized for graph processing, suitable for representing large user networks and analyzing communities belonging to these networks.

6.3 EXPERIMENTS

The experiments were performed on each of the applications subjected to the tests following the procedure described below, which consists of three distinct phases:

- 1 Execution of *SEEPEP* to obtain the path conditions and generate the launchers corresponding to each identified path.
- 2 Profiling the application subject to the test with the launchers generated and obtaining the plans for each path.
- 3 Use of tool to control the execution of the application by feeding it with input data sets larger by one order of magnitude compared to those used to generate the launchers.

We have generated at least one large data set for each profiled path. We also set a reasonable deadline, that is 20% longer than the minimum deadline, measured by running the application on Spark configured to use all the resources available on the clusters, and with the same data sets used in the experiments .

The datasets were randomly generated using *SEEPEP*.

The comparison of *xSpark_{SEEPEP}* with the original xSpark version has been done, for each application under test, identifying the best and worst cases, i.e. the paths with the lowest and highest number of stages³. In this way, we have quantified the error that xSpark can generate due to the fact that it ignores which plan was used in the profiling phase.

6.3.1 Results

Results produced by *xSpark_{SEEPEP}* for PromoCalls and Louvain tools, up to the profiling phase, are shown in Tables 6.1 and 6.2. Column

² <https://github.com/Sotera/spark-distributed-louvain-modularity>

³ In case two paths have the same number of stages, we chose the path with the shortest/longest execution time respectively for the best/worst case.

Path lists the paths found by *SEEPEP*: 8 unique paths were discovered in both cases.

Column Found? shows whether or not *SEEPEP* succeeded in generating a test case (and thus a corresponding profiling launcher) for the identified paths: it was successful in generating test cases for all path of *PromoCalls*, instead it was able to identify 6 out of 8 possible paths of *Louvain*⁴. As a matter of fact, we currently have no clue if *Louvain* can execute these program paths or not.)

Column Jobs and column Stages report the number of jobs and stages collected when profiling the launchers with *xSpark*, which range between 3 and 9 jobs, 3 and 9 stages for *Promocalls*, and 11 and 17 jobs, 73 and 364 stages for *Louvain*, respectively. These data are not available for the two paths of *Louvain* for which *SEEPEP* did not generate a launcher. In both tables, we marked with • and † the paths that correspond to the best and worst case, respectively, of each application.

The effectiveness of *xSpark_{SEEPEP}* to control the execution of the tested applications, whose inputs were fed with large datasets, is measured by the results of our experiments that are summarized in Tables 6.3 and 6.4. These tables also include the results obtained with the original version of *xSpark*, tuned on the worst and best case datasets mentioned above, and allow us to compare *xSpark_{SEEPEP}* against *xSpark*. The meaning of each column of the tables is explained here below.

For each profiled path P_i , column Experiment indicates the data obtained with *xSpark_{SEEPEP}* (*xSpark_s*), and *xSpark* configured with the worst-case dataset (*xSpark_w*) and with the best-case dataset (*xSpark_b*), respectively. In column deadline we show the set deadline in seconds. In column exec_time the actual execution time of the application is reported in seconds, as the average of 5 iterations of the experiments (for a total of 120 executions of *PromoCalls* (8 paths \times 5 repetitions \times 3 modes) and 90 executions of *Louvain*). In column violation we show the deadline violations (i.e., $exec_time > deadline$). The error is quantified in column error, defined as:

$$\text{error} = \frac{|deadline - exec_time|}{deadline} \cdot 100\%$$

that is, the percentage change vs. the deadline of the distance between the actual execution time and the deadline itself. In general the smaller error, the more efficient is the resource allocation, provided that the deadline is not violated, since less resources were used to meet the goal. On the other hand, if the deadline is violated, to smaller errors correspond shorter delays. Note that if the deadline were to be

⁴ The two paths of *Louvain* for which *SEEPEP* was not successful in identifying a corresponding test case were manually inspected. In any case, the proof that either these paths are infeasible, or an input datasets can be identified that exercise these paths, is missing.

Path	Found?	#Jobs	#Stages
0	Yes	6	6
1•	Yes	3	3
2	Yes	7	7
3	Yes	6	6
4	Yes	8	8
5	Yes	7	7
6†	Yes	9	9
7	Yes	8	8

Table 6.1: *PromoCalls* paths.

Path	Found?	#Jobs	#Stages
0	Yes	11	149
1	Yes	17	364
2†	Yes	17	364
3	Yes	11	149
4	No	-	-
5	No	-	-
6	Yes	17	364
7•	Yes	8	73

Table 6.2: *Louvain* paths.

Experiment	deadline[s]	exec_time[s]	violation	error	core_alloc[$\frac{\text{core}}{\text{s}}$]	penalty
xSpark _s	91.4	90.3	No	1.2%	41.3	—
P ₀ xSpark _w	91.4	88.0	No	3.8%	53.0	28.3%
xSpark _b	91.4	143.0	Yes	56.4%	30.6	∞
xSpark _s	56.4	46.0	No	18.4%	56.2	—
P ₁ xSpark _w	56.4	45.3	No	19.6%	56.5	0.5%
xSpark _b	56.4	55.3	No	1.9%	38.2	-33.6%
xSpark _s	107.8	106.3	No	1.3%	39.2	—
P ₂ xSpark _w	107.8	104.0	N	3.5%	52.1	32.9%
xSpark _b	107.8	175.0	Yes	62.4%	29.0	∞
xSpark _s	87.5	86.0	No	1.7%	42.4	—
P ₃ xSpark _w	87.5	83.0	No	5.1%	53.5	26.2%
xSpark _b	87.5	138.0	Yes	57.8%	30.1	∞
xSpark _s	147.6	146.0	No	1.1%	37.0	—
P ₄ xSpark _w	147.6	130.0	No	11.9%	51.4	38.9%
xSpark _b	147.6	228.0	Yes	54.5%	28.4	∞
xSpark _s	77.0	75.3	No	2.2%	41.0	—
P ₅ xSpark _w	77.0	70.0	No	9.1%	53.7	30.9%
xSpark _b	77.0	122.0	Yes	58.4%	29.7	∞
xSpark _s	122.2	120.3	No	1.5%	39.2	—
P ₆ xSpark _w	122.2	120.7	No	1.2%	43.6	11.3%
xSpark _b	122.2	204.0	Yes	67.0%	27.9	∞
xSpark _s	112.1	110.7	No	1.3%	39.2	—
P ₇ xSpark _w	112.1	100.0	No	10.8%	53.0	35.2%
xSpark _b	112.1	180.0	Yes	60.6%	28.8	∞

Table 6.3: Results for PromoCalls

Experiment	deadline[s]	exec_time[s]	violation	error	core_alloc[$\frac{\text{core}}{\text{s}}$]	penalty
xSpark _s	184.3	180.1	No	2.3%	35.9	—
P ₀ xSpark _w	184.3	142.0	No	23.0%	46.1	28.4%
xSpark _b	184.3	222.3	Yes	20.6%	15.2	∞
xSpark _s	228.0	227.0	No	0.4%	32.3	—
P ₁ xSpark _w	228.0	222.0	No	2.6%	33.2	2.8%
xSpark _b	228.0	329.3	Yes	44.4%	7.4	∞
xSpark _s	292.8	290.7	No	0.7%	32.3	—
P ₂ xSpark _w	292.8	289.0	No	1.3%	32.5	0.5%
xSpark _b	292.8	429.0	Yes	46.5%	7.0	∞
xSpark _s	228.7	226.0	No	1.2%	35.5	—
P ₃ xSpark _w	228.7	211.3	No	7.6%	41.4	16.6%
xSpark _b	228.7	292.0	Yes	27.7%	16.0	∞
xSpark _s	163.0	159.4	No	2.2%	38.4	—
P ₆ xSpark _w	163.0	158.0	No	3.0%	39.8	3.8%
xSpark _b	163.0	242.0	Yes	48.5%	8.5	∞
xSpark _s	156.0	139.0	No	10.9%	33.2	—
P ₇ xSpark _w	156.0	131.5	No	15.7%	43.6	31.4%
xSpark _b	156.0	152.7	No	2.1%	30.9	-7.0%

Table 6.4: Results for Louvain

considered strict, the penalty for a violation would be considered of infinite value [75]. In column `core_alloc` we show the average core allocation during the execution, that is defined as:

$$\text{core_alloc} = \frac{\sum_{s=0}^{\text{exec_time}} \text{coresAllocatedAtSecond}(s)}{\text{exec_time}}$$

We remark that the maximum value of `core_alloc` is 64 core/second since 64 is the number of cores provided by the cluster used for these experiments.

In the last column `penalty` we quantify the performance of $x\text{Spark}_{\text{SEEP}}_{\text{PEP}}$ compared to xSpark when executing the same experiment: `penalty` is defined as:

$$\text{penalty} = \begin{cases} \frac{\text{ru}_{\text{WORST}} - \text{ru}_{\text{BEST}}}{\text{ru}_{\text{SEEP}}_{\text{PEP}}} \cdot 100\%, & \text{if Violation} = \text{No} \\ \infty, & \text{if Violation} = \text{Yes} \end{cases}$$

Remarkably in our experiments, $x\text{Spark}_{\text{SEEP}}_{\text{PEP}}$ never violated the deadline, hence `penalty` measures the amount of resources that either $x\text{Spark}_w$ or $x\text{Spark}_b$ have used with respect to $x\text{Spark}_{\text{SEEP}}_{\text{PEP}}$. E.g. a penalty of 30% means that xSpark used 30% more resources than $x\text{Spark}_{\text{SEEP}}_{\text{PEP}}$. Instead, a negative penalty means that xSpark used less resources than $x\text{Spark}_{\text{SEEP}}_{\text{PEP}}$. Lastly, when the deadline is violated by either $x\text{Spark}_w$ or $x\text{Spark}_b$, we consider `penalty` to be infinite [75].

As evidenced by the data in Tables 6.3 and 6.4, $x\text{Spark}_b$ violated the deadline in 7 out of 8 paths in the experiments with PromoCalls, and 5 out of 6 paths in the experiments with Louvain. This is due to the optimistic, yet wrong, estimations made in the profiling phase. To explain this behaviour, we have to consider that, for Promocalls, $x\text{Spark}_b$ computes the local deadlines and the resource allocation as if the *PEP* always consisted of 3 stages. As a consequence, in all the experiments but P_1 (which is the actual best-case path) xSpark under allocates resources, resulting in the execution time eventually exceeding the deadline by as much as 51.9%. We measured the highest error displacement, equals to 67.0%, with the worst-case path P_6 where xSpark faces the widest gap between the profiling estimations and the actual runtime workload.

In the same tables we can appreciate that $x\text{Spark}_w$ does not violate any deadline, on the contrary it causes the earlier termination of the applications in most of the cases, with an error between 1.2% and 19.6% in the case of PromoCalls, and between 1.3% and 23% in the case of Louvain. This behaviour depends on the pessimistic, yet wrong, estimations made in the profiling phase, that mistakenly consider the worst-case path of the applications to represent all the possible paths. When considering paths P_1 , P_4 , P_5 and P_7 of PromoCalls, and P_0 , P_3 and P_7 of Louvain, the error is greater than 7%, leading to significantly sub-optimal resource allocations.

Remarkably, and on the contrary, $xSpark_{SEEP}$ does not violate any deadline and successfully provides an efficient resource allocation. The average error measured in our experiments is equal to 3.6% for PromoCalls, where for $xSpark_b$ and $xSpark_w$ is 52.4% and 8.1%, respectively, and equal to 2.9% for Louvain, where $xSpark_b$ and $xSpark_w$ have an average error of 31.6% and 8.9%, respectively.

The data in columns `core_alloc` confirm that $xSpark_{SEEP}$ outperforms the performance of $xSpark_b$ and $xSpark_w$. $xSpark_b$ underestimates allocated resources so as to make xSpark violate the deadlines in all experiments, except for path P_1 in PromoCalls and path P_7 in Louvain (the best cases). In the latter two cases, profiled data match the runtime workload, leading $xSpark_b$ to outperform both $xSpark_w$ and $xSpark_{SEEP}$, and minimizes the error and used resources.

$xSpark_{SEEP}$ allocates an average of 25.5% fewer resources in PromoCalls, and 13.9% in Louvain, with respect to $xSpark_w$.

From the results of our experiments we can evince a positive answer to both research questions RQ₁ and RQ₂, because $xSpark_{SEEP}$ effectively and efficiently controls the allocation of resources during the execution of PromoCalls and Louvain, keeping the execution times within considered deadlines with significantly smaller errors and consuming a lower amount of resources than the original version of xSpark.

6.4 THREATS TO VALIDITY

We came up with a considerable trial effort, which has led us to perform a total of 226 experiments on two different applications: a paradigmatic example and a real-world application taken from GitHub. We have shown that *SEEPEP* was able to find a test-case for 14 out 16 of the application paths statically identified with symbolic execution, and demonstrated how xSpark could take advantage from the integration with *SEEPEP*. In this section, we highlight the threats that may constrain the validity of our current results [90]:

Internal Threats. The test cases generated by *SEEPEP* have been slightly modified to increase the size of the datasets (without breaking the path conditions) for the execution of the experiments. This was done to ensure that we could test the desired paths and reliably obtain different repetitions of the experiments. However, data sets were not created in a totally random way. For this reason, we have preliminary executed some experiments with completely randomly created data sets and have obtained a similar result to the one presented. A broader set of experiments could be done to address this aspect as a matter for future developments.

External Threats. A limit to the generality of the experiments and of the solutions tested may derive from the assumption underlying the choice of the Spark applications that we have considered. We have

chosen two applications: one that uses Spark's core transformations and one that uses the *GraphX* entry for graph analysis. To increase the degree of generality of evaluations, an extension of the number and variety of Spark applications is desirable, taking into consideration different types, such as those that exploit machine learning solutions and use SQL.

Profiling as it is currently done constitutes an additional limitation of the approach used. In fact, a profiling session is required for each test case (launcher) found by *SEEPEP*, but this could become a problem when the number of program paths becomes too high. A desirable direction for a future work could be the improvement of this part of the tool chain. The improvement could be achieved through the adoption of a branch-based criterion to select profiling executions, replacing the current path-based criterion. The underpinning rationale is that since the execution of all the program branches assures that any possible *PEP* of the program has been profiled at least once; then, after profiling, we can map the profiled *PEPs* on the paths (and then the path conditions) that include the corresponding program branches. This optimization might also mitigate the issues with paths that our test generator fails to exercise, such as the 2 uncovered paths that we experienced in Louvain.

Concept and Conclusion Threats. The experiments demonstrate the validity of our idea, namely that the knowledge of the different *PEPs* generated by the Spark applications helps to analyze and control their performance and execution time. Furthermore, we show that an approach based simply on knowledge of the worst case may be sufficient to limit the number of deadline violations, but is significantly less efficient than the one we obtain with our solution. The results obtained are statistically robust and they show only a small variance.

CONCLUSION

IN the previous chapters we have presented the work done in order to support deadline-based QoS constrained multi-*PEP* Spark applications, i.e. applications whose execution flow cannot be represented with a single *Parallel Execution Plan (PEP)*, and whose actual execution flow is only known at runtime.

This chapter summarizes the conclusion of our work and the future works to improve the solution by extending the field of applicability and optimizing each individual components, to obtain a more complete and efficient solution.

7.1 CONCLUSION

The contribution provided by this thesis consists of the design and the development of *xSpark_{SEEPEP}*, a solution composed by an original and lightweight application of the principles of symbolic execution to detect the parallel plans of the execution of Spark multi-*PEP* applications, create the related execution profiles, and by the integration in *xSpark* of additional features that exploit the knowledge of the execution plans to manage the allocation of resources at runtime to help achieve the goal of maintaining the execution time of the application within a deadline specified by the user.

The evaluation results show that *xSpark_{SEEPEP}* meets all the expectations identified by the research questions formulated in Section 1.3, as it misses fewer deadlines and allocate resources more efficiently than *xSpark*.

7.2 FUTURE WORK

As the applicability of the profiling contained in *xSpark_{SEEPEP}* is currently subject to the limitation related to the number of paths to be profiled, that requires the profiling of the entire application by a specific launcher for each path identified by *SEEPEP*, and can become practically unfeasible if the number of paths is too high, a desirable future work could be directed at improving this part of the tool chain. This future work would be focused to moving away from the current profiling phase which uses path-based selection criteria in favour of a profiling that uses branch-based criteria.

Another future work could consist in the research along the lines of the applicability of the proposed solution to execute non-strict deadlines QoS-constrained multi-*PEP* applications.

BIBLIOGRAPHY

- [1] .
- [2] Omid Alipourfard, Hongqiang Harry Liu, Jianshu Chen, Shivaram Venkataraman, Minlan Yu, and Ming Zhang. “CherryPick: Adaptively Unearthing the Best Cloud Configurations for Big Data Analytics.” In: *14th Symposium on Networked Systems Design and Implementation*. USENIX Association, 2017 (cit. on p. 5).
- [3] *Apache Flink*. 2019 (cit. on p. 17).
- [4] *Apache Flink Architecture*. 2019 (cit. on p. 18).
- [5] *Apache Hadoop*. 2019 (cit. on pp. x, xiii).
- [6] *Apache Hadoop*. <http://hadoop.apache.org>. 2019 (cit. on p. 104).
- [7] *Apache Spark*. <http://spark.apache.org>. 2017 (cit. on p. 104).
- [8] *Apache Spark*. 2019 (cit. on pp. x, xiii, 2, 13).
- [9] *Apache Storm*. 2019 (cit. on p. 19).
- [10] *Apache Tez*. 2019 (cit. on p. 2).
- [11] Apache.org. *RDD Programming Guide - Spark 2.4.0 documentation*. <https://spark.apache.org/docs/latest/rdd-programming-guide.html>. Accessed: 2019-02-20.
- [12] Anonymous Author(s). “Symbolic Execution-driven Extraction of the Parallel Execution Plans of Spark Applications.” In: *Proceedings of The 27th ACM Joint European Software Engineering Conference and Symposium on the Foundations of Software Engineering (ESEC/FSE 2019)*. (2019), pp. 1–11 (cit. on pp. 6, 59, 60, 62).
- [13] AWS. 2019 (cit. on p. 2).
- [14] Roberto Baldoni, Emilio Coppa, Daniele Cono D’elia, Camil Demetrescu, and Irene Finocchi. “A Survey of Symbolic Execution Techniques.” In: *ACM Comput. Surv.* 51.3 (May 2018), 50:1–50:39. ISSN: 0360-0300 (cit. on pp. 38–45).
- [15] L. Baresi and G. Quattrocchi. “Towards Vertically Scalable Spark Applications.” In: *Proc. of Euro-Par 2018: Parallel Processing Workshops*. Springer, 2018 (cit. on pp. 5, 58, 59).
- [16] L. Baresi, S. Guinea, A. Leva, and G. Quattrocchi. *Fine-grained Dynamic Resource Allocation for Big-Data Applications*. Tech. rep. 2018 (cit. on pp. 5, 58, 59).

- [17] Clark Barrett, Daniel Kroening, and Thomas Melham. *Problem solving for the 21st century: Efficient solver for satisfiability modulo theories*. Knowledge Transfer Report, Technical Report 3. London Mathematical Society, Smith Institute for Industrial Mathematics, and System Engineering, 2014 (cit. on pp. 39, 41).
- [18] *Big Data Companies* 2019. 2019 (cit. on p. 2).
- [19] Vincent D Blondel, Jean-Loup Guillaume, Renaud Lambiotte, and Etienne Lefebvre. “Fast Unfolding of Communities in Large Networks.” In: *Journal of statistical mechanics: theory and experiment* (2008) (cit. on pp. 7, 106).
- [20] Robert S. Boyer, Bernard Elspas, and Karl N. Levitt. “SELECT – A Formal System for Testing and Debugging Programs by Symbolic Execution.” In: *Proc. of Int. Conf. on Reliable Software*. Los Angeles, California: ACM, 1975, pp. 234–245 (cit. on p. 38).
- [21] Pietro Braione, Giovanni Denaro, and Mauro Pezzè. “JBSE: a symbolic executor for Java programs with complex heap inputs.” In: *Proceedings of the 24th ACM SIGSOFT International Symposium on Foundations of Software Engineering, FSE 2016*. 2016, pp. 1018–1022.
- [22] Pietro Braione, Giovanni Denaro, Andrea Mattavelli, and Mauro Pezzè. “Combining symbolic execution and search-based testing for programs with complex heap inputs.” In: *Proceedings of the 26th ACM SIGSOFT International Symposium on Software Testing and Analysis*. 2017, pp. 90–101.
- [23] Pietro Braione, Giovanni Denaro, Andrea Mattavelli, and Mauro Pezzè. “SUSHI: a test generator for programs with complex structured inputs.” In: *Proceedings of the 40th International Conference on Software Engineering: Companion Proceedings, ICSE 2018*. 2018, pp. 21–24.
- [24] Cristian Cadar, Daniel Dunbar, and Dawson R. Engler. “KLEE: Unassisted and Automatic Generation of High-coverage Tests for Complex Systems Programs.” In: *Proc. 8th USENIX Conf. on Operating Systems Design and Implementation*. OSDI’08. San Diego, California: USENIX Association, 2008, pp. 209–224 (cit. on p. 46).
- [25] Cristian Cadar and Koushik Sen. “Symbolic Execution for Software Testing: Three Decades Later.” In: *Communications of the ACM* 56.2 (2013), pp. 82–90. ISSN: 0001-0782 (cit. on pp. 42, 46).
- [26] Satish Chandra, Stephen J. Fink, and Manu Sridharan. “Snuggliebug: A Powerful Approach to Weakest Preconditions.” In: *Proc. 30th ACM SIGPLAN Conf. on Programming Language Design and Implementation*. PLDI’09. Dublin, Ireland: ACM, 2009, pp. 363–374. ISBN: 978-1-60558-392-1 (cit. on p. 47).

- [27] D. Cheng, J. Rao, C. Jiang, and X. Zhou. "Resource and Deadline-Aware Job Scheduling in Dynamic Hadoop Clusters." In: *Parallel and Distributed Processing Symposium (IPDPS), 2015 IEEE International*. Institute of Electrical and Electronics Engineers (IEEE), 2015 (cit. on pp. xi, xiv).
- [28] Vitaly Chipounov, Volodymyr Kuznetsov, and George Candea. "The S2E Platform: Design, Implementation, and Applications." In: *ACM Transactions on Computer Systems (TOCS)*. TOCS 2011 30.1 (2012), 2:1–2:49 (cit. on p. 45).
- [29] Peter Dinges and Gul Agha. "Targeted Test Input Generation Using Symbolic-concrete Backward Execution." In: *Proc. 29th ACM/IEEE Int. Conf. on Automated Software Engineering*. ASE'14. Västerås, Sweden, 2014, pp. 31–36. ISBN: 978-1-4503-3013-8 (cit. on p. 47).
- [30] *Docker Documentation*. 2019 (cit. on p. 5).
- [31] Schahram Dustdar, Yike Guo, Benjamin Satzger, and Hong-Linh Truong. "Principles of Elastic Processes." In: *IEEE Internet Computing* 15.5 (2011), pp. 66–71 (cit. on pp. 22, 23).
- [32] EMC. 2019 (cit. on p. 2).
- [33] Gordon Fraser and Andrea Arcuri. "Whole Test Suite Generation." In: *IEEE Transactions on Software Engineering* 39.2 (2013), pp. 276–291.
- [34] Gartner. <https://www.gartner.com/en>. 2019 (cit. on p. 1).
- [35] Giovanni Paolo Gibilisco et al. "Stage Aware Performance Modeling of DAG Based in Memory Analytic Platforms." In: *Proc. of IEEE 9th International Conference on Cloud Computing*. IEEE. 2016 (cit. on p. 58).
- [36] Patrice Godefroid, Nils Klarlund, and Koushik Sen. "DART: Directed Automated Random Testing." In: *Proc. ACM SIGPLAN Conf. on Programming Language Design and Implementation*. PLDI'05. Chicago, IL, USA, 2005, pp. 213–223. ISBN: 1-59593-056-6 (cit. on pp. 43, 44).
- [37] Patrice Godefroid, Michael Y. Levin, and David A. Molnar. "Automated Whitebox Fuzz Testing." In: *Proc. Network and Distributed System Security Symp.* NDSS'08. 2008 (cit. on p. 44).
- [38] Patrice Godefroid, Michael Y. Levin, and David A. Molnar. "SAGE: Whitebox Fuzzing for Security Testing." In: *Queue* 10.1 (2012), 20:20–20:27. ISSN: 1542-7730 (cit. on p. 39).
- [39] Google. 2019 (cit. on p. 2).
- [40] Google Map-Reduce. 2019 (cit. on pp. x, xiii, 2).
- [41] Nikolas Roman Herbst, Samuel Kounev, and Ralf H Reussner. "Elasticity in Cloud Computing: What It Is, and What It Is Not." In: *ICAC*. 2013 (cit. on p. 22).

- [42] Benjamin Hindman, Andy Konwinski, Matei Zaharia, Ali Ghodsi, Anthony D. Joseph, Randy Katz, Scott Shenker, and Ion Stoica. "Mesos: A Platform for Fine-grained Resource Sharing in the Data Center." In: *Proc. of the 8th USENIX Conference on Networked Systems Design and Implementation*. NSDI. USENIX, 2011 (cit. on pp. xi, xiv).
- [43] William E. Howden. "Symbolic Testing and the DISSECT Symbolic Evaluation System." In: *IEEE Transactions on Software Engineering (TSE)* 3.4 (1977), pp. 266–278. ISSN: 0098-5589 (cit. on p. 38).
- [44] HPE. 2019 (cit. on p. 2).
- [45] IBM. 2019 (cit. on p. 2).
- [46] IDC. 2019 (cit. on p. 9).
- [47] Internet Live Statistics. 2019 (cit. on p. 1).
- [48] M. Islam et al. "dSpark: Deadline-based Resource Allocation for Big Data Applications in Apache Spark." In: *Proc. of the 13th IEEE International Conference on eScience*. 2017 (cit. on pp. 4, 58).
- [49] Java 8. <http://www.oracle.com/technetwork/java/javase/overview/java8-2100321.html..> 2019 (cit. on p. 104).
- [50] K. Kc and K. Anyanwu. "Scheduling Hadoop Jobs to Meet Deadlines." In: *Proc. of the IEEE 2nd International Conference on Cloud Computing Technology and Science*. IEEE, 2010.
- [51] Jeffrey O Kephart and David M Chess. "The Vision of Autonomic Computing." In: *Computer* 36 (2003) (cit. on p. 22).
- [52] James C. King. "A New Approach to Program Testing." In: *Proc. Int. Conf. on Reliable Software*. Los Angeles, California: ACM, 1975, pp. 228–233 (cit. on p. 38).
- [53] James C. King. "Symbolic Execution and Program Testing." In: *Communications of the ACM* 19.7 (1976), pp. 385–394. ISSN: 0001-0782 (cit. on p. 38).
- [54] Rogério de Lemos, Holger Giese, Hausi A. Müller, Mary Shaw, Jesper Andersson, Marin Litoiu, Bradley Schmerl, Gabriel Tamura, Norha M. Villegas, Thomas Vogel, Danny Weyns, Luciano Baresi, Basil Becker, Nelly Bencomo, Yuriy Brun, Bojan Cukic, Ron Desmarais, Schahram Dustdar, Gregor Engels, Kurt Geihs, Karl M. Göschka, Alessandra Gorla, Vincenzo Grassi, Paola Inverardi, Gabor Karsai, Jeff Kramer, Antónia Lopes, Jeff Magee, Sam Malek, Serge Mankovskii, Raffaela Mirandola, John Mylopoulos, Oscar Nierstrasz, Mauro Pezzè, Christian Prehofer, Wilhelm Schäfer, Rick Schlichting, Dennis B. Smith, João Pedro Sousa, Ladan Tahvildari, Kenny Wong, and Jochen Wuttke. "Software Engineering for Self-Adaptive Systems: A Second Research Roadmap." In: *Software Engineering for Self-Adaptive*

- Systems II: International Seminar, Dagstuhl Castle, Germany, October 24-29, 2010 Revised Selected and Invited Papers.* Ed. by Rogério de Lemos, Holger Giese, Hausi A. Müller, and Mary Shaw. Springer Berlin Heidelberg, 2013 (cit. on p. 21).
- [55] Kin-Keung Ma, Khoo Yit Phang, Jeffrey S. Foster, and Michael Hicks. "Directed Symbolic Execution." In: *Proc. 18th Int. Conf. on Static Analysis*. SAS'11. Venice, Italy, 2011, pp. 95–111. ISBN: 978-3-642-23701-0 (cit. on p. 47).
 - [56] Francesco Marconi, Giovanni Quattrocchi, Luciano Baresi, Marcello Bersani, and Matteo Rossi. "On the Timed Analysis of Big-Data Applications." In: *Proc. of the 10th NASA Formal Methods Symposium*. Springer, 2018 (cit. on p. 58).
 - [57] Patrick Wendell et al Matei Zaharia Reynold S Xin. "Apache Spark: A unified engine for big data processing." In: *Communications of the ACM* 59 (Nov. 2016), pp. 56–65 (cit. on p. 3).
 - [58] Olga Shumsky Matlin, Ewing L. Lusk, and William McCune. "SPINning Parallel Systems Software." In: *Model Checking of Software, 9th International SPIN Workshop*. 2002, pp. 213–220.
 - [59] Phil McMinn. "Search-based Software Test Data Generation: A Survey." In: *Software Testing, Verification & Reliability* 14.2 (2004), pp. 105–156. ISSN: 0960-0833 (cit. on p. 46).
 - [60] Microsoft Azure. <https://azure.microsoft.com/en-us/>. 2019 (cit. on p. 91).
 - [61] Microsoft Azure Virtual Machines. <https://azure.microsoft.com/en-us/services/virtual-machines/>. 2019 (cit. on p. 104).
 - [62] Microsoft. 2019 (cit. on p. 2).
 - [63] Leonardo Mendonça de Moura and Nikolaj Bjørner. "Z₃: An Efficient SMT Solver." In: *Tools and Algorithms for the Construction and Analysis of Systems, TACAS 2008*. 2008, pp. 337–340.
 - [64] Andrew Or. *Understanding your Apache Spark Application Through Visualization*. 2019 (cit. on p. 16).
 - [65] Oracle. 2019 (cit. on p. 2).
 - [66] Salman Pervez, Ganesh Gopalakrishnan, Robert M. Kirby, Rajeev Thakur, and William Gropp. "Formal Methods Applied to High-performance Computing Software Design: A Case Study of MPI One-sided Communication-based Locking." In: *Software Practice and Experience* 40.1 (Jan. 2010), pp. 23–43.
 - [67] J. Polo et al. "Performance-driven Task Co-scheduling for MapReduce Environments." In: *Proc. of the 20th IEEE Network Operations and Management Symposium*. IEEE, 2010.

- [68] Poornima Purohit, DR Apoorva, and et al PV Lathashree. "Big data analytics in the cloud: Spark on Hadoop vs MPI/OpenMP on Beowulf." In: *Procedia Computer Science* 53 (Dec. 2015), pp. 121–130 (cit. on p. 3).
- [69] Veselin Raychev, Madanlal Musuvathi, and Todd Mytkowicz. "Parallelizing User-defined Aggregations Using Symbolic Execution." In: *Proceedings of the 25th Symposium on Operating Systems Principles*. SOSP '15. ACM, 2015, pp. 153–167.
- [70] *Running Spark on Mesos*. 2019 (cit. on p. 31).
- [71] *Running Spark on YARN*. 2019 (cit. on p. 31).
- [72] Mythreyee S, Poornima Purohit, Apoorva D.R., Harshitha R, and Lathashree P.V. "A Study on Use of Big Data in Cloud Computing Environment." In: *International Journal of Advance Research, Ideas and Innovations in Technology* 3 (2017), pp. 1312–1318 (cit. on pp. x, xiii, 2, 21).
- [73] *SAP*. 2019 (cit. on p. 2).
- [74] Anuj Sehgal, Vladislav Perelman, Siarhei Kuryla, and Jürgen Schönwälder. "Management of Resource Constrained Devices in the Internet of Things." In: *Communications Magazine, IEEE* 50 (2012) (cit. on p. 22).
- [75] Kang G Shin and Parameswaran Ramanathan. "Real-time Computing: A New Discipline of Computer Science and Engineering." In: *Proceedings of the IEEE* 82 (1994) (cit. on p. 111).
- [76] Yan Shoshitaishvili, Ruoyu Wang, Christopher Salls, Nick Stephens, Mario Polino, Andrew Dutcher, John Grosen, Siji Feng, Christophe Hauser, Christopher Kruegel, and Giovanni Vigna. "SOK: (State of) The Art of War: Offensive Techniques in Binary Analysis." In: *IEEE Symp. on Security and Privacy*. SP'16. 2016, pp. 138–157 (cit. on p. 39).
- [77] S. Sidhanta, W. Golab, and S. Mukhopadhyay. "OptEx: A Deadline-Aware Cost Optimization Model for Spark." In: *16th IEEE/ACM International Symposium on Cluster, Cloud and Grid Computing*. 2016 (cit. on pp. 5, 58).
- [78] Stephen F. Siegel and Louis F. Rossi. "Analyzing BlobFlow: A Case Study Using Model Checking to Verify Parallel Scientific Software." In: *Recent Advances in Parallel Virtual Machine and Message Passing Interface, 15th European PVM/MPI Users' Group Meeting*. 2008, pp. 274–282.
- [79] Stephen F. Siegel and Timothy K. Zirkel. "Loop Invariant Symbolic Execution for Parallel Programs." In: *13th International Conference on Verification, Model Checking and Abstract Interpretation*. 2012, pp. 412–427.

- [80] Stephen F. Siegel, Anastasia Mironova, George S. Avrunin, and Lori A. Clarke. "Combining Symbolic Execution with Model Checking to Verify Parallel Numerical Programs." In: *ACM Transactions on Software Engineering and Methodology* 17.2 (May 2008), 10:1–10:34.
- [81] *Teradata*. 2019 (cit. on p. 2).
- [82] *Ubuntu 14.04.5 LTS (Trusty Tahr)*. <http://releases.ubuntu.com/14.04/>. 2019 (cit. on p. 104).
- [83] Vinod Kumar Vavilapalli, Arun C. Murthy, Chris Douglas, Sharad Agarwal, Mahadev Konar, Robert Evans, Thomas Graves, Jason Lowe, Hitesh Shah, Siddharth Seth, Bikas Saha, Carlo Curino, Owen O'Malley, Sanjay Radia, Benjamin Reed, and Eric Baldeschwieler. "Apache Hadoop YARN: Yet Another Resource Negotiator." In: *Proceedings of the 4th Annual Symposium on Cloud Computing*. SOCC '13. Santa Clara, California: ACM, 2013, 5:1–5:16. ISBN: 978-1-4503-2428-1 (cit. on p. 26).
- [84] Abhishek Verma, Ludmila Cherkasova, and Roy H. Campbell. "ARIA: Automatic Resource Inference and Allocation for Mapreduce Environments." In: *Proc. of the 8th ACM International Conference on Autonomic Computing*. ACM, 2011 (cit. on pp. xi, xiv).
- [85] J. Verne. *Journey to the Center of the Earth*. Classics illustrated. Huge Print Press, 1957. ISBN: 9780758311993.
- [86] *Virtual machine sizes for Azure Cloud services*. <https://docs.microsoft.com/en-us/azure/cloud-services/cloudservices-sizes-specs#dv2-series>. 2019 (cit. on p. 104).
- [87] *VMware*. 2019 (cit. on p. 2).
- [88] Danny Weyns, M. Usman Iftikhar, Sam Malek, and Jesper Andersson. "Claims and Supporting Evidence for Self-adaptive Systems: A Literature Study." In: *Proc. of the 7th International Symposium on Software Engineering for Adaptive and Self-Managing Systems*. SEAMS '12. Zurich, Switzerland: IEEE Press, 2012. ISBN: 978-1-4673-1787-0 (cit. on p. 21).
- [89] *What is Big Data*. <https://www.oracle.com/big-data/guide/what-is-big-data.html>. 2019 (cit. on p. 1).
- [90] Claes Wohlin et al. "Empirical Research Methods in Web and Software Engineering." In: *Web Engineering* (2006) (cit. on p. 112).
- [91] Tao Xie, Nikolai Tillmann, Jonathan de Halleux, and Wolfram Schulte. "Fitness-guided path exploration in dynamic symbolic execution." In: *Proc. 2009 IEEE/IFIP Int. Conf. on Dependable Systems and Networks*. DSN'09. 2009, pp. 359–368 (cit. on p. 46).

- [92] Matei Zaharia, Dhruba Borthakur, Joydeep Sen Sarma, Khaled Elmeleegy, Scott Shenker, and Ion Stoica. "Delay Scheduling: A Simple Technique for Achieving Locality and Fairness in Cluster Scheduling." In: *Proc. of the 5th European Conference on Computer Systems*. EuroSys '10. Paris, France: ACM, 2010. ISBN: 978-1-60558-577-2 (cit. on p. 3).
- [93] Qi Zhang, Lu Cheng, and Raouf Boutaba. "Cloud computing: State of the Art and Research Challenges." In: *Journal of internet services and applications* 1 (2010) (cit. on p. 22).