

Article

Exploiting Machine Learning For Improving In-memory Execution of Data-intensive Workflows on Parallel Machines

Riccardo Cantini¹ , Fabrizio Marozzo¹ , Alessio Orsino¹ , Domenico Talia^{1,*}  and Paolo Trunfio¹ 

¹ DIMES, University of Calabria, Rende, Italy

* Correspondence: talia@dimes.unical.it

Version April 29, 2021 submitted to Future Internet

Abstract: Workflows are largely used to orchestrate complex sets of operations required to handle and process huge amounts of data. Parallel processing is often vital to reduce execution time when complex data-intensive workflows must be run efficiently, and at the same time in-memory processing can bring important benefits to accelerate execution. However, optimization techniques are necessary to fully exploit in-memory processing avoiding performance drops due to memory saturation events. This paper proposes a novel solution, called Intelligent In-memory Workflow Manager (IIWM), for optimizing the in-memory execution of data-intensive workflows on parallel machines. IIWM is based on two complementary strategies: 1) a machine learning strategy for predicting memory occupancy and execution time of workflow tasks; 2) a scheduling strategy that allocates tasks to a computing node taking into account the (predicted) memory occupancy and execution time of each task, and the memory available on that node. The effectiveness of the machine learning-based predictor and the scheduling strategy are demonstrated experimentally using as a testbed Spark, a high-performance Big Data processing framework that exploits in-memory computing to speed up execution of large-scale applications. In particular, two synthetic workflows have been prepared for testing the robustness of IIWM in scenarios characterized by a high level of parallelism and a limited amount of memory reserved for execution. Furthermore, a real data analysis workflow has been used as a case study, for better assessing the benefits of the proposed approach. Thanks to high accuracy in predicting resources used at runtime, IIWM was able to avoid disk writes caused by memory saturation, outperforming a traditional strategy in which only dependencies among tasks are taken into account. Specifically, IIWM achieved up to 31% and 40% reduction of makespan and a performance improvement up to 1.45x and 1.66x on the synthetic workflows and the real case study respectively.

Keywords: Workflow, Data-intensive, In-memory, Machine Learning, Apache Spark, Scheduling.

1. Introduction

A data-intensive workflow is the description of a process that usually involves a set of computational steps implementing complex scientific functions, such as data acquisition, transformation, analysis, storage, and visualization [1]. Parallelism can be achieved by concurrently executing independent tasks by trying to make use of all computing nodes, even if, in many cases, it is necessary to execute multiple tasks on the same computing node [2]. For example, this occurs when the number of tasks is greater than the number of available nodes, or because multiple tasks use a dataset located on the same node. These scenarios are prone to memory saturation and moving data to disk may result in higher execution times, which leads to the need for a scheduling strategy able to cope with this issue [3,4].

34 In most cases, distributed processing systems use a-priori policies for handling task execution and
35 data management. For example, in the MapReduce programming model used by *Hadoop*, mappers
36 write intermediate results after each computation so performing disk-based processing with partial use
37 of memory [5] through the exploitation of the Hadoop Distributed File System (HDFS). On the other
38 hand, *Apache Spark*¹, that is a state-of-the-art data analysis framework for large-scale data processing
39 exploiting in-memory computing, relies on a Directed Acyclic Graph (DAG) paradigm and is based
40 on: *i*) an abstraction for data collections which enables parallel execution and fault-tolerance, named
41 Resilient Distributed Datasets (RDDs) [6]; *ii*) a DAG engine, that manages the execution of jobs,
42 stages and tasks. Besides, it provides different storage levels for data caching and persistence, while
43 performing in-memory computing with partial use of the disk. The Spark in-memory approach is
44 generally more efficient, but a time overhead may be caused by spilling data from memory to disk
45 when memory usage exceeds a given threshold [7]. This overhead can be significantly reduced if
46 memory occupancy of a task is known in advance, to avoid running in parallel two or more tasks that
47 cumulatively exceed the available memory, thus causing data spilling. For this reason, memory is
48 considered a key factor for performance and stability of Spark jobs and Out-of-Memory (OOM) errors
49 are often hard to fix. Recent efforts have been oriented towards developing prediction models for the
50 performance estimation of Big Data applications, although most of the approaches rely on analytical
51 models and only a few recent studies have investigated the use of supervised machine learning models
52 [8–10].

53 In this work we propose a system, named *Intelligent In-memory Workflow Manager* (IIWM), specially
54 designed for improving application performance through intelligent usage of memory resources. This
55 is done by identifying clusters of tasks that can be executed in parallel on the same node, optimizing
56 in-memory processing so avoiding the use of disk storage. Given a data-intensive workflow, IIWM
57 exploits a regression model for estimating the amount of memory occupied by each workflow task
58 and its execution time. This model is trained on a log of past executed workflows, represented in a
59 transactional way through a set of relevant features that characterize the considered workflow, such as:

- 60 • Workflow structure, in terms of tasks and data dependencies.
- 61 • Input format, such as the number of rows, dimensionality, and all other features required to
62 describe the complexity of input data.
- 63 • Type of the tasks, i.e., the computation performed by a given node of the workflow. For
64 example, in the case of data analysis workflows, we can distinguish among supervised learning,
65 unsupervised learning, and association rules discovery tasks, and also between learning and
66 prediction tasks.

67 Predictions made for a given computing node are applicable to all computing nodes of the same
68 type (i.e., having the same architecture, processor type, operating systems, memory resources), which
69 makes the proposed approach effectively usable on large-scale homogeneous HPC systems composed
70 of many identical servers. Given a data-intensive workflow, IIWM exploits the estimates coming
71 from the machine learning model for producing a scheduling plan aimed at reducing (and, in most
72 cases, avoiding) main memory saturation events, which may happen when multiple tasks are executed
73 concurrently on the same computing node. This leads to the improvement of application performance,
74 as swapping or spilling to disk caused by main memory saturation may result in significant time
75 overhead, which can be particularly costly when running workflows involving very large datasets
76 and/or complex tasks.

77 IIWM has been experimentally evaluated using as a testbed Spark, which is expected to become the
78 most adopted Big Data engine in the next few years [11]. In particular, we assessed the benefits coming
79 from the use of IIWM by executing two synthetic workflows specially generated for investigating

¹ <https://spark.apache.org/>

specific scenarios related to the presence of a high level of parallelism and a limited amount of memory reserved for execution. The effectiveness of the proposed approach has been further confirmed through the execution of a real data mining workflow as a case study. We carried out an in-depth comparison between IWM and a traditional blind scheduling strategy, which only considers workflow dependencies for the parallel execution of tasks. The proposed approach showed to be the most suitable solution in all evaluated scenarios outperforming the blind strategy thanks to high accuracy in predicting resources used at runtime, which leads to the minimization of disk writes caused by memory saturation.

The rest of the paper is organized as follows. Section 2 discusses related work. Section 3 describes the proposed system. Section 4 presents and discusses the experimental results. Section 5 concludes the paper.

1.1. Problem statement

The problem addressed in this study consists in the optimization of the in-memory execution of data-intensive workflows, evaluated in terms of makespan (i.e., the total time required to process all given tasks), and application performance. The main reason behind the drop in performance in such workflows is related to the necessity of swapping/spilling data to disk when memory saturation events occur. To cope with this issue, we propose an effective way of scheduling a workflow that minimizes the probability of memory saturation, while maximizes in-memory computing and thus performance.

A workflow \mathcal{W} can be represented using a DAG, described by a set of tasks $\mathcal{T} = \{t_1, t_2, \dots, t_n\}$ (i.e., vertices) and dependencies among them $\mathcal{A} \subseteq (\mathcal{T} \times \mathcal{T}) = \{a_1, \dots, a_m\}: a_i = (t_i, t_j), t_i \in \mathcal{T}, t_j \in \mathcal{T}$ (i.e., directed edges). Specifically, data dependencies (i.e., all the input data of a task have already been made available) have to be considered rather than control dependencies (i.e., all predecessors of a task must be terminated before it can be executed), as we refer to data-intensive workflows [12].

Formally, given a set of q computing resources $R = \{r_1, \dots, r_q\}$, workflow scheduling can be defined as the mapping $\mathcal{T} \rightarrow R$ from each task $t \in \mathcal{T}$ to a resource $r \in R$, so as to meet a set of specified constraints which influence the choice of an appropriate scheduling strategy [13]. Workflow scheduling techniques are often aimed at optimizing several factors, including makespan and overall cost that in turn depend on data transfer and compute cost [14]. In this study, a multi-objective optimization has been applied, jointly minimizing execution time and memory saturation. This is achieved by using a scheduling strategy that exploits a regression model aimed at predicting the behavior of a given workflow, in terms of resource demand and execution time (see Section 3). For the Reader's convenience, Table 1 shows the meaning of the main symbols used in the paper.

Symbol	Meaning
$\mathcal{T} = \{t_1, t_2, \dots, t_n\}$	Set of tasks.
$\mathcal{A} \subseteq (\mathcal{T} \times \mathcal{T}) = \{a_1, \dots, a_m\}$	Dependencies. $a_i = (t_i, t_j), t_i \in \mathcal{T}, t_j \in \mathcal{T}$.
d_t	Description of the dataset processed by task t .
$\mathcal{W} = (\mathcal{T}, \mathcal{A})$	Workflow.
$\mathcal{N}^{in}(t) = \{t' \in \mathcal{T} \mid (t', t) \in \mathcal{A}\}$	In-neighbourhood of task t .
$\mathcal{N}^{out}(t) = \{t' \in \mathcal{T} \mid (t, t') \in \mathcal{A}\}$	Out-neighbourhood of task t .
\mathcal{M}	Regression prediction model.
$\mathcal{S} = \langle s_1, \dots, s_k \rangle$	List of stages. $s_i \subseteq \mathcal{T} \mid (t_x \parallel t_y) \forall t_x, t_y \in s_i$.
C	Maximum amount of memory available for a computing node.
$C_s = C - \sum_{t \in s} \mathcal{M}.predict_mem(t, d_t)$	Residual capacity of a stage s .

Table 1. Meaning of the main symbols.

113 2. Related work

114 Recent studies have shown the effectiveness of machine learning-based prediction modelling in
115 supporting code optimization, parallelism mapping, task scheduling, and processor resource allocation
116 [10]. Moreover, predicting running times and memory footprint is important for estimating the cost of
117 execution and better managing resources at runtime [11]. For instance, in-memory data processing
118 frameworks like Spark can benefit from informed co-location of tasks [10]. In fact, if too many
119 applications or tasks are assigned to a computing node, such that the memory used on the host exceeds
120 the available one, memory paging to disk (i.e., swapping), data spilling to disk in Spark, or OOM
121 errors can occur with consequential drops of performance.

122 Our work focuses on improving the performance of a Spark application using machine
123 learning-based techniques. The challenge is to effectively schedule tasks in a data-intensive workflow
124 for improving resource usage and application performance, by inferring the resource demand of each
125 task, in terms of memory occupancy and time.

126 State-of-the-art techniques aimed at improving the performance of data-intensive applications
127 can be divided into two main categories: *analytical-based* and *machine learning-based*. For each category,
128 the main proposed solutions and their differences with respect to our technique are discussed.

129 2.1. Analytical-based

130 Techniques in this category use information collected at runtime and statistics in order to tune a
131 Spark application, improving its performance as follows:

- 132 • Choosing the serialization strategy for caching RDDs in RAM, based on previous statistics
133 collected on different working sets, such as memory footprint, CPU usage, RDDs size,
134 serialization costs, etc. [15,16].
- 135 • Dynamically adapting resources to data storage, using a feedback-based mechanism with
136 real-time monitoring of memory usage of the application [17].
- 137 • Scheduling jobs by dynamically adjusting concurrency through a feedback-based strategy. Taking
138 into account memory usage via garbage collection, network I/O, and Spark RDDs lineage
139 information, it is possible to choose the number of tasks to assign to an executor [18,19].

140 The aforementioned works use different strategies to improve in-memory computing of Spark that
141 exploit static or dynamic techniques able to introduce some information in the choice of configuration
142 parameters. However, no prediction models are employed and this may lead to unpredicted behaviors.
143 IIWM, instead, uses a prediction regression model to estimate a set of information about a running
144 Spark application, exploiting it to optimize in-memory execution. Moreover, unlike real time adapting
145 strategies, which use a feedback-based mechanism by continuously monitoring the execution, the
146 IIWM model is trained offline, achieving fast and accurate predictions while used for inferring the
147 resource demand of each task in a given workflow.

148 2.2. Machine learning-based

149 These techniques are based on the development of learning models for predicting performance
150 (mainly memory occupancy and execution time) of a large set of different applications in several
151 scenarios, on the basis of prior knowledge. This enables the adoption of a *performance-centric* approach
152 [8], based on an informed performance improvement, which can be beneficial for the execution of
153 data-intensive applications, especially in the context of HPC systems.

154 Several techniques use collaborative filtering to identify how well an application will run on a
155 computing node. For instance, *Quasar* [8] uses classification techniques based on collaborative filtering
156 to determine the characteristics of the running application in allocating resources and assigning tasks.
157 When submitted, a new application is shortly profiled and the collected information is combined with
158 the classification engine, based on previous workloads, to support a greedy scheduling policy that
159 improves throughput. Application is monitored throughout the execution to adjust resource allocation

160 and assignment if required, using a single model for the estimation. Adapting this technique to Spark
161 can help to assign tasks to computing nodes within the memory constraints and avoid exceeding the
162 capacity, thus causing spilling of data to disk. Another approach based on collaborative filtering has
163 been proposed by Lull et al. [9]. In this case, the task co-location problem is modelled as a cooperative
164 game and a game-theoretic framework, namely *Cooper*, is proposed for improving resource usage. The
165 algorithm builds pairwise coalitions as stable marriages to assign an additional task to a host based
166 on its available memory, and the Spark default scheduler is adopted to assign tasks. In particular, a
167 predictor receives performance information collected offline and estimates which co-runner is better,
168 in order to find stable co-locations.

169 Moving away from collaborative filtering, Marco et al. [10] present a mixture-of-experts approach
170 to model the memory behavior of Spark applications. It is based on a set of memory models (i.e., linear
171 regression, exponential regression, Napierian logarithmic regression) trained on a wide variety of
172 applications. At runtime, an expert selector based on k-nearest neighbour (kNN) is used to choose the
173 model that best describes memory behavior, in order to determine which tasks can be assigned to the
174 same host for improving throughput. The memory models and expert selector are trained offline on
175 different working sets, recording the memory used by a Spark executor through the Linux command
176 “/proc”. Finally, the scheduler uses the selected model to determine how much memory is required
177 for an incoming application, for improving server usage and system throughput.

178 Similarly to machine learning-based techniques, IIWM exploits a prediction model trained on
179 execution logs of previous workflows, however it differs in two main novel aspects: *i*) IIWM only uses
180 high-level workflow features, without requiring any runtime information as done in [8] and [10], in
181 order to avoid the overhead that could be not negligible for complex applications; *ii*) it provides an
182 algorithm for effectively scheduling a workflow in scenarios with limited computing resources.

183 As far as we know, no similar approaches in literature can be directly compared to IIWM in
184 terms of goals and requirements. In fact, differently from IIWM, *Quasar* [8] and *Cooper* [9] can be
185 seen as resource-efficient cluster management systems, aimed at optimizing QoS constraints and
186 resource usage. With respect to the most related work, presented in [10], IIWM presents the following
187 differences.

- 188 • It focuses on data-intensive workflows while in reference [10] general workloads are addressed.
- 189 • It uses high-level information for describing an application (e.g. task and dataset features), while
190 in reference [10] low-level system features are exploited, such as cache miss rate and number of
191 blocks sent, collected by running the application on a small portion (100 MB) of the input data.
- 192 • It proposes a more general approach, since the approach proposed in [10] is only applicable to
193 applications whose memory usage is a function of the input size.

194 3. Materials and Methods

195 The *Intelligent In-memory Workflow Manager* (IIWM) is based on three main steps:

- 196 1. *Execution monitoring and dataset creation*: starting from a given set of workflows, a transactional
197 dataset is generated by monitoring the memory usage and execution time of each task, specifying
198 how it is designed and giving concise information about the input.
- 199 2. *Prediction model training*: from the transactional dataset of executions, a regression model is
200 trained in order to fit the distribution of memory occupancy and execution time, according to the
201 features that represent the different tasks of a workflow.
- 202 3. *Workflow scheduling*: taking into account the predicted memory occupancy and execution time
203 of each task, provided by the trained model, and the available memory of the computing node,
204 tasks are scheduled using an informed strategy. In this way, a controlled degree of parallelism
205 can be ensured, while minimizing the risk of memory saturation.

206 In the following sections, a detailed description of each step is provided.

207 3.1. Execution monitoring and dataset creation

208 The first step in IIWM consists of monitoring the execution of different tasks on several input
 209 datasets with variable characteristics, in order to build a transactional dataset for training the regression
 210 model. The proposed solution was specifically designed for supporting the efficient execution of data
 211 analysis tasks, which are used in a wide range of data-intensive workflows. Specifically, it focuses
 212 on three classes of data mining tasks: *classification tasks* for supervised learning, *clustering tasks* for
 213 unsupervised learning and *association rules discovery*. Using Spark as a testbed, the following data
 214 mining algorithms from the MLib² library have been used: Decision Tree, Naive Bayes, and Support
 215 Vector Machines (SVM) for classification tasks; K-Means and Gaussian Mixture Models (GMM) for
 216 clustering tasks; FPGrowth for association rules tasks.

217 3.1.1. Execution monitoring within the Spark unified memory model

218 As far as execution monitoring is concerned, a brief overview of Spark unified memory model is
 219 required. In order to avoid OOM errors, Spark uses up to 90% of the heap memory, which is divided
 220 into three categories: *reserved memory* (300 MB), used to store Spark internal objects; *user memory* (40%
 221 of heap memory), used to store data structures and RDDs computed during transformations and
 222 actions; *spark memory* (60% of heap memory), divided in *execution* and *storage*. The former refers to that
 223 used for computation during shuffle, join, sort, and aggregation processes, while the latter is used for
 224 caching RDDs. It is worth noting that, when no execution memory is used, storage can acquire all the
 225 available memory and vice versa. However, storage may not evict execution due to complexities in
 226 implementation, while stored data blocks are evicted from main memory according to a Least Recently
 227 Used (LRU) strategy.

228 The occupancy of storage memory relies on the persistence operations performed natively by the
 229 algorithms. Table 2 reports some examples of data caching implemented in the aforementioned MLib
 230 algorithms. In particular, the *cache()* call corresponds to *persist(StorageLevel.MEMORY_AND_DISK)*,
 231 where MEMORY_AND_DISK is the default storage level used for the recent API based on DataFrames.

MLlib algorithm	Persist call
K-Means	<i>//Compute squared norms and cache them norms.cache()</i>
DecisionTree	<i>//Cache input RDD for speed-up during multiple passes BaggedPoint.convertToBaggedRDD(treeInput,...).cache()</i>
GMM	<i>instances.cache() ... data.map(_asBreeze).cache()</i>
FPGrowth	<i>items.cache()</i>
SVM	<i>IstanceBlock.blokifyWithMaxMemUsage(...).cache()</i>

Table 2. Examples of *persist* calls in MLib algorithms.

232 According to the Spark unified memory model, the execution monitoring was made via the Spark
 233 REST APIs, which expose executor-level performance metrics, collected in a JSON file, including peak
 234 occupancy for both execution and storage memory along with execution time.

235 3.1.2. Dataset creation

236 Using the aforementioned Spark APIs, we monitored the execution of several MLib algorithms
 237 on different input datasets, covering the main data mining tasks, i.e. classification, clustering, and

² <https://spark.apache.org/mlib/>

238 association rules. The goal of this process is the creation of a transactional dataset for the regression
239 model training, which contains the following information:

- 240 • The description of the task, such as its class (e.g., classification, clustering, etc.), type (fitting or
241 predicting task), and algorithm (e.g., SVM, K-Means, etc.).
- 242 • The description of the input dataset in terms of the number of rows, columns, categorical columns
243 and overall dataset size.
- 244 • Peak memory usage (both execution and storage) and execution time, which represent the three
245 target variables to be predicted by the regressor. In order to obtain more significant data, the
246 metrics were aggregated on median values by performing ten executions per task.

247 For the sake of clarity, table 3 shows a sample of the dataset described above.

Task Name	Task Type	Task Class	Dataset Rows	Dataset Columns	Categorical Columns	Dataset Size (MB)	Peak Storage Memory (MB)	Peak Execution Memory (MB)	Duration (ms)
GMM	Estimator	Clustering	1474971	28	0	87.0045	433.37	1413.5	108204
K-Means	Estimator	Clustering	5000000	104	0	1239.78	4624.52	4112	56233.5
DecisionTree	Estimator	Classification	9606	1921	0	84.9105	730.09	297.895	39292
NaiveBayes	Estimator	Classification	260924	4	0	13.4986	340.92	6982.82	16531.5
SVM	Estimator	Classification	5000000	129	0	1542.58	6199.11	106.6	238594.5
FPGrowth	Estimator	AssociationRules	823593	180	180	697	9493.85	1371.03	96071.5
GMM	Transformer	Clustering	165474	14	1	6.36604	2.34	1e-06	62.5
K-Means	Transformer	Clustering	4898431	42	3	648.887	3.23	1e-06	35
DecisionTree	Transformer	Classification	1959372	42	4	257.686	3.68	1e-06	65.5
NaiveBayes	Transformer	Classification	347899	4	0	17.9982	4.26	1e-06	92.5
SVM	Transformer	Classification	5000000	129	0	1542.58	2.36	1e-06	55.5
FPGrowth	Transformer	AssociationRules	136073	34	34	13.5493	1229.95	633.5	52429
...

Table 3. A sample of the training dataset.

248 Starting from 20 available datasets, we divided them into two partitions used for training and
249 testing respectively. Afterwards, an oversampling procedure was performed, aimed at increasing the
250 number of datasets contained in the partitions. Specifically, a naive random sampling approach can
251 lead to unexpected behaviors regarding the convergence of algorithms, thus introducing noise into
252 the transactional dataset used to build the regression model. To cope with this issue, we used the
253 following feature selection strategy:

- 254 • For datasets used in classification or regression tasks we considered only the k highest scoring
255 features based on:
 - 256 – analysis of variance (F-value) for integer labels (classification problems);
 - 257 – correlation-based univariate linear regression test for real labels (regression problems).
- 258 • For clustering datasets we used a correlation-based test to maintain the k features with the
259 smallest probability to be correlated with the others.
- 260 • For association rules discovery datasets no features selection is required, as the number of
261 columns refers to the average number of items in the different transactions.

262 The described procedure has been applied separately on the training and test partitions, so as to
263 avoid the introduction of bias into the evaluation process. Specifically, the number of datasets in the
264 training and test partitions has increased from 15 to 260 and from 5 to 86 respectively. Subsequently,
265 we fed these datasets to the MLlib algorithms, obtaining two final transactional datasets of 1309 and
266 309 monitored executions, used for training and testing the regressor, respectively.

267 3.2. Prediction model training

268 Once the training and test datasets with memory and time information were built, a regression
269 model can be trained with the goal of estimating peak memory occupancy and turnaround time of a
270 task in a given workflow.

271 As a preliminary step, we analyzed the correlation between the features of the training data and
272 each target variable, using the Spearman index. We obtained the following positive correlations: a

273 value of 0.30 between storage memory and the input dataset size, 0.46 between execution memory and
 274 the task class and 0.21 between execution time and the number of columns. These results can be seen
 275 in detail in Figure 1.

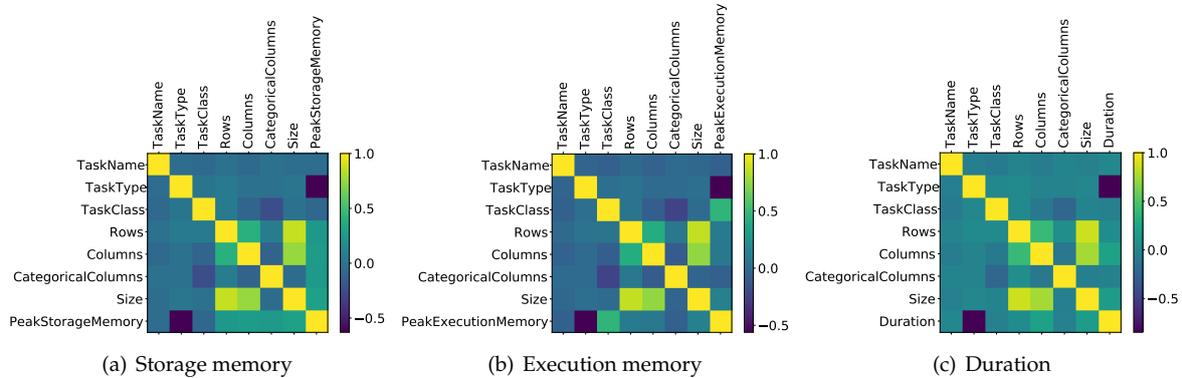


Figure 1. Correlation of target variables with the other features.

276 Afterwards we moved to the training of the regression model. Due to its complexity, the regression
 277 problem can not be faced with a simple linear regressor or its regularized variants (e.g. Ridge, Lasso or
 278 ElasticNet), but a more robust model is necessary. We experimentally evaluated this aspect by testing
 279 the forecasting abilities of these linear models achieving poor results. For this reason, an ensemble
 280 learning model has been used in order to fit the nonlinear distribution of features. Specifically, the
 281 *stacking* technique (meta learning) [20] has been used by developing a two-layer model in which a set
 282 of regressors are trained on the input dataset and a *Decision Tree* is fitted on their predictions. The first
 283 layer consists of three tree-based regressors, able to grasp different aspects of input data: a *Gradient*
 284 *Boosting*, an *AdaBoost* and an *Extra Trees* regressor. The second layer exploits a single *Decision Tree*
 285 regressor, which predicts the final value starting from the concatenation of the outputs from the first
 286 layer. The described ensemble model has been set with the hyper-parameters shown in Table 4.

Hyper-parameter	Value
<i>n_estimators</i>	500
<i>learning_rate</i>	0.01
<i>max_depth</i>	7
<i>loss</i>	<i>least squares</i>

Table 4. Hyper-parameters.

Among 20 trained models, initialized with different random states, we selected the best one by maximizing the following objective function:

$$\mathcal{O} = \bar{R}^2 - MAE$$

287 whose goal is to choose the model that best explains the variance of data, while minimizing the
 288 forecasting error. This function jointly considers the adjusted determination coefficient (\bar{R}^2), which
 289 guarantees robustness with respect to the addition of useless variables to the model compared to the
 290 classical R^2 score, and the mean absolute error (*MAE*), normalized with respect to the maximum.

291 The described model has been developed in *Python3* using the *scikit-learn*³ library and evaluated
 292 against the test set of 309 unseen executions obtained as described in Section 3.1.2. Thanks to the

³ <https://scikit-learn.org/stable/>

293 combination of different models, the ensemble technique showed to be very well suited for this task,
 294 leading to good robustness against outliers and a high forecasting accuracy, as shown in Figure 2.

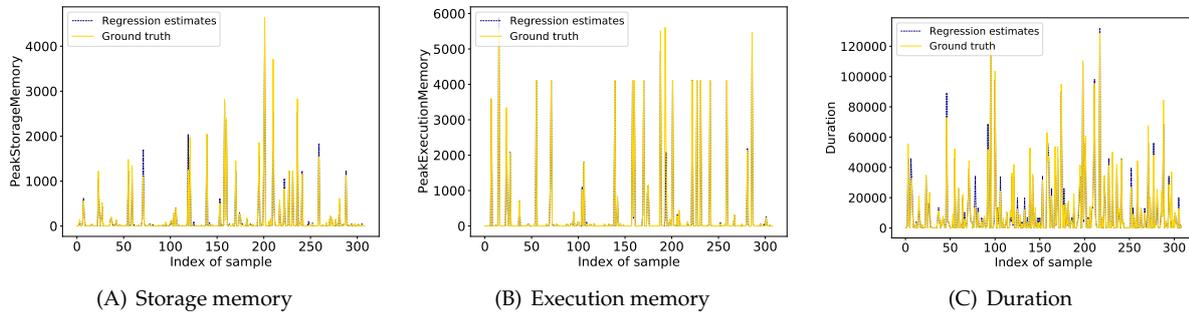


Figure 2. Meta-learner regression estimates for the different target variables.

295 These results are detailed in Table 5, which shows the evaluation metrics for each target variable,
 296 including the \bar{R}^2 score and the Pearson correlation coefficient. In particular, the Mean Absolute Error
 297 (MAE) and the Root Mean Square Error (RMSE) for the storage and execution memory represent
 298 average errors in megabytes, while for the duration they represent errors in milliseconds.

	RMSE	MAE	Adjusted R^2	Pearson Correlation
Storage Memory	108.23	26.66	0.96	0.98
Execution Memory	312.60	26.30	0.91	0.95
Duration	4443.17	2003.70	0.95	0.98

Table 5. Evaluation metrics on the test set.

299 3.3. Workflow scheduling

300 The prediction model described in Section 3.2 can be exploited to forecast the amount of memory
 301 that will be needed to execute a given task on a target computing node and its duration, based on
 302 the task features listed in Section 3.1. These predictions are then used within the scheduling strategy
 303 described in the following, whose goal is to avoid swapping to disk due to memory saturation in order
 304 to improve application performance and makespan through a better use of in-memory computing.
 305 The results discussed below refer to a static scheduling problem, as the scheduling plan is generated
 306 before the execution. In typical static scheduling the workflow system has to predict the execution
 307 load of each task accurately, using heuristic-based methods [21]. Likewise, in the proposed method the
 308 execution load of each task of a given workflow is predicted by the model trained on past executions.
 309 Moreover, we investigated how workflow tasks can be scheduled and run on a single computing node,
 310 but this approach can be easily generalized to a multi-node scenario. For example, a data-intensive
 311 workflow can be decomposed into multiple sub-workflows to be run on different computing nodes
 312 according to their features and data locality. Each sub-workflow is scheduled locally to the assigned
 313 node using the proposed strategy.

314 In IIWM, we modelled the scheduling problem as an *offline Bin Packing (BP)*. This is a well-known
 315 problem, widely used for resource and task management or scheduling, such as load balancing in
 316 mobile cloud computing architectures [22], energy-efficient execution of data-intensive applications
 317 in clouds [23], DAGs real-time scheduling in heterogeneous clusters [24] and task scheduling in
 318 multiprocessor environments [25]. Its classical formulation is as follows [26]. Let n be the number
 319 of items, w_j the weight of the j -th item and c the capacity of each bin: the goal is to assign each item
 320 to a bin without exceeding the capacity c and minimizing the number of used bins. The problem is
 321 \mathcal{NP} -complete and a lot of effort went into finding fast algorithms with near-optimal solutions. We
 322 adapted the classical problem to our purposes as follows:

- 323 • An item is a *task* to be executed.
- 324 • A bin identifies a *stage*, i.e. a set of tasks that can be run in parallel.
- 325 • The capacity of a bin is the maximum amount C of available memory in a computing node.
- 326 When assigning a task to a stage $s \in \mathcal{S}$, its residual available memory will be indicated with C_s .
- 327 • The weight of an item is the *memory occupancy* estimated by the prediction model. In the case of
- 328 Spark testbed, it will be the maximum of the execution and storage memory, in order to model
- 329 a peak in the unified memory. For what concerns the estimated *execution time*, it is used for
- 330 selecting the stage to be assigned when memory constraints hold for multiple stages.

331 With respect to the classical BP problem two changes were introduced:

- 332 • All workflow tasks have to be executed, so the capacity of a stage may still be exceeded if a task
- 333 takes up more memory than the available one.
- 334 • The assignment of a task t to a stage s is subjected to dependency constraints. Hence, if a
- 335 dependency exists between t_i and t_j , then the stage of t_i has to be executed before the one of t_j .

336 To solve the BP problem, modelled as described above, in order to produce the final scheduling
 337 plan, we used the *First Fit Decreasing* algorithm which assigns tasks sorted in non-increasing order of
 338 weight. However, the introduction of dependency constraints in the assignment process may cause the
 339 under-usage of certain stages. To cope with this issue, we introduced a further step of consolidation,
 340 aimed at reducing the number of stages by merging together stages without dependencies according
 341 to the available memory. The main execution flow of the IIWM scheduler is shown in Figure 3 and
 342 described by Algorithm 1. In particular, given a data-intensive workflow \mathcal{W} , described as a DAG by
 343 its tasks and dependencies, and the prediction model \mathcal{M} as input, a scheduling plan is generated in
 344 two steps: *i*) building of the stages and task assignment; *ii*) stage consolidation.

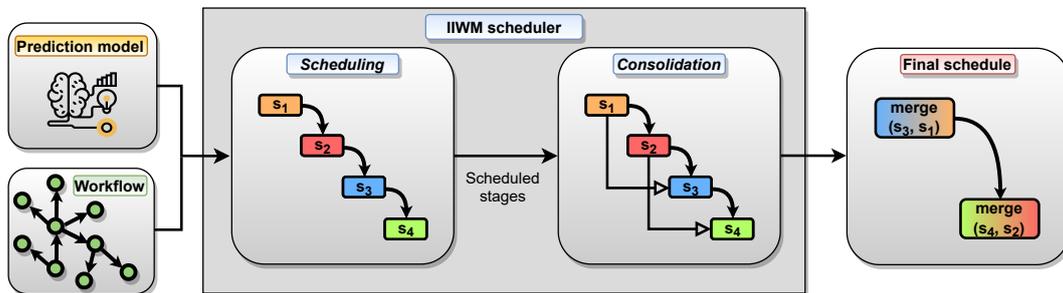


Figure 3. Execution flow of the IIWM scheduler. Given a workflow and a prediction model as input, a scheduling plan is generated in two steps: *i*) building of the stages and task assignment; *ii*) stage consolidation.

345 The algorithm is divided into two main parts: in the first part (lines 1-23), the stages are built by
 346 iteratively assigning each task according to the estimates of the prediction model; in the second part
 347 (lines 25-34), a consolidation process is performed, trying to minimize the number of stages.

348 The first part (lines 1-23) starts with the initialization of an empty list of stages \mathcal{S} , which will be
 349 filled according to a dictionary \mathcal{Q} that stores the in-degree of each task in the DAG, which is used for
 350 identifying the free tasks which can be scheduled. The prediction model \mathcal{M} is exploited to estimate
 351 the memory occupancy and execution time of each task in \mathcal{T} , according to their dataset description
 352 (lines 3-4). The dictionary \mathcal{P}_{mem} , which collects the predicted memory occupancies, is then used to
 353 sort tasks according to the First Fit Decreasing strategy (line 5). At each iteration, tasks that can be
 354 scheduled (i.e., assigned to a stage) are collected in the \mathcal{T}_{free} set. In particular, they are identified by
 355 a zero in-degree, as their execution does not depend on others (line 7). By virtue of the acyclicity of
 356 the DAG-based workflow representation, there will always exist a task $t \in \mathcal{T}$ with a zero in-degree
 357 not yet scheduled, unless set \mathcal{T} is empty. Afterwards, the task with the highest memory occupancy is

358 selected from \mathcal{T}_{free} in order to be scheduled (line 8). At this point, a list of candidate stages (\mathcal{S}_{sel}) for
 359 the selected task is identified according to the peak memory occupancy forecasted by the prediction
 360 model \mathcal{M} (lines 9-10). In particular, a stage s_i belongs to \mathcal{S}_{sel} if it satisfies the following conditions:

- 361 • The residual capacity C_{s_i} of the selected stage s_i is not exceeded by the addition of the task t .
- 362 • There not exists a dependency between t and any task t' belonging to s_i and every subsequent
 363 stage $(s_{i+1} \cup \dots \cup s_k)$, where a dependency $(t', t)^n$ is identified by a path of length $n > 0$.

364 If there exist one or more candidate stages \mathcal{S}_{sel} (line 11), the best one is chosen based on the
 365 minimum *marginal increase*. Specifically, for each of these stages, the expected increase of the execution
 366 time is estimated (lines 12-13), assigning the task t to the stage s with the lowest value (lines 14-16).
 367 Otherwise (line 17), a newly created stage is allocated for t and added to the list \mathcal{S} (lines 18-21). Once
 368 the task t is assigned to the stage s , the residual capacity C_s is updated (lines 15, 20). Then, the residual
 369 in-degree for every task in the out-neighbourhood of t (line 22) is decremented by updating the
 370 dictionary \mathcal{Q} , so as to allow the assignment of these tasks in the next iterations. Finally, the assigned
 371 task t is removed from the set of workflow nodes \mathcal{T} (line 23).

ALGORITHM 1: IIWM SCHEDULER

Input: Workflow $\mathcal{W} = (\mathcal{T}, \mathcal{A})$, Prediction model \mathcal{M}

Output: A list of stages \mathcal{S}

```

1  $\mathcal{S} \leftarrow \emptyset$ 
2  $\mathcal{Q} \leftarrow \langle t : |\mathcal{N}^{in}(t)|, \forall t \in \mathcal{T} \rangle$ 
3  $\mathcal{P}_{mem} \leftarrow \langle t : \mathcal{M}.predict\_mem(t, d_t), \forall t \in \mathcal{T} \rangle$  ▷ Memory prediction for each task in  $\mathcal{T}$ 
4  $\mathcal{P}_{time} \leftarrow \langle t : \mathcal{M}.predict\_time(t, d_t), \forall t \in \mathcal{T} \rangle$  ▷ Time prediction for each task in  $\mathcal{T}$ 
5  $\mathcal{T} \leftarrow sort\_decreasing(\mathcal{T}, \mathcal{P}_{mem})$ 
6 while  $\mathcal{T} \neq \emptyset$  do
7    $\mathcal{T}_{free} \leftarrow \{t \in \mathcal{T} \mid \mathcal{Q}[t] == 0\}$ 
8    $t \leftarrow get\_first(\mathcal{T}_{free})$ 
9    $mem_t \leftarrow \mathcal{P}_{mem}[t]$ 
10   $\mathcal{S}_{sel} \leftarrow \{s_i \in \mathcal{S} \mid mem_t \leq C_{s_i} \text{ and } \nexists (t', t)^n \in \mathcal{A}, n > 0, \forall t' \in s_i \cup s_{i+1} \cup \dots \cup s_k\}$ 
11  if  $\mathcal{S}_{sel} \neq \emptyset$  then
12     $duration \leftarrow \langle s : \max_{t' \in s} \mathcal{P}_{time}[t'], \forall s \in \mathcal{S}_{sel} \rangle$ 
13     $increase \leftarrow \langle s : \max\{\mathcal{P}_{time}[t], duration[s]\} - duration[s], \forall s \in \mathcal{S}_{sel} \rangle$ 
14     $s \leftarrow argmin_{s' \in \mathcal{S}_{sel}} increase$ 
15     $C_s \leftarrow C_s - mem_t$ 
16     $s \leftarrow s \cup \{t\}$ 
17  else
18     $s \leftarrow \emptyset$ 
19     $s \leftarrow s \cup \{t\}$ 
20     $C_s \leftarrow C_s - mem_t$ 
21     $\mathcal{S} \leftarrow \mathcal{S} \cup \{s\}$ 
22   $\mathcal{Q}[t'] = \mathcal{Q}[t'] - 1, \forall t' \in \mathcal{N}^{out}(t)$ 
23   $\mathcal{T} \leftarrow \mathcal{T} \setminus \{t\}$ 
24 // Consolidation step
25  $\mathcal{S}_{mov} \leftarrow \{s \in \mathcal{S} \mid |\mathcal{N}^{out}(t)| == 0, \forall t \in s\}$ 
26 if  $\mathcal{S}_{mov} \neq \emptyset$  then
27   for  $s_i \in \mathcal{S}_{mov}$  do
28     for  $s_j \in \mathcal{S} \mid j > i$  do
29        $mem_{s_i \cup s_j} \leftarrow \sum_{t \in s_i \cup s_j} \mathcal{P}_{mem}[t]$ 
30       if  $mem_{s_i \cup s_j} \leq C$  then
31          $s_j \leftarrow s_i \cup s_j$ 
32          $\mathcal{S} \leftarrow \mathcal{S} \setminus s_i$ 
33         break
34 return  $\mathcal{S}$ 

```

372 The second part of the algorithm (lines 25-34) performs a consolidation step with the goal of
 373 reducing the number of allocated stages by merging some of them if possible, with a consequential
 374 improvement in the global throughput. The stages involved in the consolidation step, namely the
 375 movable stages (\mathcal{S}_{mov}), are those containing tasks with a zero out-degree (line 25). This means that no
 376 task in such stages blocks the execution of another one, so they can be moved forward and merged
 377 with subsequent stages if the available capacity C is not exceeded. For each movable stage s_i (line
 378 27), another stage s_j from \mathcal{S} is searched among the subsequent ones, such that its residual capacity is
 379 enough to enable the merging with s_i (lines 28-30). The merging between s_i and s_j is performed by
 380 assigning to s_j each task of s_i (line 31), finally removing s_i from \mathcal{S} (line 32). In the end, the list of stages
 381 \mathcal{S} built by the scheduler is returned as output. Given this scheduling plan, the obtained stages will be
 382 executed in sequential order, while all the tasks in a stage will run concurrently.

383 Compared to a blind strategy where the maximum parallelism is achieved by running in parallel
 384 all the tasks not subjected to dependencies, which will be referred to as *Full-Parallel* in our experiments,
 385 IIWM can reduce both delays of parallelization (ϵ_p), due to context switch and process synchronization,
 386 and swapping/spilling to disk (ϵ_s), due to I/O operations. Delay ϵ_p is always present in all scheduling
 387 strategies when two or more tasks are run concurrently, while ϵ_s is present only when a memory
 388 saturation event occurs. Given $\epsilon = \epsilon_p + \epsilon_s$, IIWM mainly reduces ϵ_s , which is the main factor behind
 389 the drop in performance in terms of execution time, due to the slowness in accessing secondary storage.

390 As far as the Spark framework is concerned, the proposed strategy is effective for making the
 391 most of the default storage level, i.e. *MEMORY_AND_DISK*: at each internal call of the *cache()* method,
 392 data is saved in-memory as long as this resource is available, using disk otherwise. In this respect,
 393 IIWM can reduce the actual persistence of data on disk by better exploiting in-memory computing.

394 4. Results and Discussion

395 This section presents an experimental evaluation of the proposed system, aimed at optimizing the
 396 in-memory execution of data-intensive workflows. We experimentally assessed the effectiveness of
 397 IIWM using Apache Spark 3.0.1 as a testbed. In particular, we generated two synthetic workflows for
 398 analyzing different scenarios, by assessing also the benefits coming from the use of IIWM using a real
 399 data mining workflow as a case study.

400 In order to provide significant results, each experiment was executed ten times and the average
 401 metrics with standard deviations are reported. In particular, for each experiment, we evaluated the
 402 accuracy of the regression model in predicting memory occupancy and execution time.

403 We evaluated the ability of IIWM to improve application performance taking into account two
 404 different aspects:

- *Execution time*: let m_1 and m_2 be the makespan for two different executions. If $m_2 < m_1$ we can compute the improvement on makespan (m_{imp}) and application performance (p_{imp}) as follows:

$$m_{imp} = \frac{m_1 - m_2}{m_1} \times 100\% \quad p_{imp} = \frac{m_1}{m_2}$$

- *Disk usage*: we used the *on-disk usage* metric, which measures the amount of disk usage, jointly considering the volume and the duration of disk writes. Formally, given a sequence of disk writes w_1, \dots, w_k let $\tau'_i, \tau''_i \in \mathbb{T}$ be the start and end time of the w_i write respectively. Let also $W : \mathbb{T} \rightarrow \mathbb{R}$ be a function representing the amount of megabytes written to disk over time \mathbb{T} . We define on-disk usage as:

$$on\text{-}disk\ usage = \sum_{i=1}^k \frac{1}{\tau''_i - \tau'_i} \int_{\tau'_i}^{\tau''_i} W(\tau) d\tau$$

405 Specifically, for each workflow we reported: *i*) a comparison between Full-Parallel and IIWM
 406 in terms of disk usage over time; *ii*) a detailed description of the scheduling plan generated by both

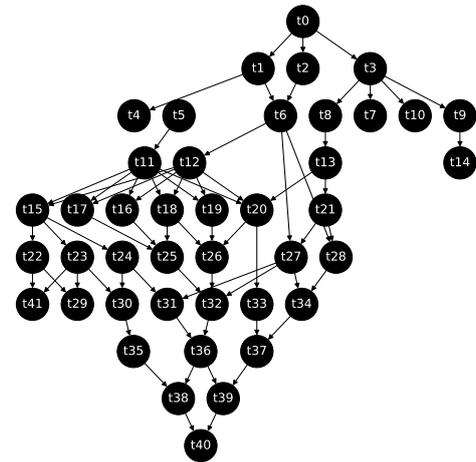
407 strategies; *iii*) the average improvement on makespan and application performance with IIWM; *iv*)
 408 statistics about the use of disk, such as the time spent for I/O operations and the *on-disk usage* metric;
 409 *v*) the execution of the workflow by varying the amount of available memory, in order to show the
 410 benefits of the proposed scheduler in different limited memory scenarios.

411 4.1. Synthetic workflows

412 We firstly evaluated our approach against two complex synthetic data analysis workflows, where
 413 the Full-Parallel approach showed its limitations due to a high degree of parallelism. The dependencies
 414 in these workflows should be understood as execution constraints. For instance, clustering has to be
 415 performed before classification for adding labels to an unlabelled dataset, or a classification task is
 416 performed after the discovery of association rules for user classification purposes.

417 The first test has been carried out on a synthetic workflow with 42 nodes. Table 6 provides a
 418 detailed description of each task in the workflow, while their dependencies are shown in Figure 4.

Node	Task Name	Task Type	Task Class	Rows	Columns	Categorical Columns	Dataset size (MB)
t ₀	NaiveBayes	Estimator	Classification	2939059	18	4	198.93904
t ₁	FPGrowth	Estimator	AssociationRules	494156	180	180	417.01007
t ₂	NaiveBayes	Estimator	Classification	5000000	27	0	321.86484
t ₃	K-Means	Estimator	Clustering	1000000	104	0	247.95723
t ₄	DecisionTree	Estimator	Classification	4000000	53	0	505.45
t ₅	DecisionTree	Estimator	Classification	4000000	27	0	257.4918
t ₆	DecisionTree	Estimator	Classification	5000000	129	0	1542.5775
t ₇	K-Means	Estimator	Clustering	2000000	53	0	252.72397
t ₈	NaiveBayes	Estimator	Classification	2000000	104	0	495.9038
t ₉	NaiveBayes	Estimator	Classification	1000000	129	0	307.56625
t ₁₀	SVM	Estimator	Classification	2000000	53	0	252.72397
t ₁₁	K-Means	Estimator	Clustering	2049280	9	2	122.02598
t ₁₂	GMM	Estimator	Clustering	2458285	28	0	145.00838
t ₁₃	K-Means	Estimator	Clustering	9169	5812	1	101.88691
t ₁₄	SVM	Estimator	Classification	2000000	27	0	128.74657
t ₁₅	K-Means	Estimator	Clustering	3000000	104	0	743.87286
t ₁₆	SVM	Estimator	Classification	3000000	53	0	379.08823
t ₁₇	SVM	Estimator	Classification	14410	1921	0	127.3811
t ₁₈	K-Means	Estimator	Clustering	5000000	53	0	631.8128
t ₁₉	K-Means	Estimator	Clustering	5000000	104	0	1239.7812
t ₂₀	K-Means	Estimator	Clustering	2000000	78	0	371.93442
t ₂₁	SVM	Estimator	Classification	3000000	104	0	743.87286
t ₂₂	K-Means	Estimator	Clustering	2939059	18	4	198.93904
t ₂₃	SVM	Estimator	Classification	19213	1442	0	123.28475
t ₂₄	DecisionTree	Estimator	Classification	3000000	129	0	922.6897
t ₂₅	K-Means	Estimator	Clustering	1959372	26	4	189.5505
t ₂₆	DecisionTree	Estimator	Classification	4898431	18	4	331.56735
t ₂₇	NaiveBayes	Estimator	Classification	4898431	18	4	331.56735
t ₂₈	K-Means	Estimator	Clustering	2939059	34	4	334.91486
t ₂₉	K-Means	Estimator	Clustering	4898431	18	4	331.56735
t ₃₀	K-Means	Estimator	Clustering	1966628	42	0	170.48729
t ₃₁	NaiveBayes	Estimator	Classification	1959372	18	4	132.62437
t ₃₂	K-Means	Estimator	Clustering	3000000	78	0	557.9056
t ₃₃	DecisionTree	Estimator	Classification	3000000	53	0	379.08823
t ₃₄	DecisionTree	Estimator	Classification	14410	2401	0	159.71497
t ₃₅	K-Means	Estimator	Clustering	2939059	42	4	386.53033
t ₃₆	DecisionTree	Estimator	Classification	2939059	34	4	334.91486
t ₃₇	DecisionTree	Estimator	Classification	4000000	129	0	1230.2445
t ₃₈	NaiveBayes	Estimator	Classification	1000000	53	0	126.36286
t ₃₉	GMM	Estimator	Clustering	1000000	53	0	126.36286
t ₄₀	DecisionTree	Estimator	Classification	2939059	18	4	198.93904
t ₄₁	K-Means	Estimator	Clustering	4898431	18	4	331.56735



4.24 example with IIWM, focusing on its main steps: *i*) the scheduling of tasks based on their decreasing
 4.25 memory weight; *ii*) the allocation of a new stage; *iii*) the exploitation of the estimated execution time
 4.26 while computing the marginal increase. This last aspect can be clearly observed in iteration 17, where
 4.27 task t_{17} is assigned to stage s_7 , which presents a marginal increase equal to zero. This is the best
 4.28 choice compared to the other candidate stage (s_6), whose execution time would be increased by 12496
 4.29 milliseconds by the assignment of t_{17} , with a degradation of the overall makespan.

Iteration	State	Stages
It. 0	$\mathcal{T}_{free}^0 = \{t_0\}$ Create s_0 and assign t_0 Unlock $\{t_1, t_2, t_3\}$	$s_0 = \{t_0\}$
It. 1	$\mathcal{T}_{free}^1 = \{t_1, t_3, t_2\}$ Create s_1 and assign t_1 Unlock $\{t_4\}$	$s_0 = \{t_0\}, s_1 = \{t_1\}$
It. 2	$\mathcal{T}_{free}^2 = \{t_3, t_4, t_2\}$ Create s_2 and assign t_3 Unlock $\{t_7, t_8, t_9, t_{10}\}$	$s_0 = \{t_0\}, s_1 = \{t_1\},$ $s_2 = \{t_3\}$
It. 3	$\mathcal{T}_{free}^3 = \{t_7, t_4, t_{10}, t_2, t_8, t_9\}$ Create s_3 and assign t_7	$s_0 = \{t_0\}, s_1 = \{t_1\},$ $s_2 = \{t_3\}, s_3 = \{t_7\}$
It. 4	$\mathcal{T}_{free}^4 = \{t_4, t_{10}, t_2, t_8, t_9\}$ $\mathcal{S}_{sel} = \{s_2, s_3\}$ $increase = \{0, 0\}$ Assign t_4 to s_2	$s_0 = \{t_0\}, s_1 = \{t_1\},$ $s_2 = \{t_3, t_4\}, s_3 = \{t_7\}$
...
It. 17	$\mathcal{T}_{free}^{17} = \{t_{17}, t_{23}, t_8, t_9\}$ $\mathcal{S}_{sel} = \{s_6, s_7\}$ $increase = \{12496.363, 0\}$ Assign t_{17} to s_7 Unlock $\{t_{25}\}$	$s_0 = \{t_0\}, s_1 = \{t_1, t_2\},$ $s_2 = \{t_3, t_4, t_5\},$ $s_3 = \{t_7, t_{10}, t_6\},$ $s_4 = \{t_{12}, t_{11}\}, s_5 = \{t_{15}, t_{18}\},$ $s_6 = \{t_{19}, t_{22}\}, s_7 = \{t_{24}, t_{16}, t_{17}\}$
...

Table 8. Example of execution of algorithm 1 at iteration level.

4.30 At the end of the process, a consolidation step is exploited for optimizing throughput and
 4.31 execution time, by merging two stages with zero out-degree with some tailing stages, so as to avoid
 4.32 the sequential execution of the two stages in favour of a parallel one.

4.33 Figure 5 shows disk occupancy throughout the execution. As a consequence of memory saturation,
 4.34 the execution of Full-Parallel resulted in a huge amount of disk writes, while IIWM achieved a null
 4.35 disk usage since no swapping occurred thanks to intelligent task scheduling. Thus, this translates into
 4.36 better use of in-memory computing.

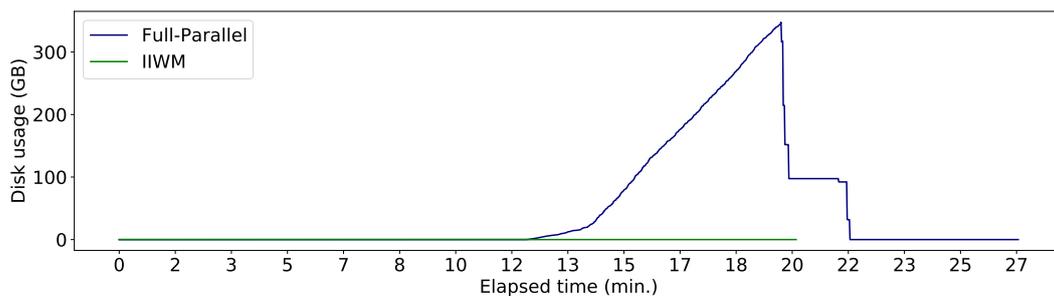


Figure 5. Disk usage over time for Full-Parallel and IIWM.

437 These results can be clearly seen also in Table 9, which shows the scheduling plan produced
 438 by the IIWM scheduler, together with some statistics about execution times and the use of the disk.
 439 In particular, given the curves representing disk writes over time shown in Figure 5, *on-disk usage*
 440 graphically represents the sum, for each disk write, of the ratio between the area under the curve
 441 identified by a write and its duration. Compared to the Full-Parallel strategy, IIWM achieved better
 442 execution times and an improvement in application performance, with a boost of almost $1.45x$ (p_{imp})
 443 and a 31.15% reduction in time (m_{imp}) on average.

Strategy	Task-scheduling plan	Number of stages	Time (min.)	Peak disk usage (MB)	Writes duration (min.)	On-disk usage (MB)
Full-Parallel	$(t_0), (t_1 \parallel t_2 \parallel t_3),$ $(t_4 \parallel t_5 \parallel t_6 \parallel t_7 \parallel t_8 \parallel t_9 \parallel t_{10}),$ $(t_{11} \parallel t_{12} \parallel t_{13} \parallel t_{14}),$ $(t_{15} \parallel t_{16} \parallel t_{17} \parallel t_{18} \parallel t_{19} \parallel t_{20} \parallel t_{21}),$ $(t_{22} \parallel t_{23} \parallel t_{24} \parallel t_{25} \parallel t_{26} \parallel t_{27} \parallel t_{28}),$ $(t_{29} \parallel t_{30} \parallel t_{31} \parallel t_{32} \parallel t_{33} \parallel t_{34} \parallel t_{41}),$ $(t_{35} \parallel t_{36} \parallel t_{37}), (t_{38} \parallel t_{39}), (t_{40})$	10	31.52 ± 0.6	356106.601	11.56	126867.065
IIWM	$(t_0), (t_1 \parallel t_2), (t_3 \parallel t_4 \parallel t_5),$ $(t_7 \parallel t_{10} \parallel t_6 \parallel t_8 \parallel t_9),$ $(t_{12} \parallel t_{11} \parallel t_{13}), (t_{15} \parallel t_{18}), (t_{19} \parallel t_{22} \parallel t_{23}),$ $(t_{24} \parallel t_{16} \parallel t_{17} \parallel t_{29}), (t_{25} \parallel t_{41}), (t_{30} \parallel t_{20}),$ $(t_{35} \parallel t_{21} \parallel t_{14}), (t_{28} \parallel t_{26} \parallel t_{27}),$ $(t_{33} \parallel t_{32} \parallel t_{34} \parallel t_{31}), (t_{37} \parallel t_{36}),$ $(t_{39} \parallel t_{38}), (t_{40})$	16	21.70 ± 0.63	0	0	0

Table 9. Scheduling plan and statistics about execution times and disk usage with 14 GB of RAM.

444 With different sizes of available memory, the Full-Parallel approach showed higher and higher
 445 execution times and disk writes as memory decreased, while IIWM was able to adapt the execution
 446 to available resources as shown in Figure 6, finding a good trade-off between the maximization of
 447 the parallelism and the minimization of the memory saturation probability. At the extremes, with
 448 unlimited available memory, or at least greater than that required to run the workflow, IIWM will
 449 perform as a full concurrent strategy, producing the same scheduling of Full-Parallel.

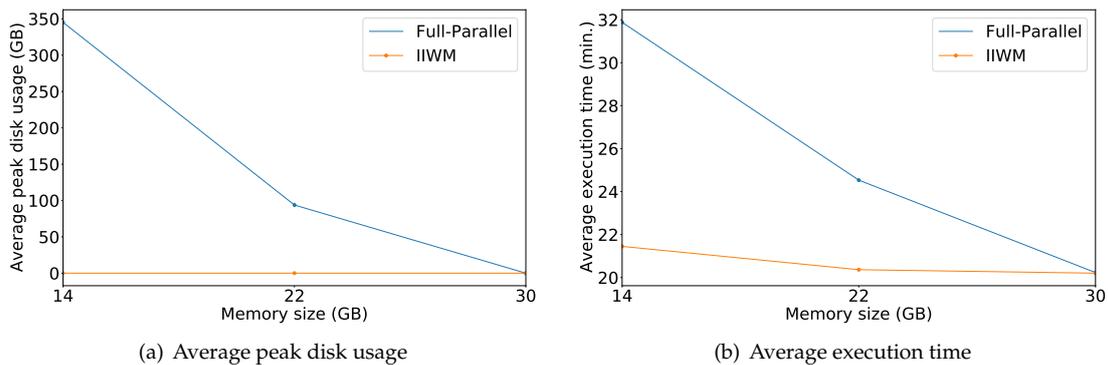


Figure 6. Average peak disk usage and execution time, varying the size of available memory.

450 The second synthetic workflow consists of the 27 tasks described by Table 10 and their
 451 dependencies, shown in Figure 7. This scenario is characterized by highly heavy tasks and very
 452 low resources, where the execution of a single task can exceed the available memory. In particular, task
 453 T_{18} has an estimated peak memory occupancy higher than Spark available unified memory of 5413.8
 454 MB (i.e., corresponding to a heap size of 9.5 GB): this will bring the IIWM scheduling algorithm to
 455 allocate the task to a new stage, but memory will be saturated anyway.

456 In such a situation, data spilling to disk cannot be avoided, but IIWM tries to minimize the
 457 number of bytes written and the duration of I/O operations. Even in this scenario, the prediction
 458 model achieved very accurate results, shown in Table 11, confirming its forecasting abilities.

471 are provided to the scheduler. Nevertheless, in the proposed work we dealt with a static scheduling
 472 problem, where all tasks are known in advance and the task-set is not modifiable at runtime.

Strategy	Task-scheduling plan	Number of stages	Time (min.)	Peak disk usage (MB)	Writes duration (min.)	On-disk usage (MB)
Full-Parallel	$(t_0), (t_1 \parallel t_2 \parallel t_3 \parallel t_4),$ $(t_5 \parallel t_6 \parallel t_7 \parallel t_8 \parallel t_9 \parallel t_{10} \parallel t_{11} \parallel t_{12}),$ $(t_{13} \parallel t_{14} \parallel t_{15} \parallel t_{16} \parallel t_{17} \parallel t_{18} \parallel t_{19} \parallel t_{20}),$ $(t_{21}), (t_{22} \parallel t_{23} \parallel t_{24} \parallel t_{25}), (t_{26})$	7	29.42 ± 1.88	27095.837	20.6	10593.79
IWM	$(t_0), (t_4 \parallel t_2), (t_{11} \parallel t_7), (t_8 \parallel t_3),$ $(t_{15} \parallel t_{10} \parallel t_9 \parallel t_{16}), (t_{18}),$ $(t_{17} \parallel t_1 \parallel t_{12}), (t_6 \parallel t_5), (t_{14}),$ $(t_{13} \parallel t_{20} \parallel t_{19}), (t_{21}),$ $(t_{23} \parallel t_{24}), (t_{25}), (t_{22}), (t_{26})$	15	22.68 ± 1.65	304.5	3.6	60.82

Table 12. Scheduling plan and statistics about execution times and disk usage with 9.5 GB of RAM.

4.2. Real case study

474 In order to assess the performance of the proposed approach against a real case study, we used
 475 a data mining workflow [27] that implements a model selection strategy for the classification of an
 476 unlabelled dataset. Figure 9 shows a representation of the workflow designed by the visual language
 477 VL4Cloud [28]. A training set is divided into n partitions and k classification algorithms are fitted on
 478 each partition for generating $k \times n$ classification models. The $k \times n$ fitted models are evaluated by a
 479 model selector on a test set to choose the best model. Afterwards, the n predictors use the best model to
 480 generate n classified datasets. The following k classification algorithms provided by the MLlib library
 481 were used: Decision Tree with C4.5 algorithm, Support Vector Machines (SVM), and Naive Bayes. The
 482 training set, test set and unlabelled dataset provided as input for the workflow have been generated
 483 from the *Physical Unclonable Functions (PUFs)* [29] simulation through a n -fold-cross strategy. In this
 484 scenario, IWM can be used to optimize the data processing phase regarding the execution of the $k \times n$
 485 classification algorithms (estimators first, transformers then) concurrently. The other phases, such as
 486 data acquisition and partitioning, are out of our interest. The red box in Figure 9 shows the tasks of the
 487 workflow that will be analyzed.

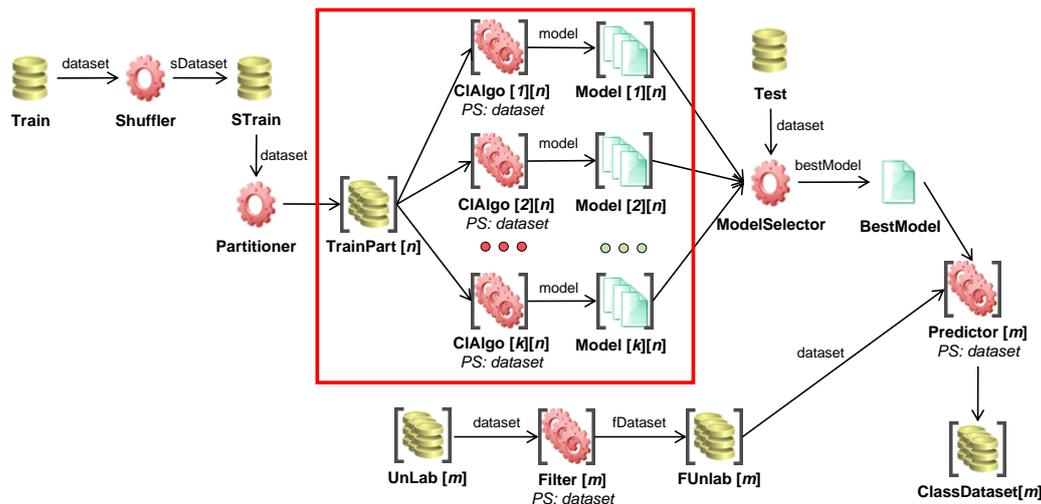


Figure 9. Ensemble learning workflow.

488 Figure 10 shows disk occupancy over time with 14 GB of RAM. Also in this case, IWM avoided
 489 disk writes, while Full-Parallel registered a high level of disk usage. In particular, during the training
 490 phase, the parallel execution of the $k \times n$ models (with $k = 3$ and $n = 5$) saturates memory with 15
 491 concurrent tasks and generates disk writes up to 124 GB.

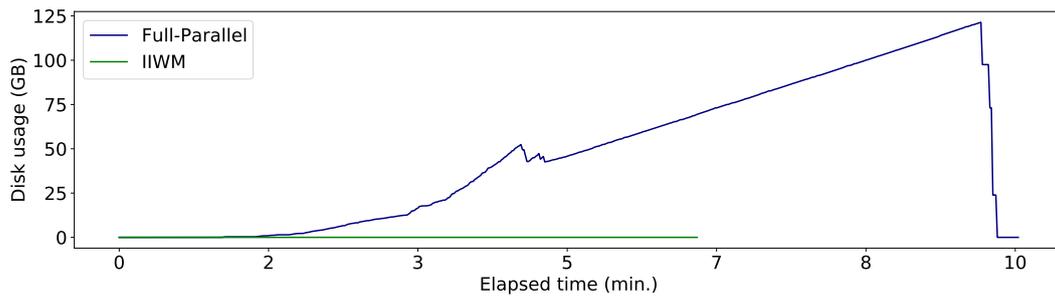


Figure 10. Disk usage over time for Full-Parallel and IIWM.

492 The results are detailed in Table 13, which shows a boost in execution time of almost 1.66x (p_{imp})
 493 and a 40% time reduction (m_{imp}) with respect to Full-Parallel.

Strategy	Task-scheduling plan	Number of stages	Time (min.)	Peak disk usage (MB)	Writes duration (min.)	On-disk usage (MB)
Full-Parallel	$(t_0 \parallel t_2 \parallel t_4 \parallel t_6 \parallel t_8 \parallel t_{10} \parallel t_{12} \parallel t_{14} \parallel t_{16} \parallel t_{18} \parallel t_{20} \parallel t_{22} \parallel t_{24} \parallel t_{26} \parallel t_{28}),$ $(t_1 \parallel t_3 \parallel t_5 \parallel t_7 \parallel t_9 \parallel t_{11} \parallel t_{13} \parallel t_{15} \parallel t_{17} \parallel t_{19} \parallel t_{21} \parallel t_{23} \parallel t_{25} \parallel t_{27} \parallel t_{29})$	2	11.42 ± 0.27	124730.874	9.6	54443.186
IIWM	$(t_{10} \parallel t_{12} \parallel t_{14} \parallel t_{16} \parallel t_{20} \parallel t_{22} \parallel t_{24} \parallel t_{26} \parallel t_{28}),$ $(t_{18} \parallel t_0 \parallel t_2 \parallel t_4 \parallel t_6 \parallel t_8 \parallel t_{11} \parallel t_{13} \parallel t_{15} \parallel t_{17} \parallel t_{21} \parallel t_{23} \parallel t_{25} \parallel t_{27} \parallel t_{29}),$ $(t_{19} \parallel t_1 \parallel t_3 \parallel t_5 \parallel t_7 \parallel t_9)$	3	6.88 ± 0.1	0	0	0

Table 13. Scheduling plan and statistics about execution times and disk usage with 14 GB of RAM.

494 The general trends varying the amount of available resources are also confirmed with respect to
 495 the previous examples, as shown in Figure 11.

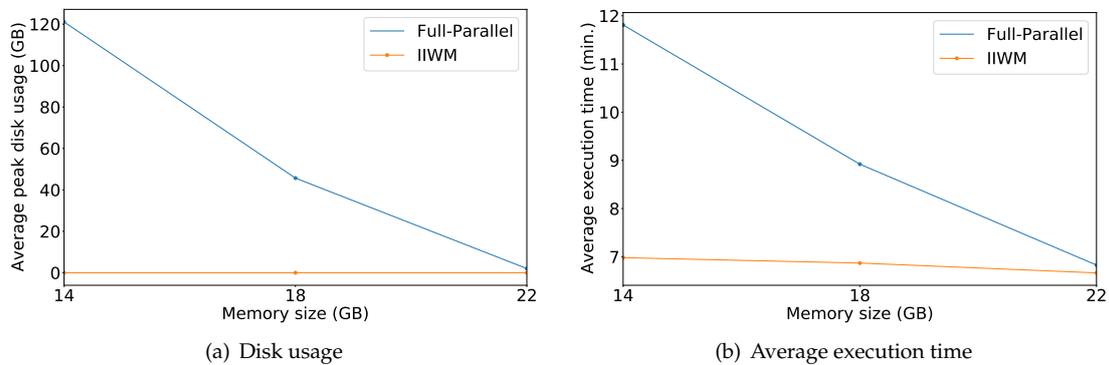


Figure 11. Average peak disk usage and execution time, varying the size of available memory.

496 5. Conclusions and Final Remarks

497 Nowadays, data-intensive workflows are widely used in several application domains, such as
 498 bioinformatics, astronomy, and engineering. This paper introduced and evaluated a system, named
 499 Intelligent In-memory Workflow Manager (IIWM), aimed at optimizing the in-memory execution of
 500 data-intensive workflows on high-performance computing systems. Experimental results suggest
 501 that by jointly using a machine learning model for performance estimation and a suitable scheduling
 502 strategy, the execution of data-intensive workflows can be significantly improved with respect to
 503 state-of-the-art blind strategies. In particular, the main benefits of IIWM resulted when has been
 504 applied to workflows having a high level of parallelism. In this case a significant reduction of memory
 505 saturation has been obtained. Therefore it can be used effectively when multiple tasks have to be
 506 executed on the same computing node, for example when they need to be run on multiple immovable
 507 datasets located on a single node or due to other hardware constraints. In these cases, an uninformed
 508 scheduling strategy will likely exceed the available memory, causing disk writes and therefore a

509 drop in performance. The proposed approach also showed to be a very suitable solution in scenarios
510 characterized by a limited amount of memory reserved for execution, thus finding possible applications
511 in data-intensive IoT workflows, where data processing is performed on constrained devices located
512 at the network edge.

513 IIWM has been evaluated against different scenarios concerning both synthetic and real data
514 mining workflows, using Apache Spark as a testbed. Specifically, by accurately predicting resources
515 used at runtime, our approach achieved up to 31% and 40% reduction of makespan and a performance
516 improvement up to 1.45x and 1.66x for the synthetic workflows and the real case study respectively.

517 In future work additional aspects of performance estimation will be investigated. For example,
518 the IIWM prediction model can be extended also to consider other common stages in workflows
519 besides data analysis, such as data acquisition, integration and reduction, and other information about
520 tasks, input data, and hardware platform features can be exploited in the scheduling strategy.

521 6. Acknowledgment

522 This work has been supported by the ASPIDE Project funded by the European Union's Horizon
523 2020 Research and Innovation Programme under grant agreement No 801091.

524 References

- 525 1. Talia, D.; Trunfio, P.; Marozzo, F. *Data Analysis in the Cloud*; Elsevier, 2015. ISBN 978-0-12-802881-0.
- 526 2. Da Costa, G.; Fahringer, T.; Rico-Gallego, J.A.; Grasso, I.; Hristov, A.; Karatza, H.D.; Lastovetsky, A.;
527 Marozzo, F.; Petcu, D.; Stavrinides, G.L.; Talia, D.; Trunfio, P.; Astsatryan, H. Exascale machines require
528 new programming paradigms and runtimes. *Supercomputing Frontiers and Innovations* **2015**, *2*, 6–27.
- 529 3. Li, M.; Tan, J.; Wang, Y.; Zhang, L.; Salapura, V. SparkBench: A Comprehensive Benchmarking Suite
530 for in Memory Data Analytic Platform Spark. Proceedings of the 12th ACM International Conference
531 on Computing Frontiers; Association for Computing Machinery: New York, NY, USA, 2015; CF '15.
532 doi:10.1145/2742854.2747283.
- 533 4. De Oliveira, D.C.; Liu, J.; Pacitti, E. Data-intensive workflow management: for clouds and data-intensive
534 and scalable computing environments. *Synthesis Lectures on Data Management* **2019**, *14*, 1–179.
- 535 5. Verma, A.; Mansuri, A.H.; Jain, N. Big data management processing with Hadoop MapReduce and spark
536 technology: A comparison. 2016 Symposium on Colossal Data Analysis and Networking (CDAN), 2016,
537 pp. 1–4. doi:10.1109/CDAN.2016.7570891.
- 538 6. Zaharia, M.; Chowdhury, M.; Das, T.; Dave, A.; Ma, J.; McCauly, M.; Franklin, M.J.; Shenker, S.; Stoica,
539 I. Resilient distributed datasets: A fault-tolerant abstraction for in-memory cluster computing. 9th
540 {USENIX} Symposium on Networked Systems Design and Implementation ({NSDI} 12), 2012, pp. 15–28.
- 541 7. Samadi, Y.; Zbakh, M.; Tadonki, C. Performance comparison between Hadoop and Spark frameworks
542 using HiBench benchmarks. *Concurrency and Computation: Practice and Experience* **2018**, *30*, e4367.
- 543 8. Delimitrou, C.; Kozyrakis, C. Quasar: Resource-Efficient and QoS-Aware Cluster Management.
544 Proceedings of the 19th International Conference on Architectural Support for Programming Languages
545 and Operating Systems; Association for Computing Machinery: New York, NY, USA, 2014; ASPLOS '14, p.
546 127–144. doi:10.1145/2541940.2541941.
- 547 9. Llull, Q.; Fan, S.; Zahedi, S.M.; Lee, B.C. Cooper: Task Colocation with Cooperative Games. 2017
548 IEEE International Symposium on High Performance Computer Architecture (HPCA), 2017, pp. 421–432.
549 doi:10.1109/HPCA.2017.22.
- 550 10. Marco, V.S.; Taylor, B.; Porter, B.; Wang, Z. Improving Spark Application Throughput via Memory
551 Aware Task Co-Location: A Mixture of Experts Approach. Proceedings of the 18th ACM/IFIP/USENIX
552 Middleware Conference; Association for Computing Machinery: New York, NY, USA, 2017; Middleware
553 '17, p. 95–108. doi:10.1145/3135974.3135984.
- 554 11. Maros, A.; Murai, F.; Couto da Silva, A.P.; M. Almeida, J.; Lattuada, M.; Gianniti, E.; Hosseini, M.; Ardagna,
555 D. Machine Learning for Performance Prediction of Spark Cloud Applications. 2019 IEEE 12th International
556 Conference on Cloud Computing (CLOUD), 2019, pp. 99–106. doi:10.1109/CLOUD.2019.00028.

- 557 12. Talia, D. Workflow Systems for Science: Concepts and Tools. *International Scholarly Research Notices* **2013**,
558 2013, 1–15.
- 559 13. Smanchat, S.; Viriyapant, K. Taxonomies of workflow scheduling problem and techniques in the cloud.
560 *Future Generation Computer Systems* **2015**, *52*, 1–12. Special Section: Cloud Computing: Security, Privacy
561 and Practice, doi:<https://doi.org/10.1016/j.future.2015.04.019>.
- 562 14. Bittencourt, L.F.; Madeira, E.R.; S. Da Fonseca, N.L. Scheduling in hybrid clouds. *IEEE Communications*
563 *Magazine* **2012**, *50*, 42–47. doi:10.1109/MCOM.2012.6295710.
- 564 15. Zhao, Y.; Hu, F.; Chen, H. An adaptive tuning strategy on spark based on in-memory computation
565 characteristics. 2016 18th International Conference on Advanced Communication Technology (ICACT),
566 2016, pp. 484–488. doi:10.1109/ICACT.2016.7423442.
- 567 16. Chen, D.; Chen, H.; Jiang, Z.; Zhao, Y. An adaptive memory tuning strategy with high performance for
568 Spark. *International Journal of Big Data Intelligence* **2017**, *4*, 276–286.
- 569 17. Xuan, P.; Luo, F.; Ge, R.; Srimani, P.K. Dynamic Management of In-Memory Storage for Efficiently
570 Integrating Compute-and Data-Intensive Computing on HPC Systems. 2017 17th IEEE/ACM
571 International Symposium on Cluster, Cloud and Grid Computing (CCGRID), 2017, pp. 549–558.
572 doi:10.1109/CCGRID.2017.66.
- 573 18. Tang, Z.; Zeng, A.; Zhang, X.; Yang, L.; Li, K. Dynamic memory-aware scheduling in
574 spark computing environment. *Journal of Parallel and Distributed Computing* **2020**, *141*, 10 – 22.
575 doi:<https://doi.org/10.1016/j.jpdc.2020.03.010>.
- 576 19. Bae, J.; Jang, H.; Jin, W.; Heo, J.; Jang, J.; Hwang, J.; Cho, S.; Lee, J.W. Jointly optimizing task granularity
577 and concurrency for in-memory mapreduce frameworks. 2017 IEEE International Conference on Big Data
578 (Big Data), 2017, pp. 130–140. doi:10.1109/BigData.2017.8257921.
- 579 20. Wolpert, D.H. Stacked generalization. *Neural networks* **1992**, *5*, 241–259.
- 580 21. Liu, J.; Pacitti, E.; Valduriez, P.; Mattoso, M. A Survey of Data-Intensive Scientific Workflow Management.
581 *Journal of Grid Computing* **2015**, *13*, 457–493. doi:10.1007/s10723-015-9329-8.
- 582 22. Raj, P.H.; Kumar, P.R.; Jelciana, P. Load Balancing in Mobile Cloud Computing using Bin Packing’s First
583 Fit Decreasing Method. International Conference on Computational Intelligence in Information System.
584 Springer, 2018, pp. 97–106.
- 585 23. Baker, T.; Aldawsari, B.; Asim, M.; Tawfik, H.; Maamar, Z.; Buyya, R. Cloud-SEnergy: A bin-packing based
586 multi-cloud service broker for energy efficient composition and execution of data-intensive applications.
587 *Sustainable Computing: informatics and systems* **2018**, *19*, 242–252.
- 588 24. Stavrinides, G.L.; Karatza, H.D. Scheduling real-time DAGs in heterogeneous clusters by combining
589 imprecise computations and bin packing techniques for the exploitation of schedule holes. *Future*
590 *Generation Computer Systems* **2012**, *28*, 977–988.
- 591 25. Coffman, Jr, E.G.; Garey, M.R.; Johnson, D.S. An application of bin-packing to multiprocessor scheduling.
592 *SIAM Journal on Computing* **1978**, *7*, 1–17.
- 593 26. Darapuneni, Y.J. A Survey of Classical and Recent Results in Bin Packing Problem **2012**.
- 594 27. Marozzo, F.; Rodrigo Duro, F.; Garcia Blas, J.; Carretero, J.; Talia, D.; Trunfio, P. A data-aware
595 scheduling strategy for workflow execution in clouds. *Concurrency and Computation: Practice and*
596 *Experience* **2017**, *29*, e4229, [<https://onlinelibrary.wiley.com/doi/pdf/10.1002/cpe.4229>]. e4229 cpe.4229,
597 doi:<https://doi.org/10.1002/cpe.4229>.
- 598 28. Marozzo, F.; Talia, D.; Trunfio, P. A Workflow Management System for Scalable Data Mining on Clouds.
599 *IEEE Transactions On Services Computing* **2018**, *11*, 480–492. ISSN: 1939-1374.
- 600 29. Aseeri, A.O.; Zhuang, Y.; Alkathairi, M.S. A Machine Learning-Based Security Vulnerability Study on XOR
601 PUFs for Resource-Constraint Internet of Things. 2018 IEEE International Congress on Internet of Things
602 (ICIOT), 2018, pp. 49–56. doi:10.1109/ICIOT.2018.00014.

603 **Publisher’s Note:** MDPI stays neutral with regard to jurisdictional claims in published maps and institutional
604 affiliations.

605 © 2021 by the authors. Submitted to *Future Internet* for possible open access publication
606 under the terms and conditions of the Creative Commons Attribution (CC BY) license
607 (<http://creativecommons.org/licenses/by/4.0/>).