

Infrastructures for High-Performance Computing: Cloud Computing Development

Environments

Cristian Cosentino

c/o DIMES – University of Calabria

Via P. Bucci 41c

87036 Rende (CS)

Italy

Email: cristian.cosentino@dimes.unical.it

Telephone: +39 0984 494782

Fabrizio Marozzo (corresponding author)

c/o DIMES – University of Calabria

Via P. Bucci 42c

87036 Rende (CS)

Italy

Email: fmarozzo@dimes.unical.it

Telephone: +39 0984 494782

Paolo Trunfio

c/o DIMES – University of Calabria

Via P. Bucci 41c

87036 Rende (CS)

Italy

Email: paolo.trunfio@unical.it

Telephone: +39 0984 494788

Keywords

Apache Airflow, Apache Beam, Atom, Azure ML, BigML, Data-analytics development environments, DMCF, Eclipse, Google Colab, Hadoop, Integrated development environments, IntelliJ, NetBeans, Parallel-processing development environments, PyCharm, Spark, Swift, Visual Studio and Workflow development environments

Abstract

This chapter describes some of the most representative cloud computing development environments classified into four types. *Integrated development environments* are used to code,

debug, deploy and monitor cloud applications that are executed on a cloud infrastructure. *Parallel-processing development environments* are used to define parallel applications for processing large amounts of data that run on a cluster of virtual machines provided by a cloud infrastructure. *Workflow development environments* are used to define workflow-based applications that are executed on a cloud infrastructure. *Data-analytics development environments* are used to define data analysis applications through machine learning and data mining tools provided by a cloud infrastructure.

Key points

- This chapter provides an overview of the main cloud computing development environments in use.
- The content of this chapter discusses the main Integrated Development Environments (IDEs) used to simplify the development, deployment, and monitor of cloud applications.

Introduction

Developing cloud applications may be a complex task, with specific issues that go beyond those of stand-alone application programming. For instance, cloud programming must deal with deployment, scalability and monitoring aspects that are not easy to handle without the use of ad-hoc environments (Talía, Trunfio and Marozzo, 2015). In fact, to simplify the development of cloud applications, cloud computing development environments are often used. This article describes some of the most representative cloud computing development environments currently in use. The environment presented in this paper are classified into four types:

- *Integrated development environments*, which are used to code, debug, deploy and monitor cloud applications that are executed on a cloud infrastructure. The environments discussed in this

article are *IntelliJ, Visual Studio, Eclipse, NetBeans, PyCharm and Atom*.

- *Parallel-processing development environments*, which are used to define parallel applications for processing large amounts of data that are run on a cluster of virtual machines provided by a cloud infrastructure. The environments presented here are *Hadoop* and *Spark*.
- *Workflow development environments*, which are used to define workflow-based applications that are executed on a cloud infrastructure. The examples discussed here are *Swift, DMCF, Apache Airflow and Apache Beam*.
- *Data-analytics development environments*, which are used to define data analysis applications through machine learning and data mining tools provided by a cloud infrastructure. The examples presented in this article are *Azure ML, BigML and Google Colab*.

Integrated Development Environments

Eclipse

Eclipse (see Section Relevant Websites) is one of the most popular integrated development environments (IDEs) for software programmers that can be used to define applications in C++, Java, JavaScript, PHP, Python, R and so on, which can be run and deployed on multiple operating systems and computing platforms, including the most popular cloud computing infrastructures.

The Eclipse platform can be extended by installing plug-ins, such as development toolkits for novel programming languages and/or systems. Plug-ins can be programmed using Eclipse APIs and can be run on any of the supported operating systems. At the core of Eclipse is an architecture for discovering, loading, and running plug-ins. In addition to providing a development environment for programming languages, Eclipse supports development for most popular application servers

(e.g., Tomcat, GlassFish) and is often capable of installing the required server directly from the IDE. It supports remote debugging that allows programmers to debug the code of applications running on servers.

Eclipse provides three types of products:

- *Desktop IDE*, for defining and running Java applications, C/C++ software, PHP web pages and so on in a desktop PC;
- *Cloud IDEs*, a Cloud IDE to develop software using a browser;
- *IDE Platforms*, a set of frameworks and common services to support the use of Eclipse as a component model.

Eclipse allows to program cloud applications for the main public and private cloud infrastructures. For example, *AWS Toolkit for Eclipse* is an open source plug-in that allows developers to define, debug, and deploy Java applications on Amazon Web Services. *IBM Eclipse Tools for Bluemix* enables the deployment and integration of many services from Bluemix into applications. Finally, *Google Cloud Tools for Eclipse* simplifies the development of web applications that utilize Google cloud technology.

Visual Studio

Microsoft Visual Studio (see Section Relevant Websites) is a Microsoft IDE that allows users to develop, test, and deploy applications for the web, desktop, cloud, mobile, and game consoles. It is fully integrated with Microsoft technologies such as Windows API, Microsoft Office, Microsoft Azure and Microsoft Store.

Visual Studio includes a code editor for supporting code completion and refactoring. An integrated debugger helps to observe the run-time behavior of programs and find problems. The functionality

of the IDE can be enhanced through plug-ins, such as visual tools aiding in the development of GUI of web, desktop and mobile applications.

Visual Studio supports the most popular programming languages. For example, C++ for performance across a wide range of devices; Python for cross-platform scripting; R, for data processing; Node.js for scalable applications in JavaScript; C# as a multi-paradigm programming language for a variety of platforms, including the cloud.

Visual Studio allows the programming of cloud applications for Microsoft Azure and other cloud infrastructures. For example, *Visual Studio Tools for Azure* allows building, managing, and deploying cloud applications on Azure, whilst the *AWS Toolkit for Visual Studio* is a plugin that permits to develop, debug, and deploy .NET applications that use Amazon Web Services.

IntelliJ

IntelliJ (see Section Relevant Websites) is an IDE for web, mobile, cloud and enterprise development that supports languages like Java, JavaScript, Groovy and Scala. It is developed by the JetBrains company and is available in an Apache Licensed community edition, and in a commercial edition. The community edition allows Java and Android development. The commercial edition extends the community edition with support for web and enterprise development. Both community and commercial editions support cloud development with sponsorship of the main cloud providers.

One important feature of IntelliJ is code completion made by analyzing the source code of a user. Specifically, IntelliJ indexes user source code, for providing relevant code suggestions and on-the-fly code analysis. As the previous two systems discussed above, IntelliJ supports plugins for adding additional functionality to the IDE, such as version control systems (e.g., GIT), databases

(e.g., Microsoft SQL Server) and automatic task runners (e.g., Grunt).

IntelliJ IDEA supports the most popular Java application servers, such as Tomcat, JBoss and Glassfish. A developer can deploy, debug and monitor an application onto an application server. Moreover, IntelliJ IDEA provides dedicated plug-ins that allows programmers to manage Docker virtual machines.

IntelliJ allows developers to create and interact with cloud applications using the API of the most popular cloud infrastructures. For example, through *AWS Manager Plugin* it provides integration with AWS services like EC2, RDS and S3, whilst *Cloud Tools for IntelliJ* is a Google-sponsored plugin that allows IntelliJ developers to interact with Google Cloud Platform services.

Atom

Atom is an open-source IDE developed by GitHub, which supports different programming languages like JavaScript, HTML, CSS, Python and many more. Atom is highly customizable thanks to a wide range of packages available in its library. Additionally, Atom provides an intuitive and easy-to-use coding experience for software developers.

Atom is often used as an integrated development environment for building cloud applications. It comes equipped with debugging and performance profiling tools that help identify and fix coding issues. Atom supports integration with cloud services such as AWS, Heroku, and Azure, making deploying clouds on these platforms easy. Some of Atom's features include:

- *Customizable user interface*, which allows users to customize the visual interface according to their preferences, choosing from different layout options, themes, and extension packages;
- *Git integration*, offers native integration with Git, the most popular source code versioning system, allowing to easily manage code changes and collaborate with other team;

- *Extension packages*, has a vast ecosystem of extension packages that add additional functionality to the editor. These packages can be installed directly from the editor and customized to meet each user's specific needs;
- *Advanced search functions*, offers advanced search functions, including text search and regular expression search, to help quickly and easily find specific parts of code;

Atom can be a great choice for the development of cloud applications for public and private cloud infrastructures, thanks to its advanced features for code versioning and its integration with hosting platforms such as GitHub. Moreover, it is highly customizable thanks to its platform's wide range of packages and themes, making it a flexible and extensible editor.

NetBeans

NetBeans is a multilingual open-source integrated development environment that supports desktop, mobile and web development. One of its main features is its extensibility, which allows the use of different plugins. NetBeans supports integration with various cloud computing services, including Amazon Web Services (AWS) and Oracle Cloud. By integrating with these cloud services, NetBeans allows programmers to build, test and deploy cloud applications directly from the IDE. Specifically, it provides a plugin for AWS that allows developers to create and manage AWS resources directly from the IDE, including Amazon EC2, Amazon S3, and Amazon RDS. In addition, the AWS plugin provides tools for managing access keys and configuring security rules. It allows easy integration with Oracle Cloud, it provides a specific plugin that allows you to use the Oracle Cloud Infrastructure directly from the IDE. The plugin allows programmers to create and manage Oracle Cloud resources, such as virtual machine instances and databases, and

to deploy your cloud applications on Oracle Cloud. NetBeans also offers features for developing Java-based cloud applications, such as integration with the JavaServer Faces (JSF) framework and the ability to create and manage Java EE web applications.

In general, NetBeans is a highly customizable and versatile IDE that supports the development of cloud applications in multiple programming languages and integration with various cloud platforms.

PyCharm

PyCharm is an integrated development environment specifically designed to help developers build and manage Python projects. In terms of cloud development, PyCharm offers a range of features that make it an excellent choice for developing and deploying Python applications on cloud platforms. Here are some of its key features:

- *Remote Development*, offers remote development capabilities that allow developers to work on code that's stored on remote servers or in the cloud;
- *Integration with Cloud Platforms*, provides out-of-the-box integration with major cloud platforms, such as Amazon Web Services, Microsoft Azure, and Google Cloud Platform. This allows you to easily deploy and manage Python applications on cloud infrastructure;
- *Cloud Debugging*, includes advanced debugging features that are specifically designed for cloud development. With PyCharm, can debug applications running on remote servers or in cloud environments, and can even debug distributed systems with multiple components.
- *Testing and Continuous Integration*, supports a range of testing frameworks, such as pytest and unittest, and integrates with popular continuous integration tools, such as Jenkins and Travis CI. This makes it easy to create and run automated tests and ensure that their code is always ready for

deployment;

- *Support for serverless architectures*, PyCharm supports the development of serverless applications on AWS, GCP, and Azure. With PyCharm, can easily create and deploy serverless functions and applications that automatically scale to meet demand;

PyCharm with its cloud integration, remote development capabilities, serverless support, and advanced debugging and testing features makes it a powerful and flexible IDE that provides a range of features and tools to support cloud-based Python development.

Parallel-Processing Development Environments

Hadoop

Apache Hadoop (see Section Relevant Websites) is commonly used to develop parallel applications that analyze big amounts of data. It can be adopted for developing parallel applications using many programming languages (e.g., Java, Ruby, Python, C++) based on the MapReduce programming model (Dean and Ghemawat, 2008) on a cluster or on a cloud platform. Hadoop relieves developers from having to deal with classical distributed computing issues, such as load balancing, fault tolerance, data locality, and network bandwidth saving. The Hadoop project is not only about the MapReduce programming model (Hadoop MapReduce module), as it includes other modules such as:

- *Hadoop Distributed File System (HDFS)*: a distributed file system providing fault tolerance with automatic recovery, portability across heterogeneous commodity hardware and operating systems, high-throughput access and data reliability.
- *Hadoop YARN*: a framework for cluster resource management and job scheduling.
- *Hadoop Common*: common utilities that support the other Hadoop modules.

With the introduction of YARN in 2013, Hadoop turned from a batch processing solution into a platform for running a large variety of data applications, such as streaming, in-memory, and graphs analysis. As a result, Hadoop became a reference for several other frameworks, such as: Giraph (see Section Relevant Websites) for graph analysis; Storm (see Section Relevant Websites), for streaming data analysis; Hive (see Section Relevant Websites), which is a data warehouse software for querying and managing large datasets; Pig (see Section Relevant Websites), which is as a dataflow language for exploring large datasets; Tez (see Section Relevant Websites) for executing complex directed-acyclic graph of data processing tasks; Oozie (see Section Relevant Websites), which is a workflow scheduler system for managing Hadoop jobs.

Hadoop is available in most cloud infrastructures. For example, the *HDInsight* service by Microsoft Azure, the *Amazon Elastic MapReduce (EMR)* service by AWS, and Google Cloud Dataproc service by Google Cloud.

Spark

Apache Spark (see Section Relevant Websites) is an open-source framework for in-memory data analysis and machine learning developed at UC Berkeley in 2009. It can process distributed data from several sources, such as HDFS, HBase, Cassandra, and Hive. It has been designed to efficiently perform both batch processing applications (similar to MapReduce) and dynamic applications like streaming, interactive queries, and graph analysis (Belcastro, Cantini, Marozzo et al. (2022)). Spark is compatible with Hadoop data and it can run in Hadoop clusters through the YARN module. However, in contrast to Hadoop's two-stage MapReduce paradigm in which intermediate data are always stored in distributed file systems, Spark stores data in a cluster's memory and queries it repeatedly so as to obtain better performance for several classes of

applications (e.g., interactive jobs, real-time queries, and stream data) (Xin, Rosen, Zaharia *et al.*, 2013). The Spark project has different components:

- *Spark Core* contains the basic functionalities of the library such as for manipulating collections of data, memory management, interaction with distributed file systems, task scheduling, and fault recovery.
- *Spark SQL* provides an API to query and manipulate structured data using standard SQL or Apache Hive variant of SQL.
- *Spark Streaming* provides an API for manipulating streams of data.
- *GraphX* is a library for manipulating and analyzing big graphs.
- *MLlib* is a scalable machine learning library on top of Spark that implements many common machine learning and statistical algorithms.

Several big companies and organizations use Spark for big data analysis purposes: for example, Ebay uses Spark for log transaction aggregation and analytics, Kelkoo for product recommendations, SK Telecom analyzes mobile usage patterns of customers.

Similarly to Hadoop, most cloud infrastructures provide Spark as a service, like *IBM Analytics for Apache Spark*, *Azure HDInsight* and *Google Cloud Dataproc*.

Workflow Development Environments

Swift

Swift (Wilde, Hategan, Wozniak *et al.*, 2011) is an implicitly parallel scripting language that runs workflows across several distributed systems, like supercomputers, clusters, grids and clouds. The Swift language has been designed at the University of Chicago and at the Argonne National Lab

to provide users with a workflow-based language for cloud computing.

Swift separates the application workflow logic from runtime configuration. This approach allows a flexible development model. The Swift language allows invocation and running of external application code and allows binding with application execution environments without extra coding from the user. Swift/K is the previous version of the Swift language that runs on the Karajan grid workflow engine across wide area resources. Swift/T is a new implementation of the Swift language for high-performance computing. In this implementation, a Swift program is translated into an MPI program that uses the Turbine and ADLB runtime libraries for scalable dataflow processing over MPI. The Swift-Turbine Compiler (STC) is an optimizing compiler for Swift/T and the Swift Turbine runtime is a distributed engine that maps the load of Swift workflow tasks across multiple computing nodes. Users can also use Galaxy (Giardine, Riemer, Hardison *et al.*, 2005) to provide a visual interface for Swift.

The Swift language provides a functional programming paradigm where workflows are designed as a set of code invocations with their associated command-line arguments and input and output files. Swift is based on a C-like syntax and uses an implicit data-driven task parallelism. In fact, it looks like a sequential language, but being a dataflow language, all variables are futures, thus execution is based on data availability. When input data is ready, functions are executed in parallel. Moreover, parallelism can be exploited through the use of the *for each* statement. The Turbine runtime comprises a set of services that implement the parallel execution of Swift scripts exploiting the maximal concurrency permitted by data dependencies within a script and by external resource availability. Swift has been used for developing several scientific data analysis applications, such as prediction of protein structures, modeling the molecular structure of new materials, and decision making in climate and energy policy.

Data Mining Cloud Framework

The Data Mining Cloud Framework (DMCF) is a software system developed at the University of Calabria for designing and executing data analysis workflows on clouds (Marozzo, Talia and Trunfio (2018)). A Web-based user interface allows users to compose their applications and submit them for execution over cloud resources, according to a Software-as-a-Service (SaaS) approach.

The DMCF architecture has been designed to be deployed on different cloud settings. Currently, there are two different deployments of DMCF: (1) on top of a Platform-as-a-Service (PaaS) cloud, i.e., using storage, compute, and network APIs that hide the underlying infrastructure layer; (2) on top of an Infrastructure-as-a-Service (IaaS) cloud, i.e., using virtual machine images (VMs) that are deployed on the infrastructure layer.

The DMCF software modules can be grouped into *web components* and *compute components*. DMCF allows users to compose, check, and run data analysis workflows through a HTML5 web editor. The workflows can be defined using two languages: Visual Language for Cloud (VL4Cloud) (Marozzo, Duro, Blas et al. (2017)) and JavaScript for Cloud (*JS4Cloud*) (Marozzo, Talia and Trunfio (2015)). Both languages use three key abstractions:

- *Data* elements, representing input files (e.g., a dataset to be analyzed) or output files (e.g., a data mining model).
- *Tool* elements, representing software tools used to perform operations on data elements (partitioning, filtering, mining, etc.).
- *Tasks*, which represent the execution of Tool elements on given input Data elements to produce some output Data elements.

The DMCF editor generates a JSON descriptor of the workflow, specifying what are the tasks to be executed and the dependency relationships among them. The JSON workflow descriptor is

managed by the DMCF workflow engine that is in charge of executing workflow tasks on a set of workers (virtual processing nodes) provided by the cloud infrastructure. The workflow engine implements a data-drive task parallelism that assigns workflow tasks to idle workers as soon as they are ready to execute and tries to minimize data movement between different workers (Duro, Marozzo, Blas et al. 2016).

Apache Airflow

Apache Airflow is an open-source workflow platform that allows for the definition, scheduling, and monitoring of complex data processing tasks. Airflow was designed to handle workflow management in distributed data processing environments, such as those based in the cloud. Airflow allows to define workflows through a Python-based programming language, offering a high level of modularity and reproducibility. In addition, it permits interrupt the execution of workflows based on a specific calendar or events, automatically managing the repetition of workflows in case of errors or failures. Airflow also provides a web interface to monitor the status of workflows and the individual activities that compose them, allowing you to check progress in real-time and intervene in case of problems. Airflow supports integration with different data processing infrastructures, such as Hadoop, Spark, and Kubernetes. In the cloud environment, Airflow is especially useful for managing distributed workflows across multiple compute nodes. It can run on cloud infrastructures such as Amazon Web Services (AWS) or Google Cloud Platform (GCP), and it is used to orchestrate data processing across multiple computing instances. Additionally, Airflow can be used to automate the management of data upload and download processes to and from a cloud infrastructure.

Apache Airflow is a very powerful and flexible workflow orchestration platform capable of

handling complex workflows in cloud-distributed data processing environments. Thanks to its modular architecture and wide range of integration tools, Airflow represents a very useful tool for automating and monitoring data processing processes on cloud infrastructures.

Apache Beam

Apache Beam is an open-source data processing framework that provides a unified programming model for batch and streaming data processing. The main feature of Beam is to provide a high-level abstraction that allows developers to write data processing logic without worrying about the underlying infrastructure. This means that developers can write code once and then run it on different data processing platforms, including Apache Flink, Apache Spark, and Google Cloud Dataflow. Beam provides a Java and Python API for writing data processing applications. Beam applications are composed of a set of data transformations that define the processing logic. A transformation is a function that receives a set of data as input, processes the data, and produces a set of data as output. Transformations can be chained together to create a data processing pipeline. Apache Beam can be used to define workflow-based applications that run on a cloud infrastructure. It provides a higher level abstraction on the underlying data processing platforms, which means high portability of the code that involves running on different platforms without having to modify it. In addition, Beam offers greater flexibility in choosing the data processing platform.

Data-Analytics Development Environments

Azure Machine Learning

Azure Machine Learning (see Section Relevant Websites) is a SaaS that provides a Web-based machine learning environment for the creation and automation of machine learning workflows.

Through its user-friendly interface, data scientists and developers can perform several common data analysis/mining tasks on their data and automate their workflows.

Using its drag-and-drop interface, users can import their data in the environment or use special readers to retrieve data from several sources, such as Web URL (HTTP), OData Web service, Azure Blob Storage, Azure SQL Database, Azure Table. After that, users can compose their data analysis workflows where each data processing task is represented as a block that can be connected with each other through direct edges, establishing specific dependency relationships among them. Azure ML includes a rich catalog of processing tool that can be easily included in a workflow to prepare/transform data or to mine data through supervised learning (regression e classification) or unsupervised learning (clustering) algorithms. Optionally, users can include their own custom scripts (e.g., in R or Python) to extend the tools catalog. When workflows are correctly defined, users can evaluate them using some testing dataset. Users can easily visualize the results of the tests and find very useful information about models accuracy, precision and recall. Finally, in order to use their models to predict new data or perform real time predictions, users can expose them as Web services. Always through a Web-based interface, users can monitor the Web services load and use by time.

Azure Machine Learning is a fully managed service provided by Microsoft on its Azure platform; users do not need to buy any hardware/software nor manage virtual machines manually. One of the main advantages of working with Azure Machine Learning is its auto-scaling feature: models are deployed as elastic Web services so users do not have to worry about scaling them if the models' usage increases.

BigML

BigML (see Section Relevant Websites) is provided as a Software-as-a-Service (SaaS) for

discovering predictive models from data sources and using data classification and regression algorithms. The distinctive feature of BigML is that predictive models are presented to users as interactive decision trees. The decision trees can be dynamically visualized and explored within the BigML interface, downloaded for local usage and/or integration with applications, services, and other data analysis tools. Extracting and using predictive models in BigML consists in multiple steps, as detailed in the following:

- *Data source setting and dataset creation.* A data source is the raw data from which a user wants to extract a predictive model. Each data source instance is described by a set of columns, each one representing an instance feature, or field. One of the fields is considered as the feature to be predicted. A dataset is created as a structured version of a data source in which each field has been processed and serialized according to its type (numeric, categorical, etc.).
- *Model extraction and visualization.* Given a dataset, the system generates the number of predictive models specified by the user, who can also choose the level of parallelism for the task. The interface provides a visual tree representation of each predictive model, allowing users to adjust the support and confidence values and to observe in real time how these values influence the model.
- *Prediction making.* A model can be used individually, or in a group (the so-called ensemble, composed of multiple models extracted from different parts of a dataset), to make predictions on new data. The system provides interactive forms to submit a predictive query for new data using the input fields from a model or ensemble. The system provides APIs to automate the generation of predictions, which is particularly useful when the number of input fields is high.
- *Model evaluation.* BigML provides functionalities to evaluate the goodness of the predictive models extracted. This is done by generating performance measures that can be applied to the

kind of extracted model (classification or regression).

Google Colab

Google Colab is a cloud-based Jupyter notebook environment that allows users to write, run, and share Python code in a browser. It is commonly used as a development environment for data analysis and machine learning applications, as it provides users with access to powerful computing resources and a range of data processing and analysis tools. Google Colab is a cloud-based platform that enables massive data processing, which means users can leverage the computing power of Google's cloud infrastructure without having to set up and manage their own hardware. Its features make it easy to use for both individuals and teams to work that will be able to easily manage large datasets and complex machine learning models. Additionally, Google Colab offers a range of built-in data analysis tools and libraries, including NumPy, Pandas, and Matplotlib, as well as access to powerful machine learning frameworks like TensorFlow and PyTorch. Google Colab is a powerful tool for developing data analytics and machine learning applications in a cloud-based environment. Its ease of use, powerful processing resources, and integrated data processing and analysis tools make it a popular choice for individuals and teams working on data-driven projects.

Closing Remarks

This article described four categories of cloud computing development environments. *Integrated development environments* are used to code, debug, deploy and monitor cloud applications that are executed on a cloud infrastructure. *Parallel-processing development environments* are used to define parallel applications for processing large amounts of data that are run on a cluster of virtual

machines provided by a cloud infrastructure. *Workflow development environments* are used to define workflow-based applications that are executed on a cloud infrastructure. *Data-analytics development environments* are used to define data analysis applications through machine learning and data mining tools provided by a cloud infrastructure. For each category, the article described some of the most representative cloud computing development environments currently in use, including their main features, provided services and available implementations.

See also

Dedicated Bioinformatics Analysis Hardware. Infrastructure for High-Performance Computing: Grids and Grid Computing. Infrastructures for High-Performance Computing: Cloud Computing. Infrastructures for High-Performance Computing: Cloud Infrastructures. MapReduce in Computational Biology via Hadoop and Spark. Models and Languages for High-Performance Computing. Text Mining Applications.

References

- Belcastro, L., Cantini, R., Marozzo, F., Orsino, A., Talia, D. and Trunfio, P. (2022). Programming Big Data Analysis: Principles and Solutions. *Journal of Big Data* **9**(4).
- Dean, J. and Ghemawat, S. (2008). MapReduce: simplified data processing on large clusters. *Communications of the ACM* **51**(1), 107-113.
- Duro, F. R., Marozzo, F., Garcia Blas, J., Talia, D., & Trunfio, P. (2016). Exploiting in-memory storage for improving workflow executions in Cloud platforms. *Journal of Supercomputing*, **72**(11), 4069-4088. ISSN: 1573-0484.
- Giardine, B., Riemer, C., Hardison, R. C., et al. (2005). Galaxy: A platform for interactive large-scale genome analysis. *Genome Research*, **15**, 1451-1455.

Marozzo, F., Talia, D., & Trunfio, P. (2018). A Workflow Management System for Scalable Data Mining on Clouds. *IEEE Transactions on Services Computing*, 11(3), 480-492. ISSN: 1939-1374.

Marozzo, F., Talia, D., & Trunfio, P. (2015). JS4Cloud: Script-based Workflow Programming for Scalable Data Analysis on Cloud Platforms. *Concurrency and Computation: Practice and Experience*, 27(17), 5214-5237. ISSN: 1532-0626.

Marozzo, F., Duro, F. R., Garcia Blas, J., Carretero, J., Talia, D. and Trunfio, P. (2017). A Data-aware Scheduling Strategy for Workflow Execution in Clouds". *Concurrency and Computation: Practice and Experience*, vol. 29, n. 24, Wiley InterScience, 2017. Note: ISSN: 1532-0634.

Talia, D., Trunfio, P. and Marozzo, F. (2015). *Data Analysis in the Cloud*. The Netherlands: Elsevier.

Wilde, M., Hategan, M., Wozniak, J. M., Clifford, B., Katz, D. S. and Foster, I. (2011). A language for distributed parallel scripting. *Parallel Computing*, 37(9), 633-652.

Wozniak, J. M., & Foster, I. (2014). Language features for scalable distributed-memory dataflow computing. In *Proceedings Data-flow Execution Models for Extreme-scale Computing at PACT*.

Xin, R. S., Rosen, J., Zaharia, M., et al. (2013). Shark: SQL and rich analytics at scale. In *Proceedings of the 2013 ACM SIGMOD International Conference on Management of Data (SIGMOD'13)*, New York, NY.

Relevant Websites

<http://airflow.apache.org/> – Apache Airflow: The Apache Software Foundation.

<http://beam.apache.org/> – Apache Beam: The Apache Software Foundation.

<http://giraph.apache.org/> – Apache Giraph: The Apache Software Foundation.

<http://hadoop.apache.org/> – Apache Hadoop: The Apache Software Foundation.

<http://hive.apache.org> – Apache Hive: The Apache Software Foundation.

<http://pig.apache.org> – Apache Pig: The Apache Software Foundation.

<http://spark.apache.org> – Apache Spark: The Apache Software Foundation.

<http://storm.apache.org> – Apache Storm: The Apache Software Foundation.

<https://github.com/atom> – Atom.

<https://bigml.com> – BigML, Inc.

<https://eclipse.org/> – Eclipse.

<https://colab.research.google.com/> – Google Colab.

<https://www.jetbrains.com/idea/> – IntelliJ.

<https://netbeans.apache.org/> – Apache NetBeans: The Apache Software Foundation.

<https://azure.microsoft.com/en-us/services/machine-learning/> – Microsoft Azure.

<https://www.visualstudio.com> – Microsoft Visual Studio.

<http://oozie.apache.org/> – Oozie: The Apache Software Foundation.

<https://www.jetbrains.com/pycharm/> – PyCharm.

<http://tez.apache.org> – TEZ: The Apache Software Foundation.

Online Contents