

DataCloud: Enabling the Big Data Pipelines on the Computing Continuum

Dumitru Roman¹, Nikolay Nikolov¹, Brian Elvester¹, Ahmet Soylu², Radu Prodan³, Dragi Kimovski³, Andrea Marrella⁴, Francesco Leotta⁴, Dario Benvenuti⁴, Mihhail Matskin⁵, Giannis Ledakis⁶, Anthony Simonet-Boulogne⁷, Fernando Perales⁸, Evgeny Kharlamov⁹, Alexandre Ulisses¹⁰, Arnor Solberg¹¹, and Raffaele Ceccarelli¹²

¹ SINTEF AS, Norway

² Oslo Metropolitan University, Norway

³ University of Klagenfurt, Austria

⁴ Sapienza University of Rome, Italy

⁵ KTH Royal Institute of Technology, Sweden

⁶ UBITECH, Greece

⁷ iExec, France

⁸ JOT, Spain

⁹ Bosch Center for Artificial Intelligence, Germany

¹⁰ MOG, Portugal

¹¹ Tellu, Norway

¹² Ceramica Catalano, Italy

1 Summary of the project

With the recent developments of Internet of Things (IoT) and cloud-based technologies, massive amounts of data are generated by heterogeneous sources and stored through dedicated cloud solutions. Often organizations generate much more data than they are able to interpret, and current Cloud Computing technologies cannot fully meet the requirements of the Big Data processing applications and their data transfer overheads [1]. Many data are stored for compliance purposes only but not used and turned into value, thus becoming *Dark Data*, which are not only an untapped value, but also pose a risk for organizations [3].

To guarantee a better exploitation of Dark Data, the **DataCloud project**¹³ aims to realize novel methods and tools for effective and efficient management of the Big Data Pipeline lifecycle encompassing the Computing Continuum.

Big Data pipelines are composite pipelines for processing data with non-trivial properties, commonly referred to as the Vs of Big Data (e.g., volume, velocity, value, etc.) [4]. Tapping their potential is a key aspect to leverage Dark Data, although it requires to go beyond the current approaches and frameworks for Big Data processing. In this respect, the concept of *Computing Continuum*

¹³ DataCloud is a Research and Innovation project funded by the European Commission under the Horizon 2020 program (Grant number 101016835). The project runs for three years, between 2021-2023.

extends the traditional centralised Cloud Computing with Edge¹⁴ and Fog¹⁵ computing in order to ensure low latency pre-processing and filtering close to the data sources. This will prevent to overwhelm the centralised cloud data centres enabling new opportunities for supporting Big Data pipelines.

2 Objectives and Expected Results

The main objective of the project is to develop a software ecosystem for managing Big Data pipelines on the Computing Continuum. The ecosystem consists of new languages, methods and infrastructures for supporting Big Data pipelines on heterogeneous and untrusted resources. Six lifecycle phases are covered:

1. *Pipeline discovery* concerns discovering Big Data pipelines from various data sources.
2. *Pipeline definition* deals with specifying pipelines featuring an abstraction level suitable for pure data processing.
3. *Pipeline simulation* aims to evaluate the performance of individual steps to test and optimise deployments.
4. *Resource provisioning* is concerned about securely provisioning a set of (trusted and untrusted) resources.
5. *Pipeline deployment* is concerned with deployment of pipelines across the provisioned resources.
6. *Pipeline adaptation* deals with optimised run-time provisioning of computational resources.

The ecosystem separates the design-time from the run-time deployment of pipelines and complements modern serverless approaches [2]. A set of research challenges related to each pipeline phase will be tackled within the project, such as the advancement of process mining techniques to learn the structure of pipelines, the definition of proper DSLs for pipelines, novel approaches for pipeline containerisation and blockchain-based resource marketplaces, etc.

The expected impact of DataCloud is to lower the technological entry barriers for the incorporation of Big Data pipelines in organizations' workflows and make them accessible to a wider set of stakeholders regardless of the hardware infrastructure. To achieve this, DataCloud will validate its results through a strong selection of complementary business cases offered by four SMEs and a large company targeting higher mobile business revenues in smart marketing campaigns, reduced live streaming production costs of sport events, trustworthy eHealth patient data management, and reduced time to production and better analytics in Industry 4.0 manufacturing.

Acknowledgments. This work has been supported by the Horizon 2020 project DataCloud (Grant number 101016835).

¹⁴ Edge Computing is a paradigm that brings computation and data storage closer to the location where it is needed to improve response times and save bandwidth.

¹⁵ Fog Computing uses edge devices to carry out a substantial amount of computation, storage, and communication locally and routed over the Internet backbone.

References

1. Barika, M., Garg, S., Zomaya, A.Y., Wang, L., Moorsel, A.V., Ranjan, R.: Orchestrating big data analysis workflows in the cloud: research challenges, survey, and future directions. *ACM Computing Surveys (CSUR)* **52**(5), 1–41 (2019)
2. Castro, P., Ishakian, V., Muthusamy, V., Slomiński, A.: The Rise of Serverless Computing. *Communications of the ACM* **62**(12) (2019)
3. Gimpel, G.: Bringing dark data into the light: Illuminating existing IoT data lost within your organization. *Business Horizons* **63**(4), 519–530 (2020)
4. Plale, B., Kouper, I.: The centrality of data: data lifecycle and data pipelines. In: *Data Analytics for Intelligent Transportation Systems*, pp. 91–111. Elsevier (2017)