

An holistic approach for high-level programming of next-generation data-intensive applications targeting distributed heterogeneous computing environment

Emanuele Carlini, *Patrizio Dazzi*, Matteo Mordacchini

National Research Council of Italy, CNR



National Research Council of Italy



Data, data, data

- ▶ Data-intensive applications are one of the largest customers of high throughput computing environments.
- ▶ BigData hype shows that there is a widely recognized urge for processing large amount of data
 - ▶ universally considered the “new oil”, from which to extract wisdom, knowledge and, more in general, information of very different kinds
 - ▶ ranging from commercial to financial, from medical to political hints
- ▶ This extreme intrinsic richness and heterogeneity of data is paired with the extreme complexity in its storing and processing in an efficient manner.



Managing (large amount of) data is complex

- ▶ Software tools for programming applications dealing with very large amount of data is an hot research topic.
 - ▶ technological issues, economic concerns and legal constraints.
- ▶ Data is more and more decentralized and localized,
 - ▶ moving it is not always simple due to costs, performance and legal aspects.
 - ▶ Recent processing platforms characterized by high levels of dynamicity and heterogeneity
 - ▶ no a static set of resources available for the computation
 - ▶ depending on the context of a particular application and relative data
 - ▶ orchestrating a computation in such environments is complex, time-consuming, error-prone.
 - ▶ dynamic nature of users, applications and resources makes this task
 - ▶ extremely complex without the support of automatic tools.



A few approaches have been proposed so far

- ▶ Recently, many approaches, models and solutions have been proposed to tackle these issues.
- ▶ The problem has been faced from very different perspectives.
 - ▶ Creation of infrastructures enabling a seamless exploitation of very different kind of resources
 - ▶ hardware heterogeneity
 - ▶ clouds, cloudlets and edge devices
 - ▶ Other approaches focused on smart brokerage solutions
 - ▶ easing the task of finding the most suitable resources
 - ▶ depending on several factors: user and application requirements, location, etc.
 - ▶ Others on the exploitation of heterogeneous and dynamic resources
 - ▶ by leveraging proper programming model abstractions



Specific solutions are good, holistic ones are better

- ▶ More recently, both the research and industrial cloud communities are trying:
 - ▶ to conceive and develop holistic approaches aimed at providing vertical solutions
 - ▶ ease the programming of applications targeting heterogeneous environments
 - ▶ simplify the management of large computational infrastructures.
 - ▶ users specify their requested service level (e.g. a given throughput)
 - ▶ the supporting environment ensure their satisfaction by adopting the proper resources



A programming eco-system...

- ▶ We propose our vision for a programming ecosystem aimed at
 - ▶ organizing the computation of data-intensive applications in heterogeneous platform that goes beyond the state of the art
- ▶ An ecosystem whose ability is not limited to a seamless exploitation of a set of heterogeneous and distributed resources
 - ▶ able to address users' needs about data processing
 - ▶ adopting solutions providing a dynamically differentiated set of features
 - ▶ depending on the actual running environment,
 - ▶ the hosted application,
 - ▶ the user requirements.



...proposing a different perspective

- ▶ An ecosystem in which
 - ▶ the functional logic realizing the applications self-adapts with respect to the available resources
 - ▶ properly defined user requirements and the actual context.
- ▶ Aimed at going beyond the traditional application development approach
 - ▶ That is usually realized according to a well-defined process:
 - ▶ problem definition, algorithm selection, software implementation, data preparation and application deployment.
 - ▶ Accordingly to this schema, the application management support, provided by the target execution environments **is involved only in the last stage**, ensuring that the application is deployed on a proper (set of) resource(s).

No longer developers but problem-composers (or application scientists)

- ▶ Aimed at supporting the realization of next-generation data-intensive applications as a whole
 - ▶ relying to application developers mainly for the problem definition and the high-level orchestration logic.
- ▶ Application programmers are no longer requested
 - ▶ to select and implement the specific algorithms realizing the application logic
 - ▶ or to explicitly define data movements among resources.
- ▶ Developers (application scientists) assumes a new role, they define their applications as a composition of “problems”
 - ▶ for which exist different solutions, that are eventually adopted and composed
 - ▶ no longer involved in the actual implementation of (most of) their own applications, but mainly concerned with a detailed identification of the problems to solve in order to realize applications

Infrastructure Manager + Application Manager

- ▶ Composition results from a co-operative work jointly conducted by the
 - ▶ infrastructure run-time management
 - ▶ and the application manager
- ▶ Infrastructural manager is entitled of determining the best computational, network and storage resources among the available ones
- ▶ The Application manager analyses the user requirements on:
 - ▶ the expected quality on the results of the computation
 - ▶ the available “solutions” for the target “problem”
 - ▶ the amount and nature of the data to be processed



Applications eventually composed by software modules

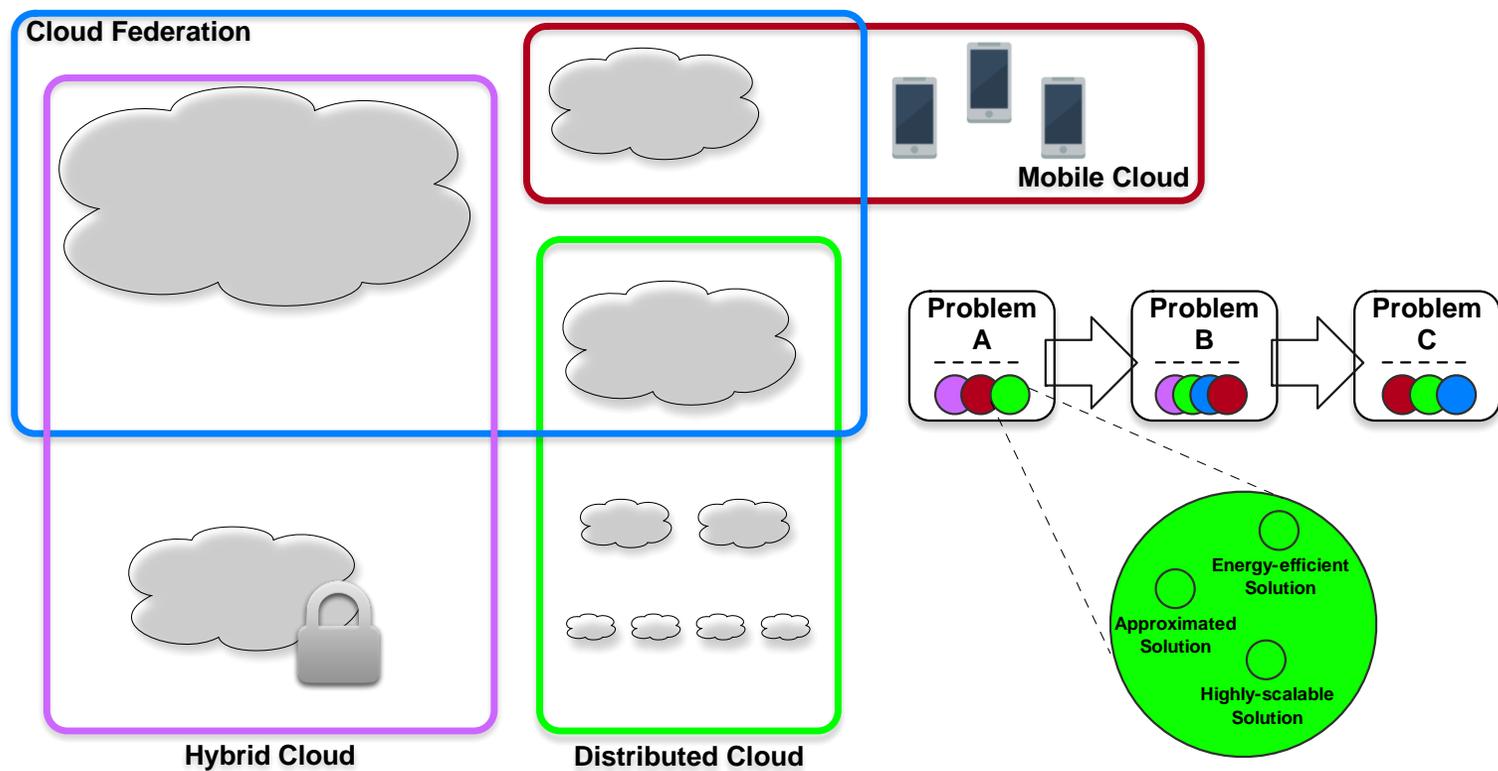
- ▶ From the interplay of such *active detection entities* is derived the actual, material, logic of the application.
- ▶ At this stage applications are effectively composed by software modules specialized w.r.t.
 - ▶ The actual context
 - ▶ the running environment
 - ▶ the available resources,
 - ▶ the input data and the quality expected by the users/customers.

Matchmaking no longer performed with resources but with solutions

- ▶ The goal for the infrastructure is no longer limited to find the most suitable resources, given an application and its requirements
- ▶ The aim is to manage applications defined as a composition of problems
 - ▶ characterized both by quantitative and qualitative requirements
 - ▶ affecting the selection of the solutions that will be adopted for their materialization
- ▶ Selection performed by an advanced execution environment, resulting from the interplay of both the application and infrastructure manager
 - ▶ by analyzing the user and application requirements, the nature and distribution of the input data, the actual context as well as the available resources,
 - ▶ taking into account the desiderata of users describing their desired trade-off on qualitative aspects (e.g., performance, cost, security, precision, etc.)
 - ▶ support decides about the solutions to instantiate to address the problems identified by the application developers.



Graphical representation of our ecosystem



Interplay with Cloud patterns

- ▶ I see a potential interesting relationship with Cloud patterns
- ▶ Yesterday's presentation at Inter-cloud cluster
 - ▶ Beniamino Di Martino
 - ▶ Agnostic patterns
- ▶ In our case non-functional requirements can have an impact on functional ones



Conclusion

- ▶ Need for a new generation of effective and efficient computing ecosystems supporting data-intensive applications
- ▶ To tackle these issues we proposed an approach assuming that developers will work at high-level
 - ▶ by identifying and defining the set of “problems” that should be solved by their applications.
- ▶ The ecosystem will provide intelligent *active detection entities*, that will be in charge of selecting the most suitable implementation for every application “problem”
 - ▶ considering the overall environment,
 - ▶ ranging from the context of the user and that of the data to be used and of the available resources,
 - ▶ leading to final deployment and orchestration of the various components of the application



Conclusion (2)

- ▶ The goal is to achieve an effective and efficient solution for the challenges posed by modern data-intensive application
 - ▶ by providing a new, simpler way to realise them
 - ▶ by allowing to seamlessly (for both the users and the developer) deploy them in the most effective way,
 - ▶ with respect to the actual context of data and computing resources.



Acknowledgements

- ▶ *This is not an activity conducted in the context of the BASMATI H2020 project*
- ▶ *Nevertheless, we hope that this vision could somehow contribute to the definition of new paradigms for the realisation of next-gen application on clouds*
 - ▶ *Maybe also to be explored in the context of BASMATI project*
- ▶ *BASMATI has received funding from the European Union's Horizon 2020 research and innovation programme under grant agreement no. 723131 and from ICT R&D program of Korean Ministry of Science, ICT and Future Planning no. R0115-16-0001.*



Questions ?



National Research Council of Italy

Backup-slide a few solutions involving different stages exist

- ▶ Some advanced systems are involved also in a few further stages
 - ▶ by providing a specific programming approach or by leveraging solutions (like Aspect-Oriented or Meta-Programming)
 - ▶ to customize at runtime non-functional features characterizing the application (e.g., logging, channel encryption, network configuration).
 - ▶ However, these systems are usually isolated components
 - ▶ interoperability is complicated or impossible to obtain
 - ▶ tools not designed to be able to adapt the *semantics* of the application to the actual context and situation

Backup slide who implement solutions ?

- ▶ Each solution solving a problem, among the one identified by application developers are selected from a collection of pre-defined solutions
- ▶ Such solutions are provided by solution designers, that are responsible for providing highly-efficient solutions, tailored to specific execution environment
- ▶ Besides their tailoring to a certain architecture/environment, the available solutions can differentiate one each other with respect to the qualitative features provided,
 - ▶ e.g., the provided level of privacy preservation, the required network bandwidth, the quality of being able to deliver either exact or approximated results