



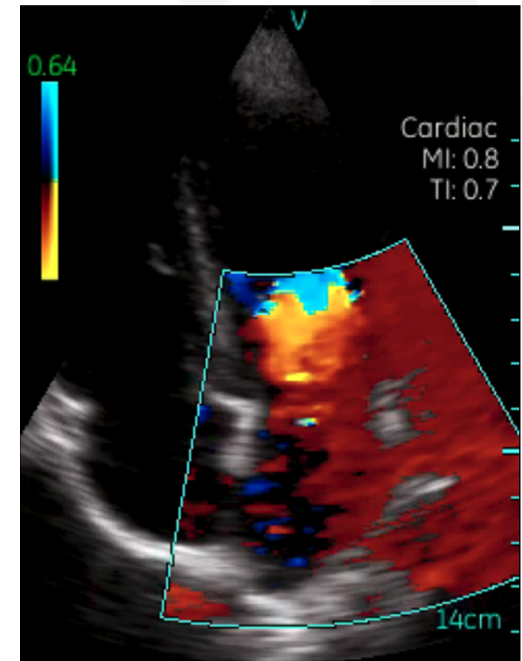
Adaptive, Trustworthy, Manageable, Orchestrated, Secure Privacy-assuring Hybrid, Ecosystem for REsilient Cloud Computing

Machine Learning Pipelines on Medical Imaging

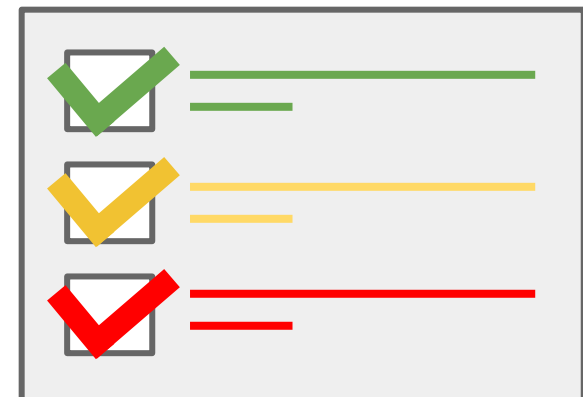
*Ignacio Blanquer (UPV), Eduardo
Camacho-Ramos, Ana Jiménez-Pastor,
Ángel Alberich-Bayarri (QUIBIM), Walter
Dos Santos, Prof. Wagner Meira Jr.
(UFMG)*



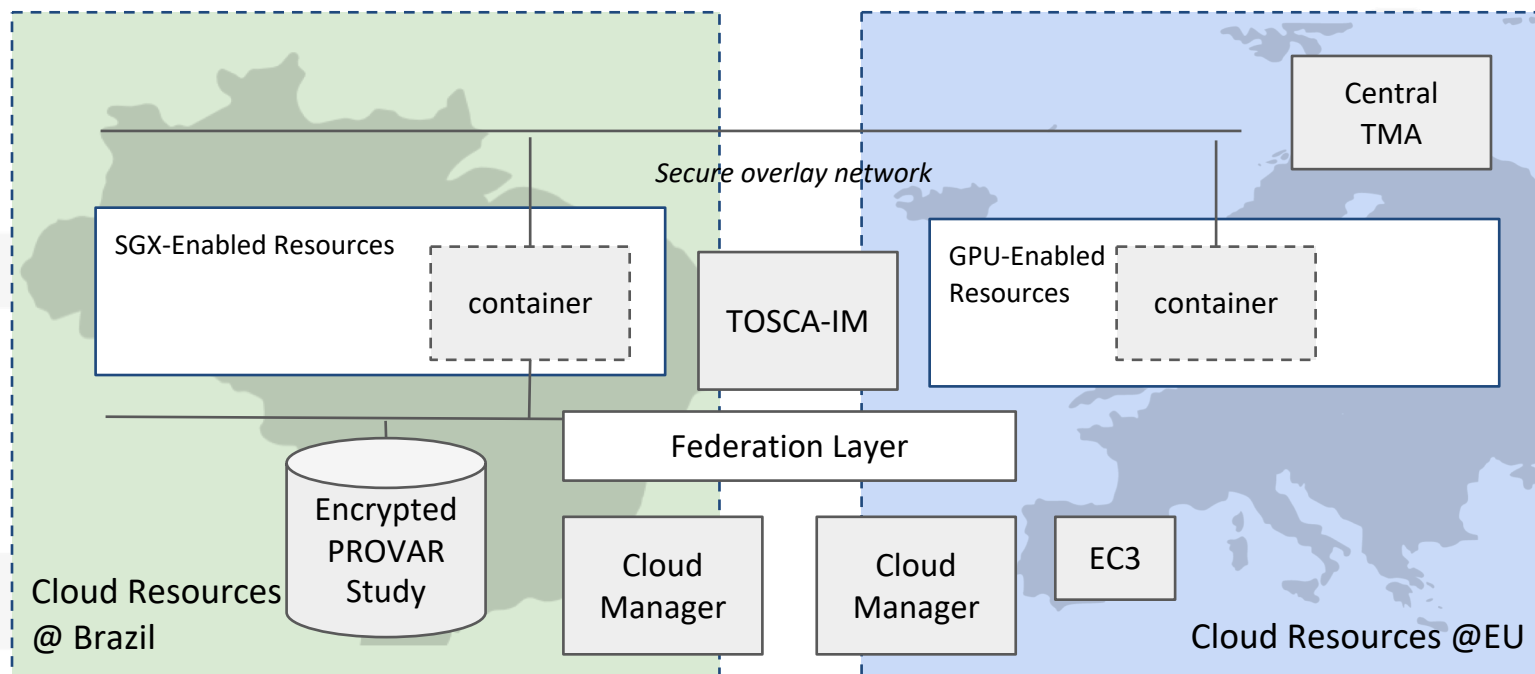
- The Rheumatic Heart Disease (RHD) is a disease that can be easily treated in its early stages, but may produce enormous damage to the heart if remains untreated, including severe sequelae and death.
- The challenge is to process a large set of medical images, along with additional metadata and clinical information, efficiently and securely, to extract features that could be used to assist and even automate diagnosis.
- Data comes from the PROVAR Echocardio data
 - 4.021 studies (4.035 Normal + 180 Borderline + 26 Definite)
 - 59.018 240×320 MP4 videos of 1-3 seconds.
 - To be classified into three categories according to the WHF criteria: Normal, Borderline and Definite RHD.
- Challenges:
 - Unbalance of the cases.
 - Noise and low quality of the echocardio images.
 - No information on the view.



- Sensitive data should remain in the Brazilian geographical boundaries and confidentiality should be preserved.
- Computing requires accelerators and may not be available within the boundaries where the sensitive data is located.
- Parallel execution should be provided.
- Repeatability and reproducibility should be a main goal.
- Flexible and dynamic environment.
- Simplified interfaces for non-ICT experts.



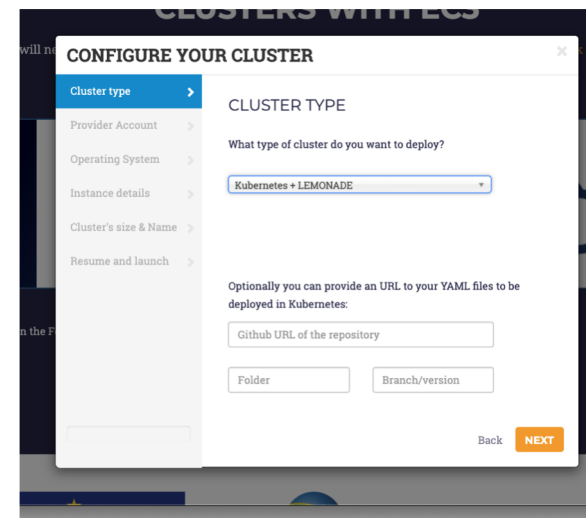
- The underlying infrastructure is a federated cloud
 - Using fogbow (www.fogbowcloud.org) on OpenStack and OpenNebula.
 - With a Federated Network to provide a coherent network space among nodes.
 - Heterogeneous resources: SGX-enabled and GPU nodes.
- Using EC3⁽¹⁾ and Infrastructure Manager⁽²⁾ to deploy a virtual infrastructure.



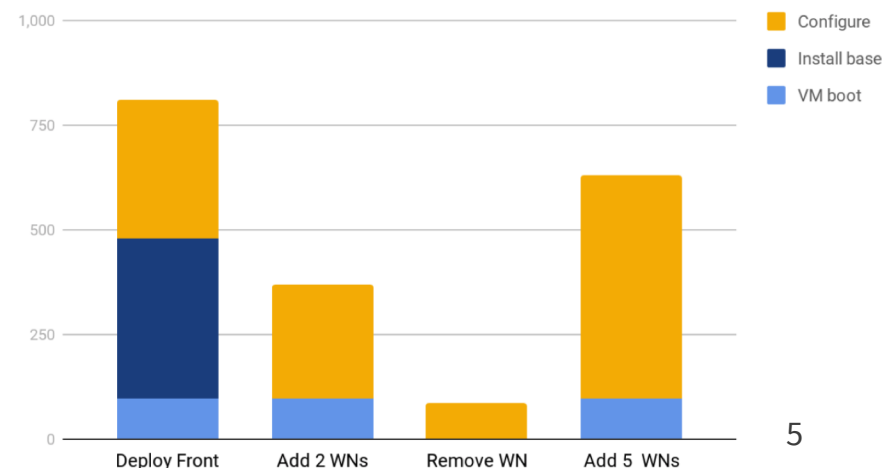
(1) <https://marketplace.eosc-portal.eu/services/elastic-cloud-compute-cluster-ec3>

(2) <https://marketplace.eosc-portal.eu/services/infrastructure-manager-im>

- The virtual infrastructure is managed by an elastic Kubernetes cluster spawn over the federated network
 - Containers and services are accessible from both sites but only through the federated network.
 - Resources are properly tagged (SGX and GPU capabilities and Brazil / Europe) so K8s applications are placed in the correct resource.
 - Infrastructure is described as code⁽³⁾.
- K8s Front-end is deployed and nodes are being powered on as the applications are deployed, creating the request for specific resources.

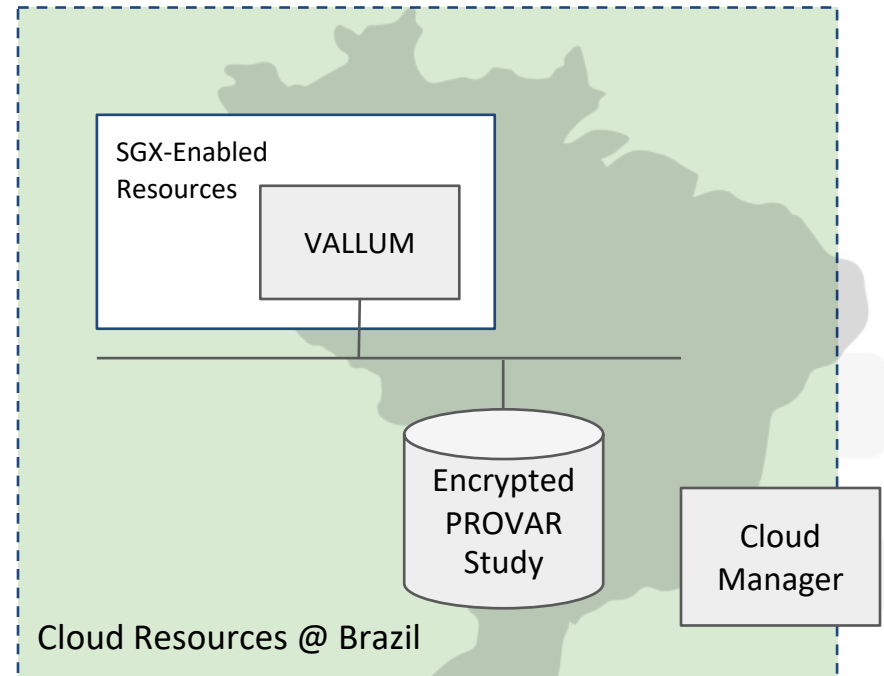


Deployment & Configuration Time (time in seconds)



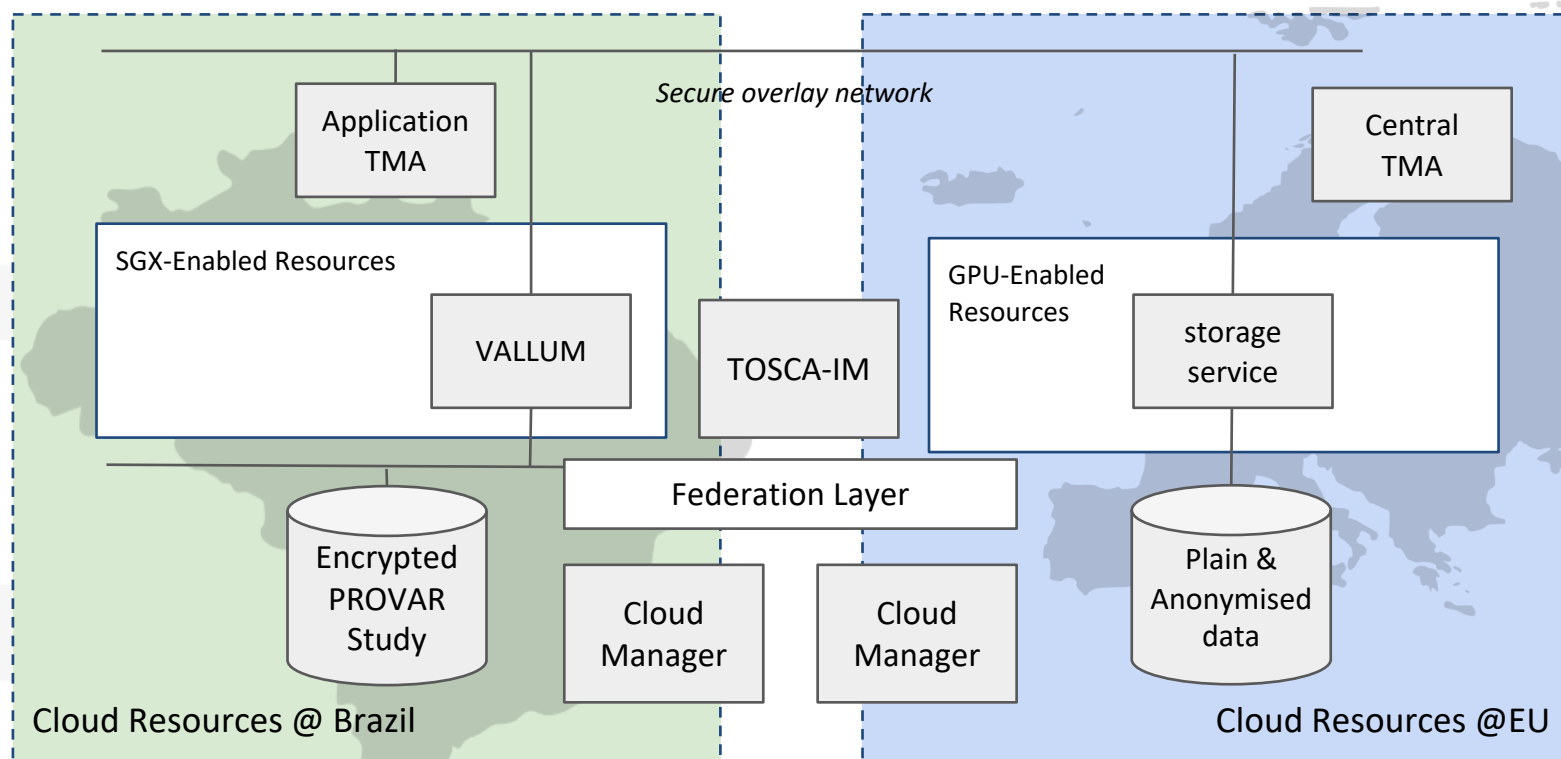
(3) <https://github.com/grycap/ec3/tree/atmosphere>

- A secure storage is deployed at the Brazilian side
 - It uses Vallum⁽⁴⁾, a service that provides on-the-fly anonymisation based on policies.
 - It masks (or blurs) the fields that are marked as sensitive to different profiles of users.
 - It relies on an HDFS filesystem for the files and on SQL databases for the structured data.
- It runs the data anonymisation and sensitive data access on enclaves running on SGX-enabled containers, so they can securely run even in untrusted cloud resources
 - Data remains encrypted in disk.

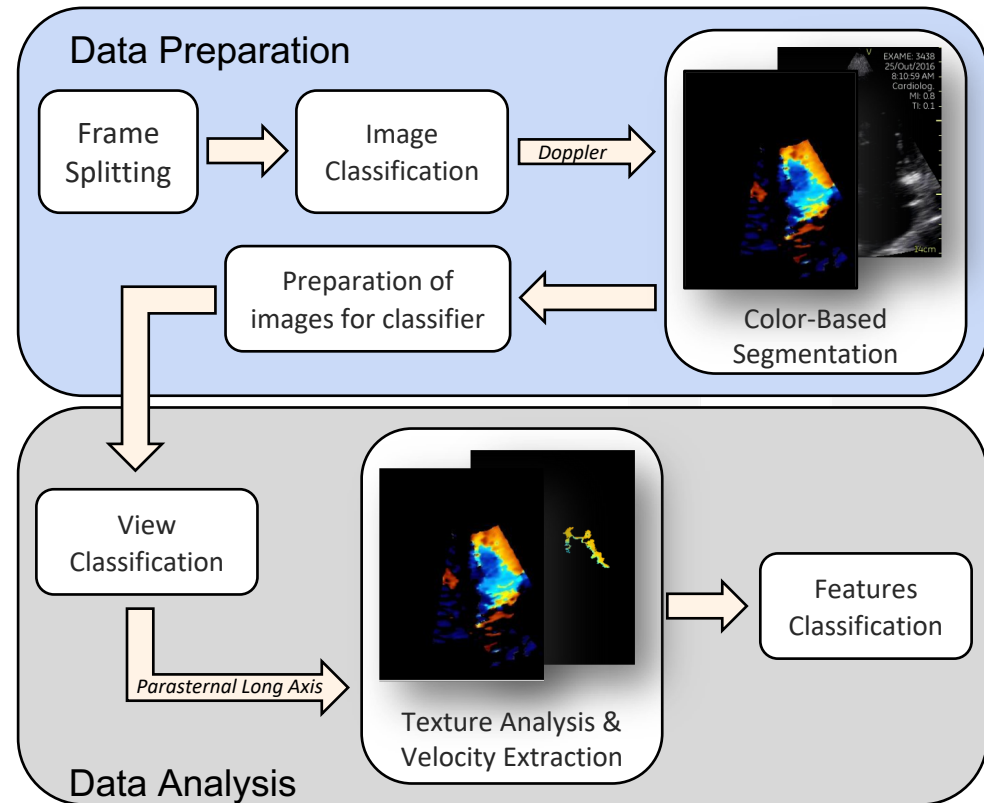


(4) <https://www.atmosphere-eubrazil.eu/vallum-framework-access-privacy-protection>

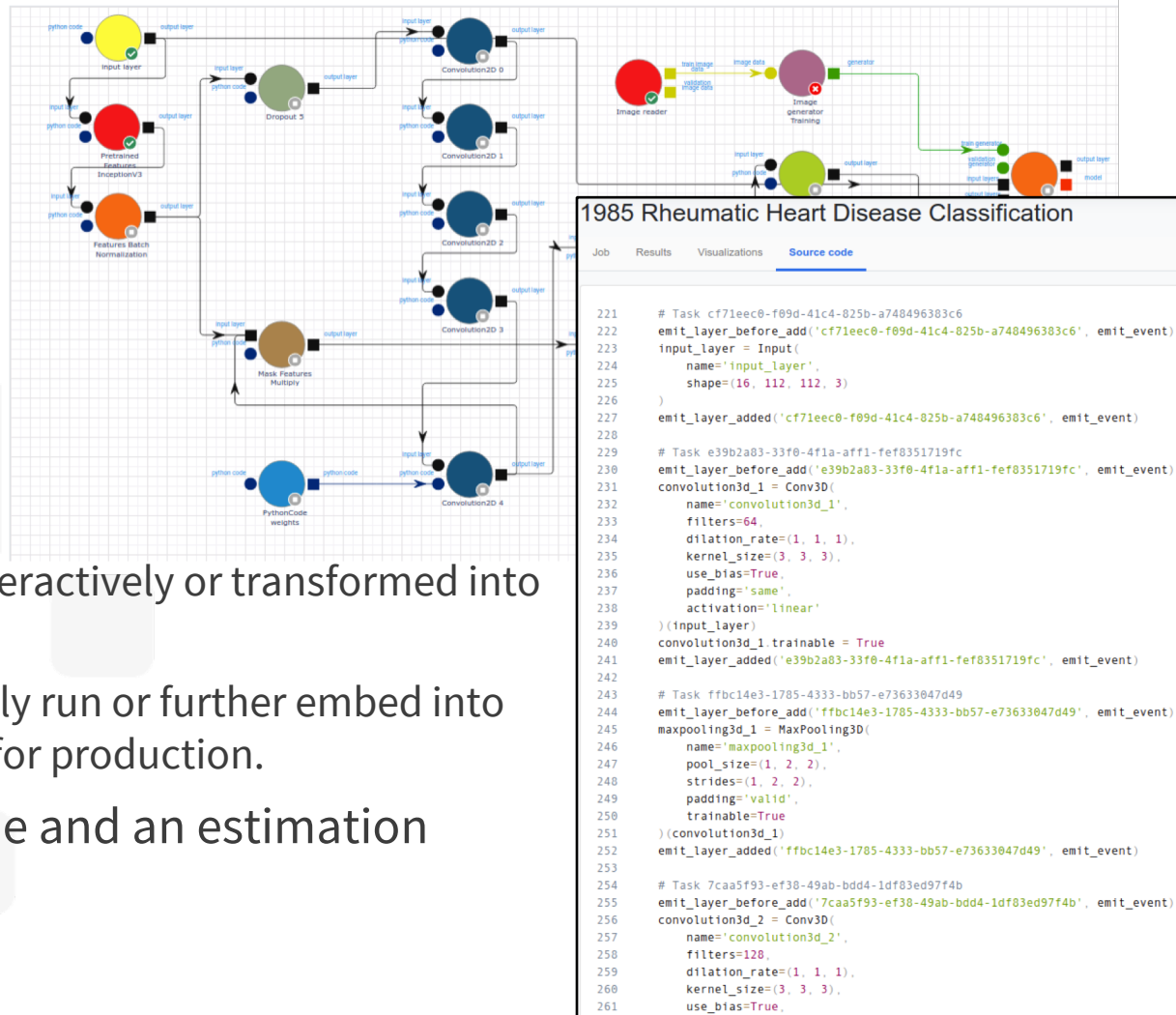
- Data is requested to Vallum from external users, but they will only access to partially anonymised data
 - Anonymised data (~1TB) is copied where the computing accelerators are placed.



- Videos are split into frames and classified by color inspection.
 - A color-based segmentation using k-means clustering extracts the color pixels from the Doppler images
- Images are classified according to their acquisition view using a CNN
 - Parasternal long axis view has proven to be relevant to obtain an accurate classification.
- First & second order texture analyses characterize the images by the spatial variation of pixel intensities.
 - Besides texture features, blood velocity information is also obtained.
- Finally, all the extracted features are classified through machine learning techniques in order to differentiate between RHD positive and healthy subjects.

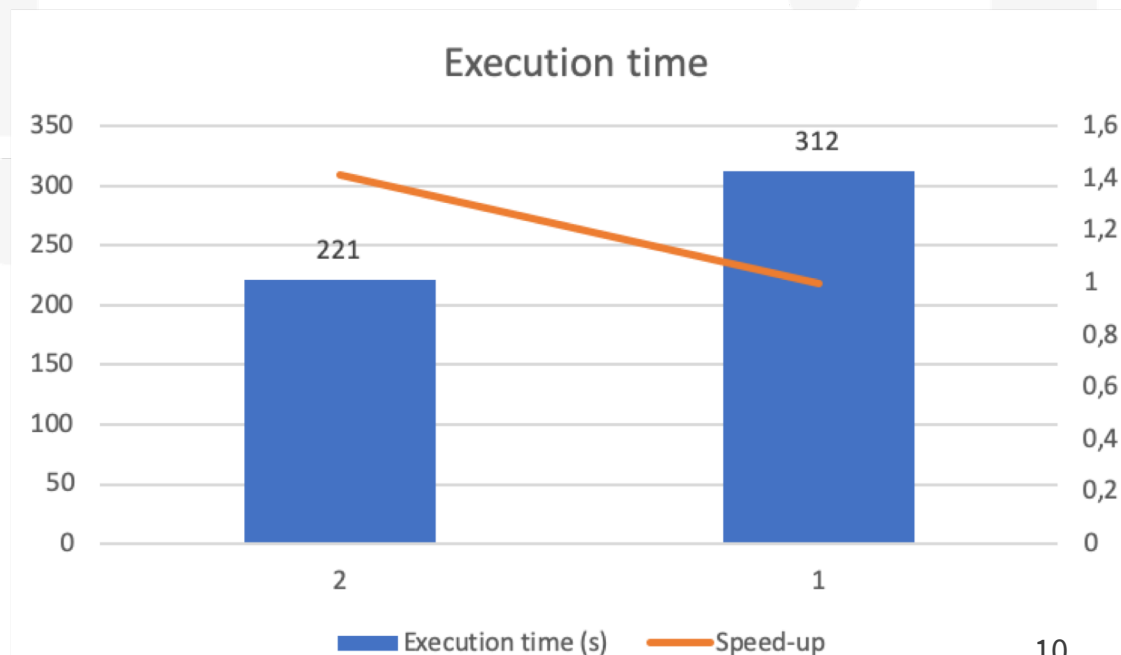


- The pipeline is developed using LEMONADE⁽⁵⁾
 - LEMONADE provides a GUI and a Machine Learning library to develop data analytics pipelines.
 - Pipelines can be run interactively or transformed into executable code.
 - Code can be interactively run or further embed into services to be exposed for production.
- A model building pipeline and an estimation pipeline are developed.

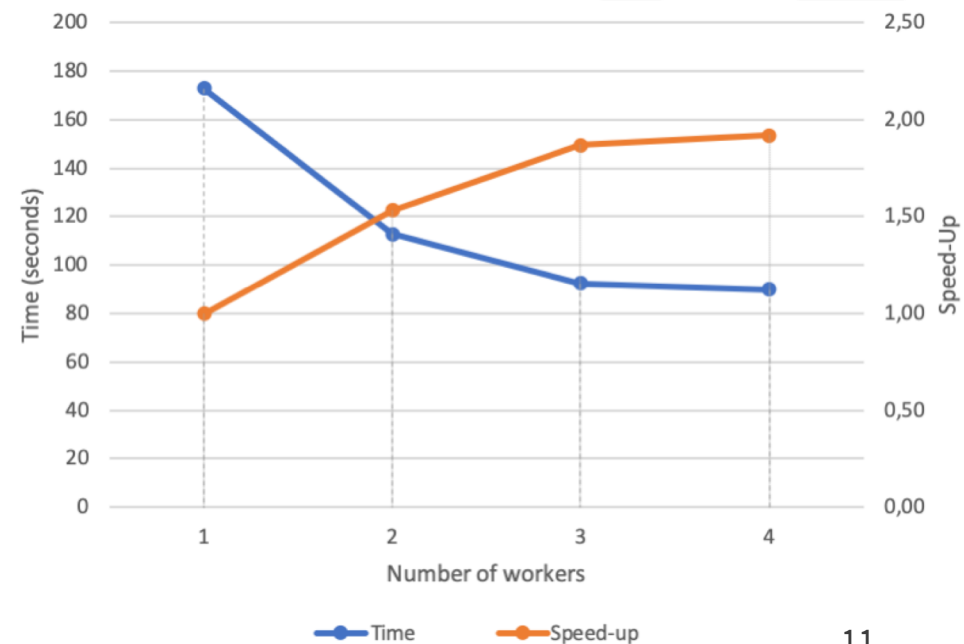


⁽⁵⁾ <https://www.atmosphere-eubrazil.eu/lemonade-live-exploration-and-mining-non-trivial-amount-data-everywhere>

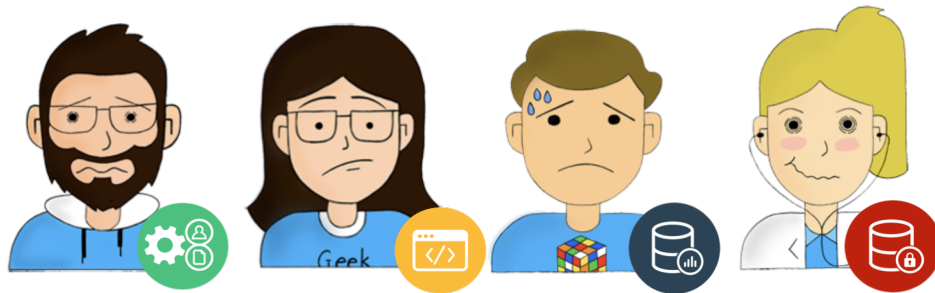
- Model building can run in parallel using MPI and Horovod
 - The model is build with keras using fp16 compression for the reduction operations.
 - Experiments have been used with 1 and 2 working nodes equipped with a TESLA V100 GPU connected through PCI Passthrough to the working nodes and the containers which run the processes.
 - Execution time shows a reduction with the addition of a second GPU but the speed-up is limited by the penalty of using an overlay network.



- An experiment has been performed for the classification of 8 patients on 1, 2 and 4 virtual compute nodes
 - Job code is extracted and executed through Jupyter on an ipython cluster that shares the filesystem.
 - Each node is a Kubernetes Pod executing a Docker container in a different Virtual Machine to reduce resource contention.
 - Speed-up is moderated (up to 2) but usability is high.

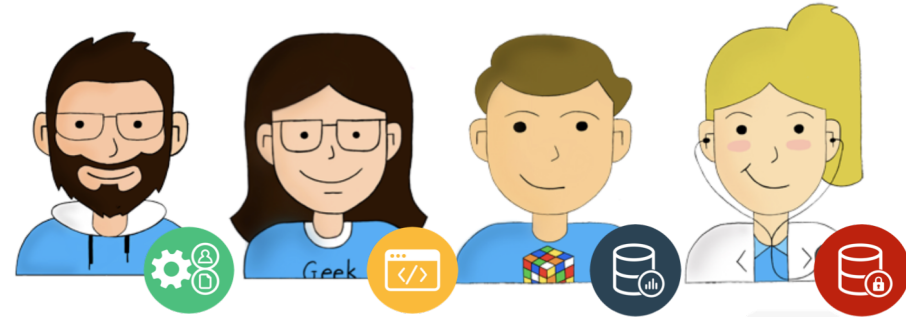


Before



- Need to manually configure the environment.
- Lack of reproducibility.
- Qualitative appraisal of the trustworthiness.
- Manual analysis of GDPR/LGDP risks
- Need to trust on the storage provider.
- Anonymisation level is qualitative.

After



- Applications templates for complex & distributed applications.
- Provide a repeatable way to deploy the whole application.
- Quantitative measure of trustworthiness
- Self-assessment of GDPR/LGDP.
- Trustable storage environment even on an untrusted provider.
- Quantitative anonymisation level.