# Optimizations to compute large correlation matrix onto GPU system of hybrid HPC clusters
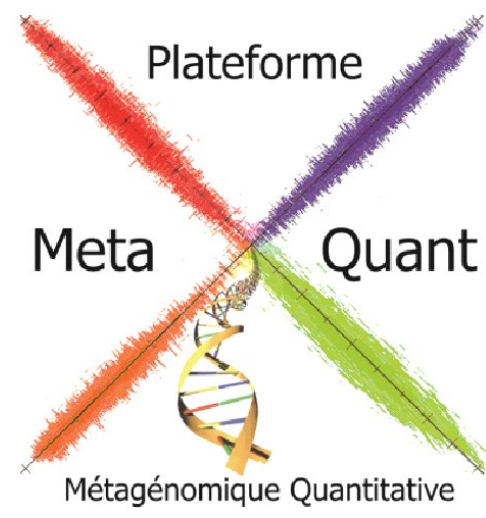
Dany Tello1, Fouad Boumezbeur2, Victor Arslan1, Vincent Ducrot1, Pierre Léonard2, Bouziane Moumen2, Sébastien Monot1, Nicolas Pons2, Tarik Saidani1, Pierre Renault2, Sean Kennedy2, Mathieu Almeida2, S.Dusko Ehrlich2 and J.M. Batto2

1**AS Plus, 22 rue René Coche 92170 Vanves, France**
2**Institut MICALIS, INRA CRJ, Domaine de Vilvert, 78352 Jouy-en-Josas, France**
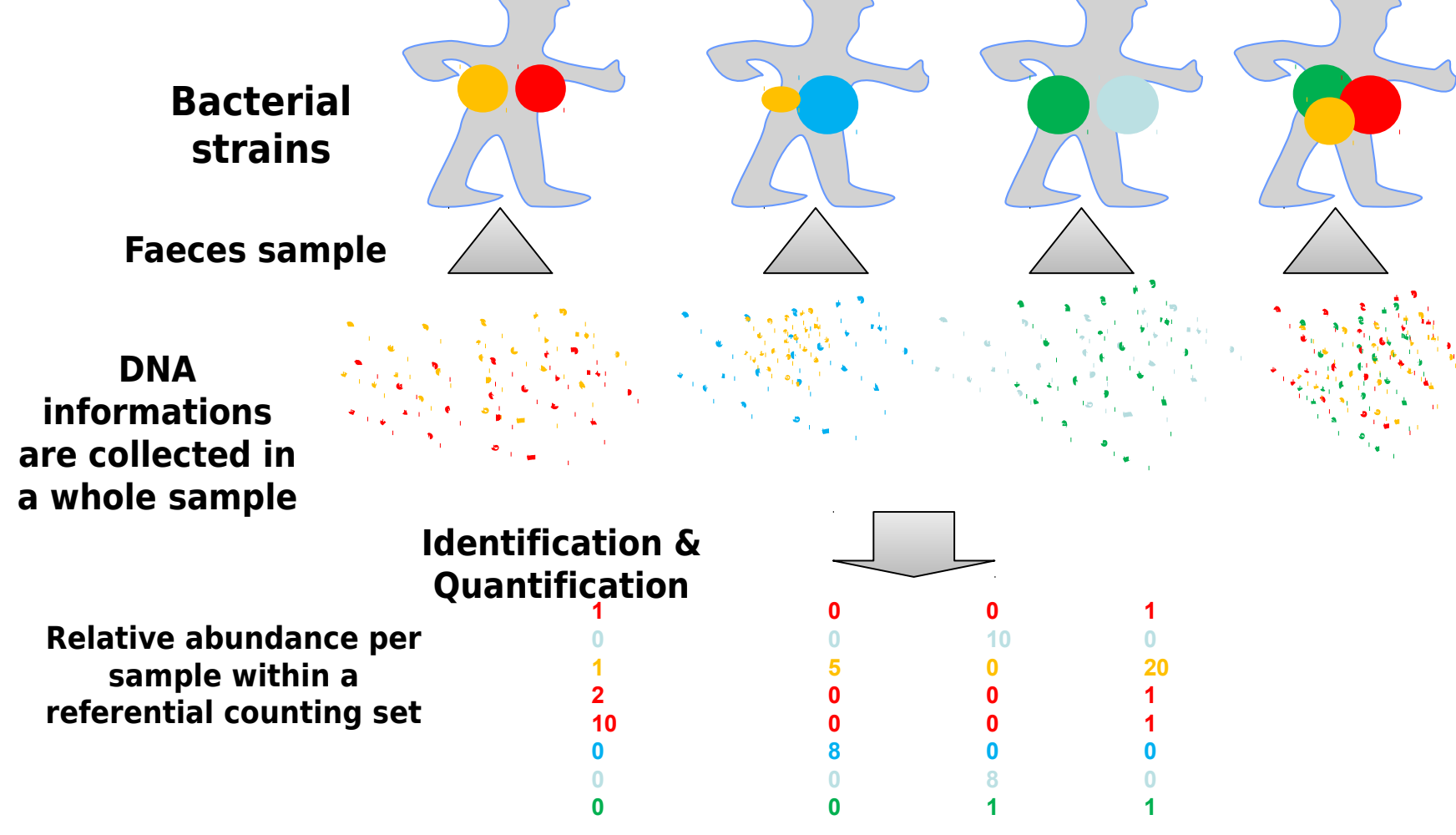
**Contact :** s.monot@asplus.fr, jean-michel.batto@jouy.inra.fr

In the MetaHIT project, large computations are performed to correlate quantitative information related to metagenomic samples. This program, named MetaProf, helps MetaQuant (INRA Micalis) in processing metagenomics analysis.

In order to achieve this computation within an acceptable time frame, we have developed a version of this program targeted for execution on GPU based clusters. Optimizations have been done to use in the most efficient way the capability of the GPU processor. The optimized program has been benchmarked on the hybrid part (GPU/nVidia M2090) of the GENCI supercomputer named Curie installed at TGCC (Très Grand Centre de Calcul du CEA - http://www-hpc.cea.fr/en/complexe/tgcc-curie.htm).

**How to explore, the Metagenomics Human Intestinal Tract (MetaHIT)?**



**Figure 1.** The computation problem consists of a classical auto-correlation on the input genes x sample matrix. Using Spearman or Pearson correlation, main issues deal with the amount of floating point operations required. Computation complexity is $O(N^2)$. Hence, for a 3M genes input matrix, it requires to compute 3,5 $10^{16}$ elementary operations.

**Figure 2.** A snapshot of Curie Supercomputer – hosted by GENCI in TGCC/CEA, Bruyères-Le-Châtel, France. The Curie Supercomputer has 144 nodes of dual Intel® Westmere® 2.66 GHz per node and dual nVidia® Tesla™ 2090 GPU (Fermi - 512 cuda cores, picture on the right side) per node – 192 Teraflops peak
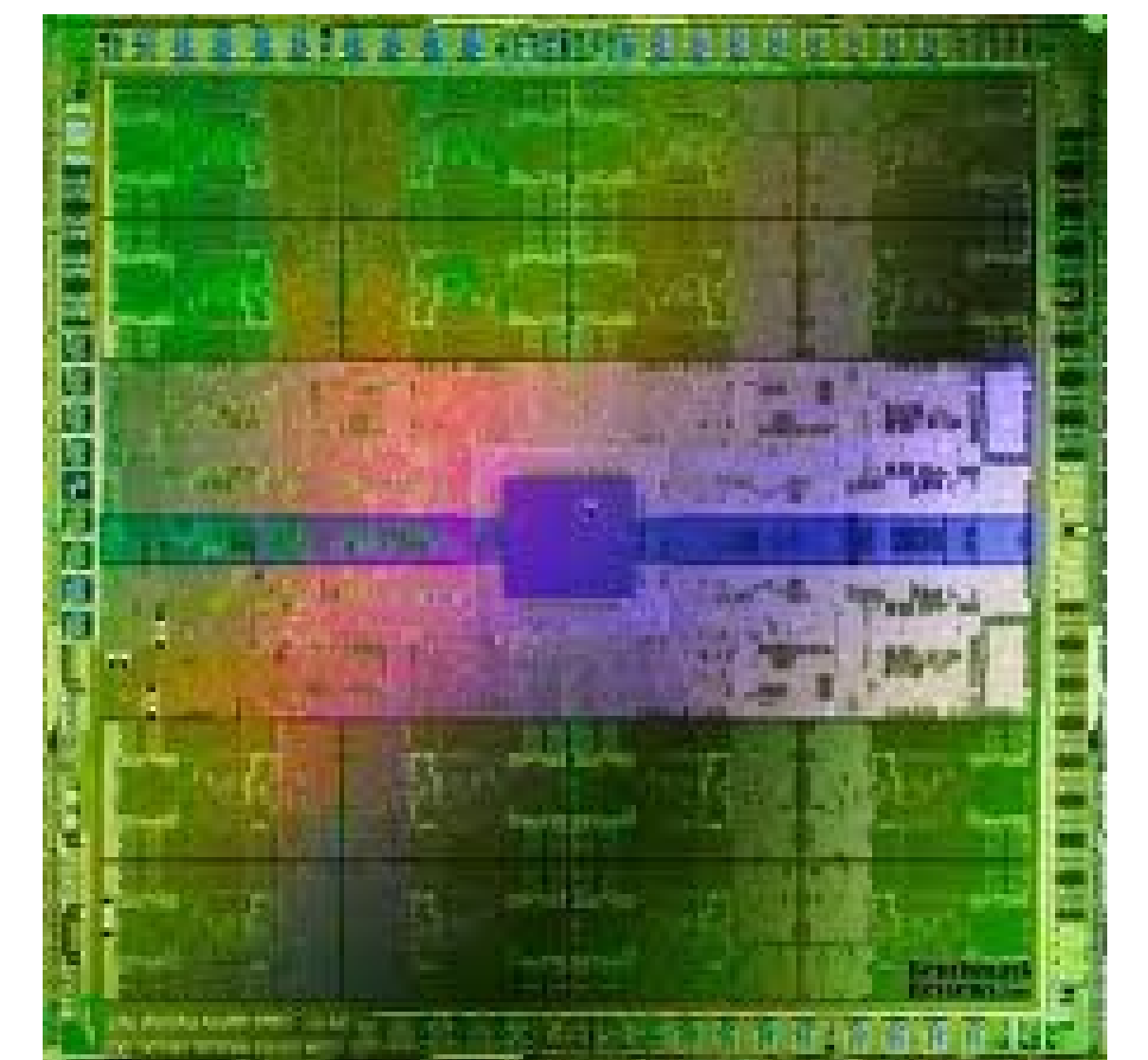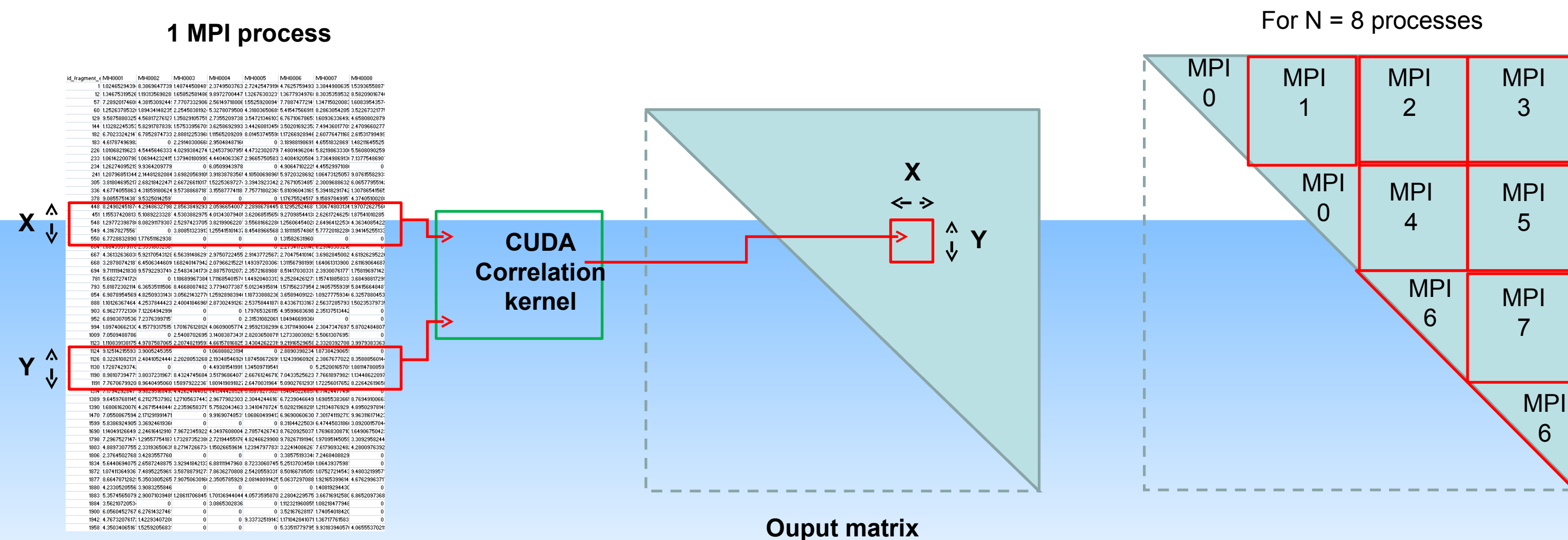
**Figure 3. Our implementation targets large clusters consisting of GPU based nodes hence a hybrid programming model : MPI at highest level aimed at distributed computing between nodes and Cuda at the node level to benefit from the massively parallel architectures of the GPU.**

A divide and conquer strategy applied to the output matrix was used to get a fair load balance between the cluster nodes and limit the memory requirement per process. (illustrated here with eight processes)
The MetaProf code is released on request with a non disclosure license.

## Quantitative Metagenomic and Big Computation

**Figure 4. A universal pipeline**
Metaprof is part of the **Meteor** processing pipeline designed by MetaQuant team to help in referencing data bank creation and in diagnostic evaluation. The MetaQuant Platform analyses 10 000 samples per year with a average size per project of 1000 samples. The MetaHIT project helps to size the pipeline and give a strong pressure to perform the computing analysis in a time below the week for a project.
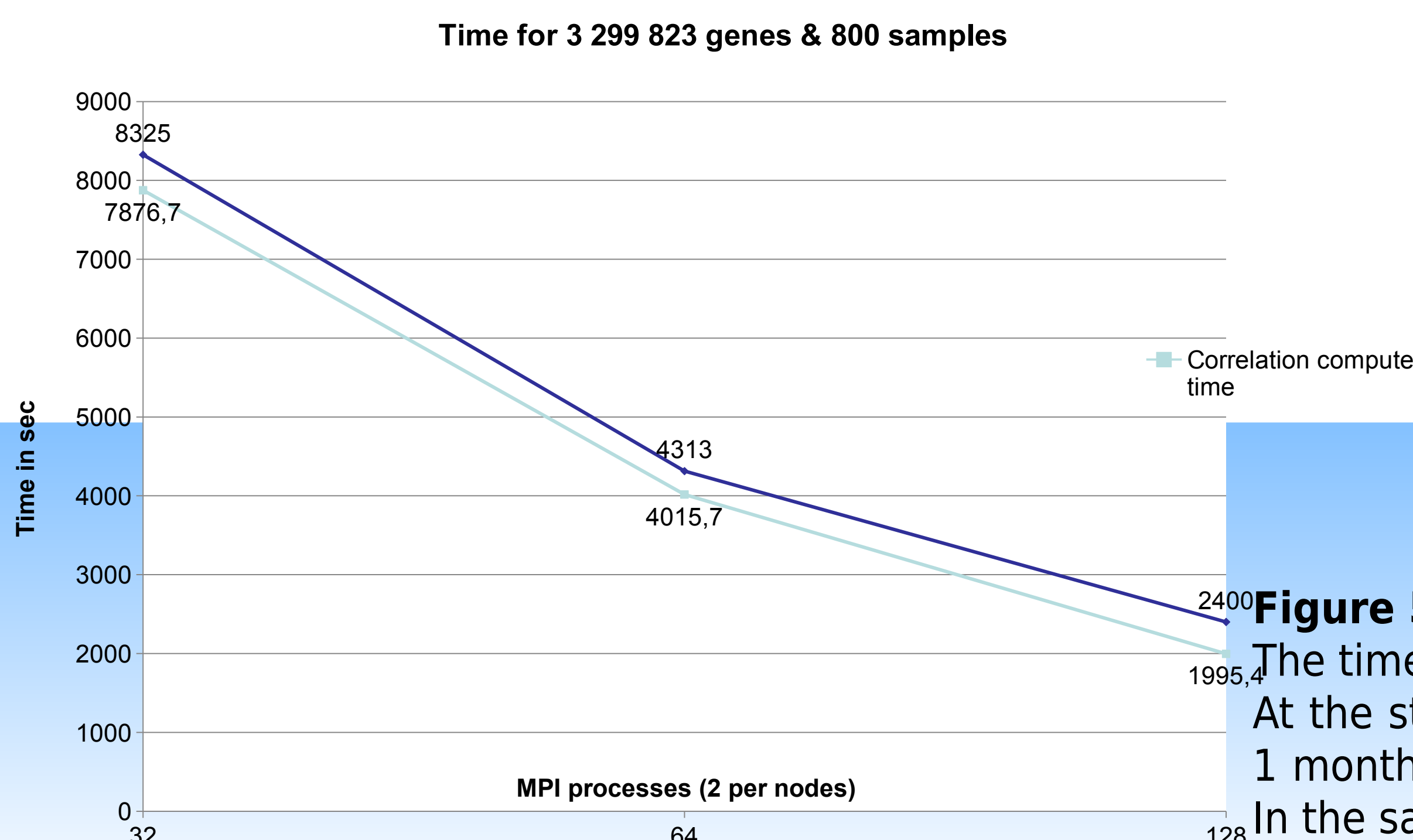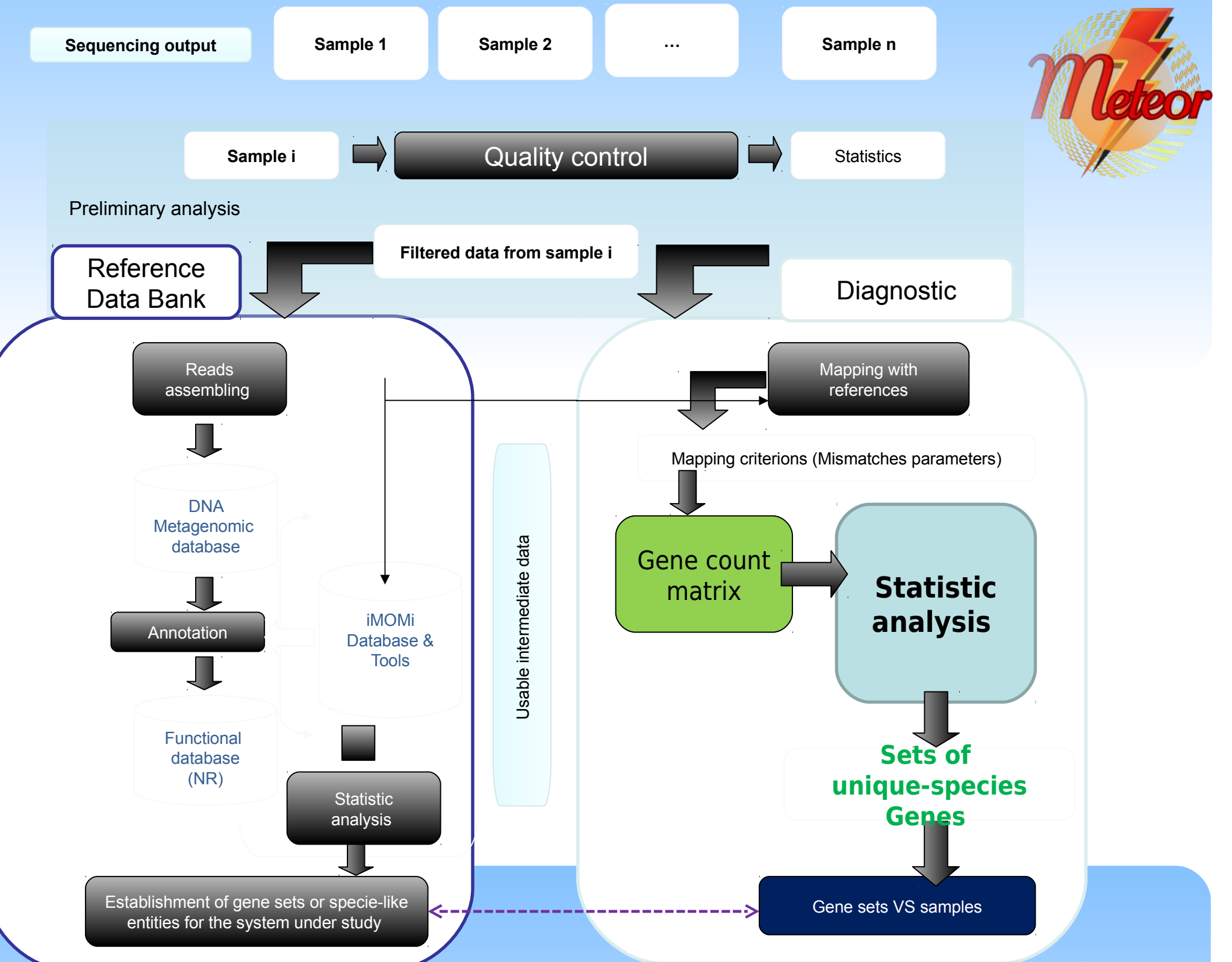
**Figure 5. Time measure for doing the computation**
The time measure benchmarked on a real use case – from the MetaHIT project.
At the starting point of the project, performing such a computation was expected to be done in 1 month. With the current optimization, the computation can be done quickly – less than 1 day. In the same time, the throughput of the NGS device is increasing faster than the Moore law.

## Expected results and potential impact

This optimization onto the Curie supercomputer proves that it is possible to use the GPU to perform huge computation and furthermore it brings GPU as a common component in the bioinformatics laboratory. It opens new opportunity : the possibility to perform more sophisticated computation on big data. Furthermore, as the benchmark proves that the time execution scales according to GPU number, it helps to design a small in house cluster. With a 16 GPU cluster, doing the MetaHIT use-case can be done in less than a day.

## Links :

MetaQuant : www.metaquant.eu     MetaHIT : www.metahit.eu     Micalis : www.micalis.fr     AS+ : www.asplus.fr     OpenGPU : www.opengpu.net/EN