

Performance Analysis and Waltime Prediction for Neuroscience Applications

Valentin Honoré, ensIIE & SAMOVAR, valentin.honore@ensieie.fr

I. CONTEXT

High Performance Computing (HPC) platforms are used for massively parallel applications requiring significant computational and memory resources. While traditional fields like astronomy and physics use monolithic codes, emerging domains such as neuroscience and bioinformatics are developing dynamic, heterogeneous workflows incorporating machine learning and AI techniques. These new applications raise particular challenges because their execution time is difficult to estimate and their memory usage is stochastic, complicating resource scheduling.

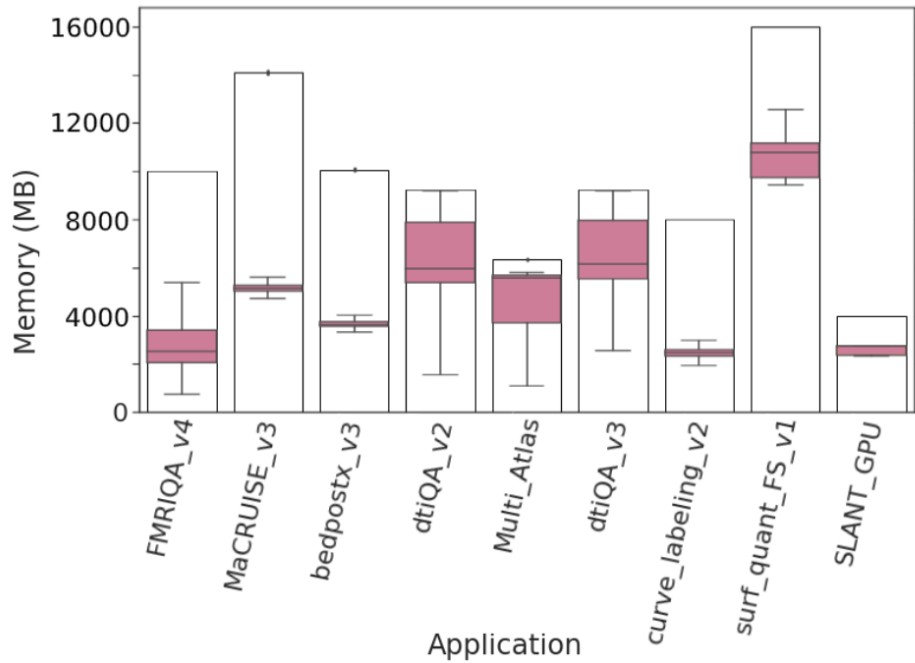


Fig. 1. Memory requests during submission and memory usage variations for nine representative medical and neuroscience applications. (extracted from [1])

New Machine Learning (ML) and AI frameworks have become important tools in exploratory domains. While progresses have been made over past years to improve these ML techniques, this progress has induced high requirements in terms of computations. For instance, Deep Learning techniques require an important training part where the quality of the model increases with the dataset size. Hence, such workflows involving ML techniques now targets HPC infrastructures that offer high computation support, as well as high memory and network performance. However, their profiles differ from classic HPC applications.

Often, the duration of these applications is difficult to estimate because they are input-independent. It is common for such an application to have walltimes between several hours to days. This characteristic is a real limitation for users for which requesting the maximum possible walltime often induces an overestimation that penalizes the total cost of the request. In addition, the stochastic memory utilization often requires users to request only high memory nodes for their execution.

Figure I presents the memory requirements and requests for nine exploratory applications from the medical and neuroscience department at the Vanderbilt University [2]. The logs are generated for a 6-month period in 2018 running on their in-house cluster. Users often utilize only fractions of the requested memory or end up with their application killed due to memory underestimation. Users tend to overestimate their resource requirements in both time and memory, which leads to these application typically waiting in the scheduler queue for days before eventually running.

II. GOAL OF THE PROJECT

In a previous work [1], we performed an extensive set of experiments to confirm the observations about the large walltime variations for a representative Neuroscience application, SLANT [3], [4]. SLANT is performing brain segmentation based on an input MRI image, as depicted in Figure II.

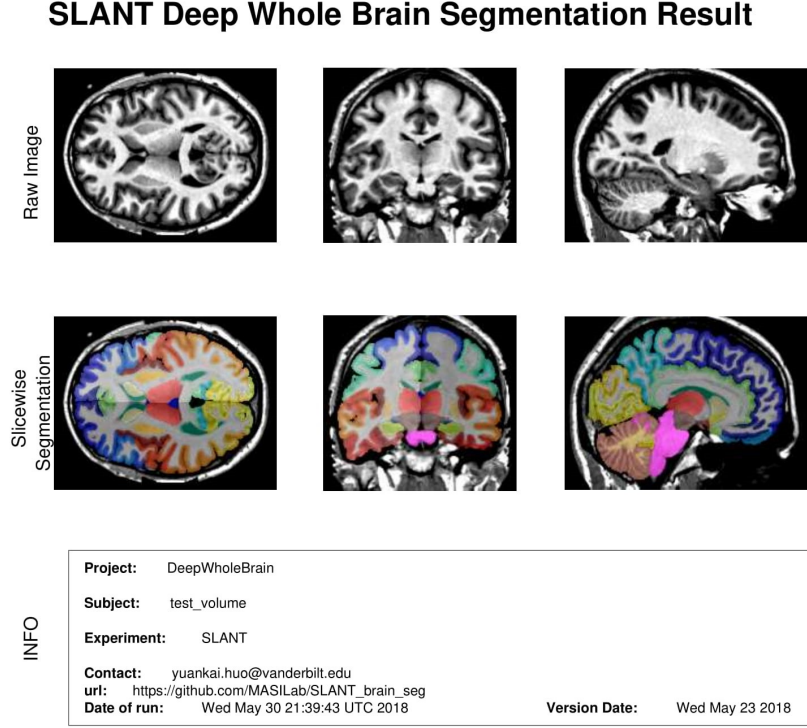


Fig. 2. Brain segmentation: a T1 MRI scan can be segmented to 133 labels based on BrainCOLOR protocol(<http://braincolor.mindboggle.info/protocols/>).

In this paper, we proposed a novel approach to extract a generic model of the runtime behavior of SLANT. We provided a first demonstration of what such an extraction would look like, along with scheduling techniques to use this model.

The goal of this project is to extend this work by using machine learning tools to predict the runtime based on image-specific characteristics. Before that, the variation in execution time will be experimented on a GPU version of SLANT, while in [1] we were restricted to a CPU version.

The main steps of this project will be:

- 1) Read the different papers related to this project [1], and also [5], [6] for context.
- 2) Starting from the reproducibility artifact [7] of our work [1], script and reproduce the experiments on a CPU machine of the Benagil team with different input sets.
- 3) Perform similar experiments using the GPU version of the application.
- 4) Compare the results between CPU and GPU versions of the application.
- 5) Use machine learning tools to predict the runtime based on image-specific characteristics.

Depending on the results and/or taste of the candidate, other directions could be envisioned such as performance analysis of the applications using tools developed in the team such as EZTrace.

III. TENTATIVE ROADMAP

For Master 2 students, the following roadmap could be envisioned:

Period	Planned work
M1	Related Work and experiment design
M2-M3	Scripting and running of experiments on both CPU and GPU
M3	Results analysis & Performance analysis of application
M4-M5	Design of ML predictive tools, Final report writing

A similar roadmap, over the academic year, should apply to Master 1 students.

IV. WORKING ENVIRONMENT

The candidate will work under the supervision of Valentin Honoré, a member of the Benagil research team. Depending on the progress and results on the project, a collaboration with Hongyang Sun from the University of Kansas could be envisioned.

The candidate will gain expertise on scripting experiments for applications running in containers, in a open-source environment. We will use Singularity to run the application(s) on the local cluster of the team. A certain appetite in bash scripting and reproducibility of experiments will be good asset for the experimental part of this project. An experience for the development of ML models is a plus but not necessary. A basic knowledge of Python should be sufficient to handle that part. It is expected that the candidate will gain experience on this subject during the project.

REFERENCES

- [1] A. Gainaru, B. Goglin, V. Honoré, and G. Pallez, “Profiles of upcoming HPC Applications and their Impact on Reservation Strategies,” *IEEE Transactions on Parallel and Distributed Systems*, 2020, Impact factor (last 2 years) : 2.6, *Article Influence Score* (February 2020) : 0.935, Core Rank A*. [Online]. Available: <https://hal.inria.fr/hal-03010676>
- [2] B. Landman, “Medical-image Analysis and Statistical Interpretation (MASI) Lab.” <https://my.vanderbilt.edu/masi/>, Accessed: 2025-09-17.
- [3] Y. Huo, Z. Xu, K. Aboud, P. Parvathaneni, S. Bao, C. Bermudez, S. M. Resnick, L. E. Cutting, and B. A. Landman, “Spatially localized atlas network tiles enables 3d whole brain segmentation from limited data,” in *Medical Image Computing and Computer Assisted Intervention – MICCAI 2018: 21st International Conference, Granada, Spain, September 16-20, 2018, Proceedings, Part III*. Berlin, Heidelberg: Springer-Verlag, 2018, p. 698–705. [Online]. Available: https://doi.org/10.1007/978-3-030-00931-1_80
- [4] —, “Spatially localized atlas network tiles enables 3d whole brain segmentation from limited data,” 2018. [Online]. Available: <https://arxiv.org/abs/1806.00546>
- [5] G. Aupy, A. Gainaru, V. Honoré, P. Raghavan, Y. Robert, and H. Sun, “Reservation Strategies for Stochastic Jobs,” in *IPDPS 2019 - 33rd IEEE International Parallel and Distributed Processing Symposium*. Rio de Janeiro, Brésil: IEEE, May 2019, pp. 166–175, average Acceptance Rate (last 5 years) : 24.5%, Core Rank A.
- [6] A. Gainaru, B. Goglin, V. Honoré, G. Pallez, P. Raghavan, Y. Robert, and H. Sun, “Reservation and Checkpointing Strategies for Stochastic Jobs,” in *IPDPS 2020 - 34th IEEE International Parallel and Distributed Processing Symposium*, La Nouvelle Orléans, USA, May 2020, average Acceptance Rate (last 5 years) : 24.5%, Core Rank A.
- [7] V. Honoré, “Reproducibility artifact,” https://gitlab.inria.fr/vhonore/stochastic_app_profiling, Accessed: 2025-09-17.