



---

Arnaud Legrand, CNRS researcher  
INRIA POLARIS project, LIG Laboratory

## Internship proposal

### Modeling and Performance Prediction of Linear Algebra Functions

**Advisors:** Arnaud Legrand (CNRS/LIG, Grenoble), Tom Cornebize (UGA/LIG, Grenoble)

**Required Skills:**

- Python and C programming, shell, ssh, git
- Basis of experiment analysis (with Python or R) is a plus
- Basis of statistics is a plus (linear regression)

## 1 Context

There is a continued need for higher compute performance: scientific grand challenges, engineering, geophysics, bioinformatics, etc. Indeed, the technological advances driven by the home PC market have contributed to achieving high performance in commodity components. For decades, computer performance had doubled every 18 months merely by increasing the clock frequency of the processors. This trend stopped last decade for reason of electricity consumption and heat. Indeed, the computational power of a computer increases nearly sub-linearly with clock frequency while the energy consumption increases more than quadratically.

As an answer to the power and heat challenges, processor constructors have increased the amount of computing units (or cores) per processor. Modern High Performance Computing (HPC) systems comprise thousands of nodes, each of them holding several multi-core processors. Most processors also have vector extensions, multiple levels of cache, dynamic voltage and frequency scaling (DVFS) and hyperthreading.

This increased hardware complexity gives a gain of performance, but also a higher variability: predicting the duration of a function call is no longer an easy task. This makes performance evaluation even more difficult and raises many software and *devops* challenges. Why do we get such performance? Is it normal or not?

## 2 Environment

The members of the POLARIS team focus their research on performance evaluation and optimization of large scale systems and parallel applications. They have a strong expertise regarding parallel applications and environment for parallel programming, performance evaluation of large scale distributed systems.

Some of POLARIS members are involved in the Joint Laboratory for Petascale Computing<sup>1</sup> between University of Illinois at Urbana-Champaign<sup>2</sup> Inria<sup>3</sup>, Argonne National Laboratory<sup>4</sup>, Illinois' Center for Extreme-Scale Computation<sup>5</sup>, the National Center for Supercomputing Applications<sup>6</sup>, the Barcelona Supercomputer Center<sup>7</sup>, Julich and the Riken. Some of their members are also involved in the European Mont-Blanc<sup>8</sup> (European scalable and power efficient HPC platform based on low-power embedded technology) that aims at prototyping HPC supercomputers from today's energy-efficient solutions used in embedded and mobile devices (e.g., ARM processors). Some of the POLARIS members are also in contact with research & development groups from Bull/ATOS that strongly influence the design and capacity planning of supercomputers.

<sup>1</sup><http://jointlab.ncsa.illinois.edu/>

<sup>2</sup><http://illinois.edu/>

<sup>3</sup><http://www.inria.fr/en/>

<sup>4</sup><http://www.anl.gov/>

<sup>5</sup><http://iacat.illinois.edu/>

<sup>6</sup><http://www.ncsa.illinois.edu/>

<sup>7</sup><http://www.bsc.es/>

<sup>8</sup><http://www.montblanc-project.eu/>

### 3 Goal

A lot of the HPC applications rely on dense linear algebra operations, using a standard interface called *Basic Linear Algebra Subprograms* (BLAS<sup>9</sup>). Several state-of-the-art implementations, such as OpenBLAS<sup>10</sup>, Atlas<sup>11</sup> and MKL<sup>12</sup>, can efficiently take advantage of the features of recent processors like multithreading, vectorization, cache and accelerators.

Modelling accurately the performance of the BLAS operations is an important task. A possible end-goal is to implement **performance non-regression tests** (most continuous integration tools only address functional non-regression tests). With a good model, we could regularly measure the duration of the operations and detect any significant but subtle changes or anomalies (possibly introduced by a modification in the software stack).

This internship can be summarized in a few steps:

1. Collect and parse the specification files (\*.h) of the BLAS library.
2. Find the list of parameter names and types for each function. For instance, the function that computes the product of two matrices has three `int` parameters (the sizes) and three `double*` parameters (the matrices).
3. Choose a possible model for the function. There is a trade-off between the complexity of the model and its accuracy. For instance, in a first approximation, we can consider that the product of two  $M \times N$  and  $N \times K$  matrices takes a time  $t = \alpha MNK + \beta + \varepsilon$ , where  $\alpha$ ,  $\beta$  are two parameters that should be determined empirically and  $\varepsilon$  is a random variable accounting for measurement variability. This approximation may not be satisfying for all cases, especially for small matrices where lower terms of the polynomial (e.g.,  $MN$ ,  $NK$ , or  $N$ ) may be non-negligible (e.g., because of cache effects). But including all terms of the polynomial in the model may also be detrimental as over-fitting induces uncertainty that precludes effective extrapolation.
4. Design and run an experiment to evaluate the model. For instance, which matrix sizes should we test? In which order?
5. Analyze the results of the experiment. This will include data visualization, linear regressions and bayesian modeling (e.g., using STAN<sup>13</sup> and/or PyMC3<sup>14</sup>).
6. Go back to step 3 if necessary.
7. Implement performance non-regression tests using the right model.

Prototypes of steps 3 to 5 have already been developed. The goal is to automate all these steps as much as possible. We also need to track the provenance of the results (e.g., when and how was the experiment done, what were the settings of the machine, etc.) and control the variability (e.g., disable the hyperthreading and turboboost features, pin the threads to the CPU cores, etc.), hence this internship will follow the reproducible research<sup>15</sup> good practices. The student will use a notebook (Jupyter<sup>16</sup> and/or Org mode<sup>17</sup>) and all the experiments will be done using Grid5000<sup>18</sup> platform.

This internship will allow you to sharpen your experimental skills both on a practical and theoretical side. We expect the intern to be curious, rigorous, and autonomous but both advisors will happily provide all the help needed on the software, hardware, theory, experimental methodology.

---

<sup>9</sup><http://www.netlib.org/blas/>

<sup>10</sup><https://github.com/xianyi/OpenBLAS>

<sup>11</sup><http://math-atlas.sourceforge.net/>

<sup>12</sup><https://software.intel.com/en-us/mkl>

<sup>13</sup><http://mc-stan.org/>

<sup>14</sup><https://docs.pymc.io/>

<sup>15</sup><https://en.wikipedia.org/wiki/Reproducibility>

<sup>16</sup><http://jupyter.org/>

<sup>17</sup><https://orgmode.org/>

<sup>18</sup><https://www.grid5000.fr>