

# IT4Innovations Directors Discretion application

**Name of the project:** Bringing Statistics to High Performance Computing

**Number of core hours requested:** 100 000

**Multi-year access:** 9

**Date:** November 18, 2021

**Name and surname of primary investigator:** George Ostrouchov

**Affiliation of primary investigator:** Oak Ridge National Laboratory and Charles University (via Fulbright)

**e-mail:** ostrouchovg@ornl.gov

**Names and surnames of other investigators<sup>1</sup>:** Jaromir Antoch

**Affiliations of other investigator<sup>1</sup>:** Charles University

**e-mail<sup>1</sup>:** antoch@karlin.mff.cuni.cz

**Research area:** Computational Statistics

## Popular abstract:

Explainable data science, which is the hallmark of statistical methods, is necessary for research that aims to discover knowledge from today's abundance of data. So far, few statisticians have ventured into supercomputing to develop scalable statistical software and to build a statistics community within supercomputing. This allocation is to support a course and a seminar at the Department of Probability and Mathematical Statistics, Faculty of Mathematics and Physics, Charles University, to build the next generation of statistics-centered and computation-savvy data scientists.

---

<sup>1</sup> Comma separated list

## Scientific readiness:

### Aims and objectives

The long-term aim is to build a statistics community within supercomputing. Currently, statisticians are largely absent among scientists who are working with large computing platforms. At the same time, statistician developed software is in high use among scientists in the biological, health, and social sciences as well as in economics and finance. In these and other disciplines, statistical methods are attractive because they emphasize explainability and understanding of uncertainty. Consequently, a statistics community within supercomputing can bring more explainability and understanding of uncertainty to scientists across many disciplines working on such large platforms. By far, the favorite software development platform of statisticians is the R language and its multitude of analysis and data exploration packages.

Specifically, this Directors Discretion allocation request results from a Fulbright Distinguished Chair appointment of the PI in the Department of Probability and Mathematical Statistics, Faculty of Mathematics and Physics, Charles University, for the Spring semester of 2022. The appointment is based on a proposed course and a seminar series in the department to teach high performance statistical computing with practical components on IT4I systems.

The course “Adventures in Supercomputing with R” is a masters level course that teaches a range of parallel computing concepts, beginning with simple pleasingly (embarrassingly) parallel computations and ends with utilization of highly scalable matrix computation libraries. This is accomplished while the students frequently navigate remote access and code development platforms needed for working with IT4I systems.

The seminar series “High Performance Statistical Computing with R” will cover most of the same concepts as the course, although it will be more tailored to each student’s interests and direction of their doctoral thesis project.

### Methods and state-of-the-art

The courses will utilize several parallel computing packages that already exist in R. This is mainly the core R parallel package for its shared memory capabilities and several distributed computing and GPU capable packages from the pbdR.org project, which the PI leads at the Oak Ridge National Laboratory. Most packages in the pbdR.org project rely on MPI in its distributed computing capable components, sometimes via libraries such as ScaLAPACK.

This parallel infrastructure will be exercised across a number of standard statistical computing problems and large data sets. The single program multiple data (SPMD) approach to distributed computing via MPI+X is the core concept that will be taught in most distributed cases. The X represents either OpenMP (for example, with OpenBLAS) or the unix fork (when the core R parallel package is employed). There are pbdR components that can utilize CUDA programmable GPUs, in which case the X is a GPU-capable code, although this is an advanced topic that we may or may not reach. In some cases, MPI can be effectively pushed into the node cores so it is MPI only, without the X.

### Impact and outlooks

It is hoped that a few of the seminar series doctoral students will gain enough knowledge to request their own allocations along with their advisors for their dissertation projects.

Knowledge of the workflow to navigate code development and job submission, or even knowing that others in your discipline have done it, lowers the entry barrier for the statistics community. It is not enough to build something: some have to be trained to lead the way in using it.

## Computational readiness:

### Computational approach, parallelization, and scalability:

Some of the computational approach is already described in the **Methods and state-of-the-art** section above. The distributed components of pbdR packages have been shown scalable on systems like Titan and Summit at OLCF, as can be seen in this [HPC Wire article](#) and this [publication](#). The PI has also taught a class on “[Using R for HPC Data Science](#)” at IT4I in 2016.

#### **Computational resources:**

Class and seminar practical assignments on IT4I systems will have a broad range. Early assignments will be just single node use with short duration under a minute. Later assignments will involve data with roughly a hundred GB on 2 to 8 nodes, requiring 2 to 10 minutes of compute time (for example, involving a singular value decomposition). We favor using larger compute resources for shorter turnaround in iterations rather than the opposite. We put more pressure on algorithm scalability and run well below memory limitations. We also note a tension that exists for data analysis on batch systems. Data analysis is a discovery process, which requires many iterations and favors interactive computing over batch. As we will use batch iterations, a statistics community used to interactive computing will initially be less efficient with batch and use more iterations.

There will be about 12 assignments. With an assumption of 30 students total across the course and the seminar, an average assignment run on 2 nodes for 5 minutes, and a need to make 40 runs to debug, we have a total of  $12 \cdot 30 \cdot 2 \cdot (5/60) \cdot 40 = 2\,400$  node hours = 86 400 Barbora core hours. We will utilize sharp compute time limits on jobs to avoid unnecessary overuse while debugging.

One reason we make the allocation computation with node hours is to experience distributed concepts. Karolina nodes have 128 cores (compared to 36 on Barbora), so a roughly equivalent request for Karolina would yield only 675 node hours, resulting in fewer multinode experiences. We round up our request to 100 000 core hours to have room for multinode experiences on Karolina. We assume that a core allocation is transferable between Barbora and Carolina.

Congratulations to Karolina on #8 Green500! While we do have some capability in pbdR built on top of the cuBLAS library, it is an advanced topic that we may not reach in either the class or seminar. If we do get some interest from our students, we will make an additional request.

### Economic readiness:

#### Socioeconomic impact:

The course and seminar series contribute to the broader aim to build a statistics community within supercomputing by training students, some of whom can become leaders in this effort. But this alone is not enough, and other avenues are being used to bring statisticians to supercomputing. For example, the PI organized and moderated a Supercomputing21 Panel Session on [High-Performance Statistical Computing: Building a Statistics Community within Supercomputing](#).

As software developed by statisticians is in high use among scientists who value the concepts of explainability and understanding of uncertainty, a statistics community within supercomputing will bring these concepts into higher use among scientists across disciplines working on such large platforms.

#### References:

Drew Schmidt, Wei-Chen Chen, Michael A. Matheson, George Ostrouchov, Programming with BIG Data in R: Scaling Analytics from One to Thousands of Nodes, Big Data Research, Volume 8, 2017, Pages 1-11, ISSN 2214-5796, <https://doi.org/10.1016/j.bdr.2016.10.002>.