

SDSC Summer Institute 2017

July 31 – August 4, 2017

SDSC Auditorium at UC San Diego

Lesson material repository: <https://github.com/sdsc/sdsc-summer-institute-2017>

MONDAY, July 31	
8:00 – 8:30AM	Registration, Coffee
8:30 – 8:45	Welcome Shawn Strande, Deputy Director, SDSC
8:45 – 9:30	Orientation Andrea Zonca, Senior Computational Scientist, SDSC - Director of the Summer Institute
9:30 – 12:15 break 10:15-10:45	How to login to Comet and launch jobs Mahidhar Tatineni, Director, User Services, SDSC
12:15 – 1:30	Lunch at Café Ventanas
1:30 – 3:00	How to use Comet filesystems: home, parallel Lustre, local SSD Manu Shantharam, Senior Computational Scientist, SDSC
3:00 – 3:30	Break
3:30 – 5:00	How to understand the performance of your software Bob Sinkovits, Director for Scientific Computing Applications, SDSC
5:30 – 8:30PM	Reception at Wayne Pfeiffer's home overlooking the Pacific <i>Sweater or jacket recommended</i> Shuttle provided from SDSC driveway

TUESDAY, August 1	
8:00 – 8:30AM	Coffee
8:30 – 10:00	How to automate a data analysis pipeline with a workflow manager Ilkay Altintas, SDSC's Chief Data Science Officer, Director, Workflows for Data Science (WorDS) Center of Excellence, SDSC
10:00 – 10:15	Break
10:15 – 12:15	Parallel sessions: basic software version control with git / advanced Github workflows and conflict management with git Andrea Zonca, Senior Computational Scientist, SDSC
12:15 – 1:30	Lunch at Café Ventanas
1:30 – 2:30	How to use Science Gateways (and how to build them) Amit Majumdar, Division Director, Data Enabled Scientific Computing, SDSC
2:30 – 3:30	How to use Singularity containers and Comet Virtual Clusters Trevor Cooper, High Performance Computing Systems Manager, SDSC
3:30 – 3:45	Break
3:45 – 4:15	SDSC Data Center Tour
4:15 – 5:00PM	Hands-on practice continues with mentors available for questions

WEDNESDAY, August 2
PARALLEL SESSIONS

8:00 – 8:30	Coffee	
	Track 1 <i>Auditorium</i>	Track 2 <i>Synthesis Center E-B143</i>
Session 1 8:30 – 12:00	GPU Computing and Programming Andreas Goetz, Research Scientist and Principal Investigator, SDSC This session provides an introduction to massively parallel computing with graphics processing units (GPUs). The use of GPUs is becoming increasingly popular across all scientific domains since GPUs can significantly accelerate time to solution for many computational tasks. Participants will be introduced to essential background of the GPU chip architecture and will learn how to program GPUs via the use of libraries, OpenACC compiler directives, and CUDA programming. The session will incorporate hands-on exercises for participants to acquire the skills to use and develop GPU aware applications.	Spark for Scientific Computing Andrea Zonca, Senior Computational Scientist, SDSC Mahidhar Tatineni, Director, User Services, SDSC Apache Spark is a cluster computing framework extensively used in Industry to process large amount of data (up to 1PB) distributed across thousands of nodes. It has been designed as a successor of Hadoop focusing on performance and usability. It provides interface in Python, Scala and Java. This session will provide an overview of the capabilities of Spark and how they can be leveraged to solve problems in Scientific Computing. Next it will feature a hands-on introduction to Spark, from batch and interactive usage on Comet to running a sample map/reduce example in Python. The final part will be devoted to two key libraries in the Spark ecosystem: Spark SQL, a general purpose query engine that can interface to SQL databases or JSON files and Spark MLlib, a scalable Machine Learning library.
12:00 – 1:30	Lunch at Café Ventanas	
Session 2 1:30 – 5:00PM	Performance Optimization Bob Sinkovits, Director for Scientific Computing Applications, SDSC This session is targeted at attendees who both do their own code development and need their calculations to finish as quickly as possible. We'll cover the effective use of cache, loop-level optimizations, force reductions, optimizing compilers and their limitations, short-circuiting, time-space tradeoffs and more. Exercises will be done mostly in C, but emphasis will be on general techniques that can be applied in any language.	Scientific visualization with VisIt and data sharing Amit Chourasia, Senior Visualization Scientist, SDSC Visualization is largely understood and used as an excellent communication tool by researchers. This narrow view often keeps scientists from fully using and developing their visualization skillset. This tutorial will provide a "from the ground up" understanding of visualization and its utility in error diagnostic and exploration of data for scientific insight. When used effectively visualization can provide a complementary and effective toolset for data analysis, which is one of the most challenging problems in computational domains. In this tutorial we plan to bridge these gaps by providing end users with fundamental visualization concepts, execution tools, customization and usage examples. Finally, a short introduction to SeedMe.org will be provided where users will learn how to share their visualization results ubiquitously.

THURSDAY, August 3
PARALLEL SESSIONS

8:00 – 8:30	Coffee	
	Track 1 <i>Auditorium</i>	Track 2 <i>Synthesis Center E-B143</i>
Session 3 8:30 – 12:00	Parallel Computing using MPI & Open MP Pietro Cicotti, Senior Computational Scientist, SDSC This session is targeted at attendees who are looking for a hands-on introduction to parallel computing using MPI and Open MP programming. The session will start with an introduction and basic information for getting started with MPI. An overview of the common MPI routines that are useful for beginner MPI programmers, including MPI environment set up, point-to-point communications, and collective communications routines will be provided. Simple examples illustrating distributed memory computing, with the use of common MPI routines, will be covered. The OpenMP section will provide an overview of constructs and directives for specifying parallel regions, work sharing, synchronization and data scope. Simple examples will be used to illustrate the use of OpenMP shared-memory programming model, and important run time environment variables. Hands on exercises for both MPI and OpenMP will be done in C and FORTRAN.	Machine Learning Overview Mai Nguyen, Lead for Data Analytics, SDSC Paul Rodriguez, Research Analyst, SDSC Machine learning is an interdisciplinary field focused on the study and construction of computer systems that can learn from data without being explicitly programmed. Machine learning techniques can be used to uncover patterns in your data and gain insights into your problem. This session provides an overview of the fundamental machine learning algorithms and techniques used to explore, analyze, and leverage data to construct data-driven solutions applicable to any domain. Topics covered include the machine learning process, data exploration, data preparation, classification, and cluster analysis. Concepts and algorithms will be introduced, followed by exercises to allow hands-on experience using R and RStudio.
12:00 – 1:30	Lunch at Café Ventanas	
Session 4 1:30 – 5:00PM	Python for HPC Andrea Zonca, Senior Computational Scientist, SDSC In this session, we will introduce four key technologies in the Python ecosystem that provide significant benefits for scientific applications run in supercomputing environments. Previous Python experience is recommended but not required. (1) The Jupyter Notebook allows users to execute code on a single compute node or cluster and export the Python web interface to the local browser for interactive data exploration and visualization. The Jupyter Notebook supports live Python code, explanatory text, LaTeX equations and plots in the same document. (2) IPython Parallel provides a simple, flexible and scalable way of running thousands of Python serial jobs by spawning IPython kernels (namely engines) on any HPC batch scheduler. It also allows interactive control of the engines from a Jupyter Notebook session along with the ability to submit more Python tasks to the engines. (3) Numba makes it possible to run pure Python code on GPUs simply by decorating functions with the data types of the input and output arguments. Pure Python prototype code can be gradually optimized by pushing the most computationally intensive functions to the GPU without the need to implement code in CUDA or OpenCL. (4) Dask is a flexible parallel computing library that allows to build a distributed computation using simple operators and then let the library automatically handle distributing data, executing the computation hierarchically and gather back the results.	Scalable Machine Learning Mai Nguyen, Lead for Data Analytics, SDSC Paul Rodriguez, Research Analyst, SDSC Machine learning is an integral part of knowledge discovery in a wide variety of applications. From scientific domains to social media analytics, the data that needs to be analyzed has become massive and complex. This session provides an introduction to approaches that can be used to perform machine learning at scale. Tools and procedures for executing machine learning techniques on HPC will be presented. Spark will also be covered. In particular, we will use Spark’s machine learning library, MLlib, to demonstrate how distributed computing can be used to provide scalable machine learning. Please note: Knowledge of fundamental machine learning algorithms and techniques is required. (See description for Machine Learning Overview.)

5:30 – 9:00PM	Beach BBQ Dinner at La Jolla Shores Hotel , <i>sweater or jacket recommended</i> 8110 Camino Del Oro, La Jolla, CA 92037 Shuttle provided from SDSC driveway
------------------	---

FRIDAY, August 4	
8:00 – 8:30	Coffee
8:30 – 9:30	Emerging Technologies in HPC Pietro Cicotti, Senior Computational Scientist, SDSC
9:30 – 11:00	Lightning Rounds
11:00 – 11:30	Wrap up
11:30AM	Adjourn Thank you for attending we hope you enjoyed the week! <i>(To-go box lunches will be available)</i>