

Introduction to SciNet

SciNet HPC Consortium
Compute Canada

May 9, 2012

Don't Panic

Outline

- ① About SciNet
 - ▶ SciNet is ...
 - ▶ How to get an account
- ② SciNet systems
- ③ Using SciNet
 - ▶ Software/Libraries
 - ▶ Compilers
 - ▶ Job submission
- ④ Data management
- ⑤ Final tips

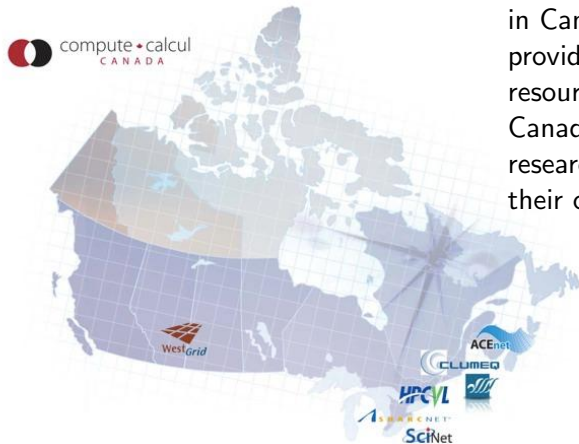
Part I

ABOUT SCINET

SciNet is ...

... a consortium for High-Performance Computing consisting of researchers at U. of T. and its associated hospitals.

One of 7 consortia in Canada that provide HPC resources to Canadian academic researchers and their collaborators.



SciNet is ...

... home to the largest supercomputer in Canada (and a bunch more)...



SciNet is ...

- where you go for courses on a wide range of computational topics.
Examples:
 - ▶ Intro to GPGPU with CUDA
 - ▶ Intro to Scientific Programming with Modern FORTRAN
 - ▶ Intro to Scientific Programming with C++
 - ▶ Parallel I/O
 - ▶ Scientific Computing Course (for credit for physics/astrophysics grads)
 - ▶ This Intro to SciNet!
- recognized by NVIDIA as both a CUDA research and teaching centre



- where you can meet with other users, and present your work, at the monthly SciNet User Group Meetings (SNUGs).

SciNet is ...

... 4 technical analysts who can work directly with you to use our resources to produce good science.

- Jonathan Dursi
- Scott Northrup
- Ramses van Zon
- Daniel Gruner



+ 7 people that make sure everything runs smoothly.

- Jaime Pinto
- Joseph Chen
- Jason Chong
- Ching-Hsing Yu
- Neil Knecht
- Leslie Groer
- Chris Loken

- + Technical director Prof. Richard Peltier
- + Business manager Teresa Henriques
- + Project coordinator Jillian Dempsey

How to get an account

Any qualified researcher at a Canadian university can get a SciNet account through the following process:

- 1 Register for a Compute Canada Database (CCDB) account
- 2 Non-faculty need a sponsor (supervisor's CCRI number), who has to have a SciNet account already.
- 3 Login to CCDB and apply for a SciNet account (click *Apply* beside SciNet on the *Consortium Accounts* page)
- 4 Agree to the Acceptable Usage Policy (e.g., don't share account, respect others, we can monitor your jobs)

How to get an account

Default account

- Allows to simultaneously run 32 jobs of 48 hours wall-time.
- Use of maximally 256 cores.
- So in a year you could use over 2 million core-hours.
- 10GB of storage (plus a bunch of temporary scratch storage)

How to get more resources on SciNet

- Users who will be needing more than the default amount of resources (compute cycles and/or storage) must have their PI apply for it through the competitively awarded.
- Resource national competition occurs in the fall of each year.
- Having an allocation also increases your priority in the queue.

Part II

SCINET SYSTEMS

Resources at SciNet

General Purpose Cluster (GPC)

- 3864 nodes with 8 Intel x86-64 cores@2.53/2.66GHz
- 328 TFlops max from 30,912 cores in total.
- 16 GB RAM per node (~ 14 GB for user jobs)
- 16 threads per node
- CentOS 6 operating system
- InfiniBand network (MPI,IO,...): $\frac{1}{4}$ DDR, $\frac{3}{4}$ QDR 5:1 (since April 19, 2012)
- #16 on the June 2009 *TOP500* supercomputer sites
- #1 in Canada

Tightly Coupled System (TCS)

- 102 nodes with 32 Power-6 cores @ 4.7GHz
- 62 TFlops max from 3264 cores
- 128 GB RAM per node
- 64 threads per node
- AIX operating system
- Interconnected by full non-blocking InfiniBand
- #80 on the June 2009 *TOP500* supercomputer sites
- *Access disabled by default. For access, email us explaining the nature of your work. Your application should scale well to over 32 procs.*

Resources at SciNet

Power 7 linux cluster (P7)



- 5 nodes with 32 Power-7 cores @ 3.3GHz
- 4.2 TFlops max from 160 cores
- 128 GB RAM per node
- 128 threads per node
- RedHat 6 operating system
- Interconnected by InfiniBand
- Accessable to TCS users

Accelerator Research Cluster (ARC)

8 GPU devel nodes and 16 NVIDIA Tesla M2070. Per node:

- 8 Intel cores (Xeon X5550) @ 2.67GHz
- 48 GB RAM
- CentOS 6 operating system
- Interconnected by DDR InfiniBand
- 2 × GPUs with Cuda capability 2.0 (Fermi) each with 448 Cuda cores @ 1.15GHz and 6 GB of RAM.

From CPUs: 683.52 GFlops max from 64 cores

From GPUs: 8.24 TFlops (double prec) from 7168 Cuda cores

Access upon request.



Disk space

- 1790 1TB disk drives, for a total of 1.4 PB of storage
- Two DCS9900 couplets, each delivering 4-5GB/s (r/w)
- Single *shared* file system GPFS on all systems
- I/O goes the infiniband network since April 19, 2012
- Your files go in /home/g/group/user and /scratch/g/group/user.

Storage space

- HPSS: Tape-backed storage expansion solution.

location	quota	block-size	time-limit	backup	devel	comp
/home	10GB	256kB	perpetual	yes	rw	ro
/scratch	20TB/1M	4MB	3 months	no	rw	rw

Part III

USING SCINET

Get on the system

1. Access systems: login.scinet.utoronto.ca

First ssh to login (not part of clusters):

```
ssh -l <username> login.scinet.utoronto.ca [-Y]
```

The login nodes are gateways, they are not part of any of the clusters and they are only to be used for small data transfer and to proceed logging into one of the devel nodes.

2. Go to the right cluster: ssh to the devel nodes

GPC: gpc01, gpc02, gpc03, gpc04

ARC: arc01

TCS: tcs01, tcs02

P7: p701

Software and Libraries

Once you log into devel nodes, what software is already installed?

- Other than essentials, all software installed as modules.
- modules set environment variables `LD_LIBRARY_PATH`, `PATH`, ...)
- Allows multiple, conflicting versions of package to be available.
- More on *Software and Libraries* page of wiki.

```
gpc-f103n084-$ module avail
```

```
-----
```

```
/scinet/gpc/Modules6/Modules/version_indepen  
3.2.6
```

```
-----
```

```
/scinet/gpc/Modules6/Modules/3.2.6/modulefiles  
dot                modules          use.own  
module-cvs         use.deprecated  
module-info        use.experimental
```

```
----- /scinet/gpc/Modules6/Modules/modulefiles  
ROOT/5.26.00  
Xlibraries/X11-32  
Xlibraries/X11-64(default)  
amber10/amber10  
autoconf/autoconf-2.64  
blast/2.2.23+  
cmake/2.8.0  
...
```

<http://wiki.scinet.utoronto.ca>

Software and Libraries

<code>module load <module-name></code>	use particular software
<code>module purge</code>	remove currently loaded modules
<code>module avail</code>	list available software packages
<code>module list</code>	list loaded modules
<code>module help <module-name></code>	describe a module

- Load frequently used modules in `.bashrc` in home directory.
- Load run-specific modules inside the job script.
- Short name gives default (e.g. `intel` → `intel/12.1.3`)
- To compile code that uses that a library from a module, add

`-I${SCINET_[shortmodulename]_INC}`

- To link, add

`-L${SCINET_[shortmodulename]_LIB}`

Software and Libraries

Dependencies

- Modules sometimes require other modules to be loaded first.
- Module will let you know if you didn't.
- For example:

```
gpc-f103n084-$ module purge
gpc-f103n084-$ module load python
python/2.6.2(11):ERROR:151: Module 'python/2.6.2' depends on one of
the module(s) 'gcc/4.6.1'
python/2.6.2(11):ERROR:102: Tcl command execution failed:  prereq gcc/4.6.1
gpc-f103n084-$ module load gcc python
gpc-f103n084-$ module list
Currently Loaded Modulefiles:
      1) gcc/4.6.1          2) python/2.6.2
gpc-f103n084-$
```

Commercial Software?

- SciNet has an extremely large and broad user base (~ 1000 users)
- ⇒ *Cannot buy everyone's favourite commercial software package.*
- Only commercial software we have installed is software that can benefit everyone:
 - ▶ GPC and ARC: Intel compilers, MKL (both in module intel)
 - ▶ TCS and P7: IBM compilers, ESSL
 - ▶ GPC, ARC, TCS and P7: DDT debugger from Allinea
 - No Matlab, Gaussian, IDL, ...
(but Octave)
 - Can work with you to install commercial software for which you have a license.

Compiling on SciNet systems

GPC compilers

- From `login.scinet.utoronto.ca`, ssh to one of the four devel nodes.

```
ssh gpc04 [-Y]
```

or

```
gpcdev -Y
```

- We recommend Intel compilers, which are

```
icc, icpc, ifort
```

for C, C++, and Fortran, respectively (from the module `intel`)

- Optimize your code for the GPC machine using of at least

```
-O3 -xhost
```

- Add `-openmp` to the command line for OpenMP
- Compile MPI code with `mpif77/mpif90/mpicc/mpicxx`.

- 1 Open MPI, in module `openmpi` (v1.4.4)
- 2 Intel MPI, in module `intelmpl` (v4.0.3)

Compiling on SciNet systems

TCS compilers and P7 compilers

- ssh to a devel node

`ssh tcs01` or `ssh tcs02` , or `ssh p701`

- Use IBM compilers: `xlc`, `xlC`, `xlF` for C, C++, and Fortran
(`module load vacpp xlf` for the latest version)
- For OpenMP, use `xlc_r`, `xlC_r`, `xlF_r`.
- For MPI, `mpicc`, `mpCC`, `mpxlf` are the mpi wrappers.
- Suggested compiler flags:

`-q64 -O3 -qhot -qarch=auto -qtune=auto`

supplemented by `-qsmp=omp` for OpenMP programs.

- On the link line we suggest using

`-q64 -bdatapsize:64k -bstacksize:64k`

also supplemented by `-qsmp=omp` for OpenMP programs.

Compiling on SciNet systems

ARC compilation

- From `login.scinet.utoronto.ca`, ssh to devel node
`ssh arc01`
- The NVIDIA cuda compiler is available (3.2 and 4.0).
You can `module load cuda/3.2` or `cuda/4.0`.
The compiler is called `nvcc`.
- Optimize for the Tesla GPUs using the compiler flags
`-arch=sm_13 -O3`
- As of version 3.0, OpenCL is included in the CUDA Toolkit so loading the CUDA module is all that is required.
- Debuggers: `cuda-gdb` or `ddt` (`module load ddt`)

Submitting jobs

SciNet=shared resource

- To run a job, you must submit to a batch queue.
- You submit jobs from a devel node in the form of a script
- Scheduling is by node!
- Best to run from the scratch directory (home=read-only)
- Copy essential results out after your runs have finished.

Limits

- Group based limits:
possible for your colleagues to exhaust group limits
- Talk to us first to run massively parallel jobs (> 2048 cores)
- While their resources last, jobs will run at a higher priority than others for groups with an allocation.

GPC queues

queue	time(hrs)	max jobs	max cores
batch	48	32/1000	256/8000 (512/16000 threads)
debug	2/0.5	1	16/64 (32/128 threads)
largemem	48	1	16 (32 threads)

- Submit to these queues with

```
qsub [options] <script>
```

- Common options (usually in script):

-l: specifies requested nodes and time, e.g.

```
-l nodes=2:ppn=8,walltime=1:00:00
```

```
-l nodes=2:qdr:ppn=8,walltime=1:00:00
```

-q: specifies the queue, e.g.

```
-q batch
```

```
-q debug
```

```
-q largemem
```

-I specifies that you want an interactive session.

GPC job script example

```
#!/bin/bash

#PBS -l nodes=1:ppn=8
#PBS -l walltime=1:00:00
#PBS -N simple-omp-job
cd $PBS_O_WORKDIR
export OMP_NUM_THREADS=8
./omp_example > output
```

GPC queues

- GPC HyperThreading: Appears as if there are 16 processors rather than 8 per node. For OpenMP applications this is the default unless OMP_NUM_THREADS is set. For MPI, try `-np 16`.
- *Always first test if this is beneficial and feasible!*
- Once the job is incorporated into the queue (this takes a minute), you can use: `showq` to show the queue, and job-specific commands such as `showstart`, `checkjob`, `canceljob`
- Always request `ppn=8`, even with hyperthreading.
- The largemem queue is exceptional, in that it provides access to two nodes (only) that have 16 processors and 128GB of ram.
- There is no queue for serial jobs, so if you have serial jobs, **YOU** will have to bunch together 8 of them to use the node's full power. GNU Parallel can help you with that.

ARC queue

time(hrs)	max jobs	max cores
48	32	64cpu / 7168 gpu

- Submit to these queues with

```
qsub [options] <script>
```

- Common options:

-l: specifies requested nodes and time, e.g.

```
-l nodes=1:gpus=2:ppn=8,walltime=1:00:00
```

-I specifies that you want an interactive session.

- See submitted jobs with `qstat`.
- Limits on this system may change.

TCS/P7 queues

queue	time(hrs)	max jobs	max cores
verylong	48	2/25	64/800 (128/1600 threads)

Submitting is done with

```
llsubmit <script>
```

and `llq` shows the queue.

- The POWER processors have Simultaneous Multi Threading. Similar to HyperThreading.
- Once your job is in the queue, you can use `llq` to show the queue, and job-specific commands such as `llcancel`, `llhold`, ...
- *Do not run serial jobs on the TCS!*
- To make your jobs start sooner, reduce the `wall_clock_limit` to be closer to the estimated run time (perhaps adding about 10 % to be sure). Shorter jobs are scheduled sooner than longer ones.

Example 1 (GPC)

```
gpc-f101n084-$ module load intel openmpi
gpc-f101n084-$ mpif90 -O3 -xhost mycode.f90 -o mycode
gpc-f101n084-$ mkdir $SCRATCH/example1
gpc-f101n084-$ cp mycode $SCRATCH/example1
gpc-f101n084-$ cd $SCRATCH/example1
gpc-f101n084-$ cat > myjob.pbs
#!/bin/bash
#PBS -l nodes=8:ppn=8,walltime=1:00:00
#PBS -N JobName
cd $PBS_O_WORKDIR
module load intel openmpi
mpirun -np 64 ./mycode > out
gpc-f101n084-$ qsub myjob.pbs
2961983.gpc-sched
gpc-f101n084-$ qstat (or checkjob 2961983, or showq -u $USER)
  Job id              Name              User Time Use S Queue
  -----
  2961983.gpc-sched JobName              rzon              0 Q batch
gpc-f101n084-$ ls
JobName.e2961983  JobName.o2961983  mycode  myjob.pbs
out
```

Example 2 (GPC)

```
gpc-f101n084-$ module load intel
gpc-f101n084-$ ifort -O3 -xhost mycode.f90 -o mycode
gpc-f101n084-$ mkdir $SCRATCH/example2
gpc-f101n084-$ cp mycode $SCRATCH/example2
gpc-f101n084-$ cd $SCRATCH/example2
gpc-f101n084-$ cat > joblist.txt
```

```
mkdir run1; cd run1; ../mycode 1 > out
mkdir run2; cd run2; ../mycode 2 > out
mkdir run3; cd run3; ../mycode 3 > out
...
mkdir run64; cd run64; ../mycode 64 > out
```

```
gpc-f101n084-$ cat > myjob.pbs
#!/bin/bash
#PBS -l nodes=1:ppn=8,walltime=24:00:00
#PBS -N ASerialJob
cd $PBS_O_WORKDIR
module load intel gnu-parallel
parallel -j 8 < joblist.txt
```

```
gpc-f101n084-$ qsub myjob.pbs
2961985.gpc-sched
```

```
gpc-f101n084-$ ls
ASerialJob.e2961985  ASerialJob.o2961985  joblist.txt
myjob.pbs            run1/
```

```
mycode
run3/
```



Part IV

DATA MANAGEMENT

File system

SciNet \neq 4000 \times your pc

- Compute nodes do not contain hard drives!
- The available disk space, /home and /scratch, all part of the GPFS file system which runs over the network.
- GPFS is a high-performance file system which provides rapid reads and writes to large data sets in parallel from many nodes.
- It performs quite poorly at accessing data sets which consist of many, small files.
- Don't keep many small files on the system.
They waste space, and are slower to access, read and write.

I/O strategies

- Do not read and write lots of small amounts of data to disk.
Reading data in from one 4MB file can be enormously faster than from 100 40KB files.
- Write your data out in binary. Faster and takes less space.
- Each process writing to a file of its own is not scalable.
A directory gets locked by the first process accessing it, so the other processes have to wait for it.
- Consider using MPI-IO (part of the MPI-2 standard), (parallel) NetCDF, HDF5, or ADIOS.
- If you must read and write a lot to disk, use ramdisk if possible.
The ramdisk can be accessed using `/dev/shm/` and is currently set to 8GB max.
- Copy back from ramdisk at end of run.

Moving large data

Moving less than 10GB through the login nodes

- Only login nodes visible from outside SciNet (1Gb/s link).
- Use scp or rsync.
- but datamover1 node is faster.

Moving more than 10GB through the datamover1 node

- Should be done from the datamover1 node (10Gb/s link).
- From any SciNet node, ssh to datamover1.
- Transfers must originate from datamover1.
Cannot copy files from the outside world to datamover1.
- Your machine must be reachable from the outside.

Moving data to HPSS

- HPSS is a tape-based storage solution.
- Available to groups with a special allocation > 5TB.

Final tips

- Test your job's requirements and scaling behaviour.
Start runs on a small scale and work your way up to larger scales.
- Accurately specify the walltime when you submit a job.
- Avoid reading and writing lots of small amounts of data to disk.
- Do not create lots of files.
- Do not submit single serial jobs.
- Do not keep lots of files in your directory (use tar).
- Read the SciNet User Guide
http://wiki.scinethpc.ca/wiki/images/5/54/SciNet_Tutorial.pdf

Useful web sites

Portal: <https://portal.scinet.utoronto.ca>

SciNet usage reports

Change password, default allocation, maillist subscriptions



Exception: gpgpu mail list:

<https://support.scinet.utoronto.ca/mailman/listinfo/scinet-gpgpu>

Useful web sites

Wiki: <http://wiki.scinet.utoronto.ca>

The screenshot shows a web browser window displaying the SciNet User Support Library wiki. The browser's address bar shows the URL https://support.scinet.utoronto.ca/wiki/index.php/SciNet_User_Support_Library. The page has a blue header with the SciNet logo and navigation tabs for 'page', 'discussion', 'view source', and 'history'. A search bar is located on the left side. The main content area is titled 'SciNet User Support Library' and includes a welcome message. Below this, there are several sections: 'System Status: Normal' with a message last updated on Mon Aug 16 13:12:08 EDT 2010; 'QuickStart Guides' listing various topics like Login essentials, GPC, TCS, Data management, Software and libraries, Job scheduling system (Moab), Performance primer, Tutorials and Manuals, FAQ, SciNet User Tutorial (UPDATED!), Usage policy, and Acknowledging SciNet; 'What's New On The Wiki' listing recent updates like SciNet class schedule, updated SciNet User Tutorial, SSH keys, disk space, and Hyperthreading; 'User-Supplied Content'; and 'News and Recent Events'. A sidebar on the left contains links to various SciNet resources, including Wiki main page, Accounts, User Support, Acknowledgment, Current events, System status, systems (Overview, Essentials, Software/libraries, Scheduler, Managing data), gpc (Quickstart, Software, MPI, Performance), tcs (Quickstart, Software, Performance), and knowledge base. The SciNet logo is also present in the bottom right corner of the page.

SciNet User Support Library

Welcome to the SciNet User Support wiki. Here you will find up-to-date manuals put together by SciNet staff and users, as well as links to external resources, to help you make use of SciNet resources for computational scientific discovery. Navigate using the links in the right style presentation below, or using the menu on the left.

System Status: Normal

message last updated on: Mon Aug 16 13:12:08 EDT 2010
(Previous messages)

QuickStart Guides

- Login essentials
- GPC: General Purpose Cluster
- TCS: Tightly Coupled System
- Data management
- Software and libraries
- Job scheduling system (Moab)
- Performance primer, and performance tools for GPC and TCS
- Tutorials and Manuals
- FAQ (frequently asked questions)
- SciNet User Tutorial **UPDATED!**
- Usage policy
- Acknowledging SciNet

What's New On The Wiki

- SciNet class schedule announced (rzon, 30 August)
- Updated and improved SciNet User Tutorial (rzon, 28 August)
- SSH keys and SciNet (jdursi, 19 August)
- How to find out how much disk space you're using and how much you have left (pinto, 10 August)
- Hyperthreading with Gromacs (cneale, 9 August)

Previous new stuff can be found in the [What's new archive](#).

User-Supplied Content

News and Recent Events

SciNet class schedule

Useful web sites

Courses: <https://support.scinet.utoronto.ca/courses>

The screenshot shows the SciNet Courses website interface. At the top, there's a navigation bar with links: All Classes (by date), SNUG Meetings, Short Courses, and Courses Forum. The main header features the SciNet logo and the text 'SciNet Courses - High Performance Education'. Below this, there's a sub-header 'SciNet Courses - High Performance Education'.

User login

Username: *

Password: *

Event Calendar

« September 2010 »

Sun	Mon	Tue	Wed	Thu	Fri	Sat
			1	2	3	4
5	6	7	8	9	10	11
12	13	14	15	16	17	18
19	20	21	22	23	24	25
26	27	28	29	30		

Current signups for

You have not signed up for any classes yet.

Welcome to the SciNet Courses Website!

Wed, 2010-09-01 14:47 — admin

Here we will post information on upcoming classes and meetings, and providing a spot where you can sign up for them. You can follow [our RSS feed](#), view all the events on a [calendar](#), or just watch the SciNet users mailing list. Course materials will usually be posted to [the Wiki](#) soon after the courses.

[admin's blog](#)

November SNUG - TechTalk: Debuggers

Mon, 2010-08-30 22:15 — admin

Start: 2010-11-10 12:00

End: 2010-11-10 13:00

Timezone: America/Montreal

The SciNet Users Group (SNUG) meetings are every month on the second Wednesday, and involve pizza, user discussion, feedback, and a half-hour talk on topics or technologies of interest to the SciNet community. November's SNUG will be on the 10th, and the TechTalk will be: "Debuggers & parallel

Upcoming events

- Intro To SciNet (3 hours)
- September SNUG - TechTalk: GPFS (5 days)
- Intro To SciNet (7 days)
- Intro to Parallel Programming (19 days)
- Parallel I/O (33 days)
- October SNUG - TechTalk: Version Control (40 days)

[more](#)

Poll

What other courses would you most like to attend?

SciNet
compute + calcul
CANADA

Links

Wiki: <http://wiki.scinet.utoronto.ca>

Courses: <https://support.scinet.utoronto.ca/courses>

Portal: <https://portal.scinet.utoronto.ca>

Technical support: support@scinet.utoronto.ca