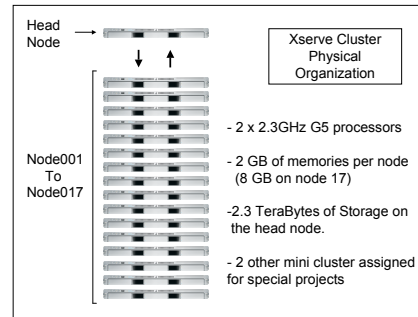


Bioinformatics Facility of the Biotechnology

The Do and Don't of the Xserve Cluster

Pascal.Lapierre@uconn.edu

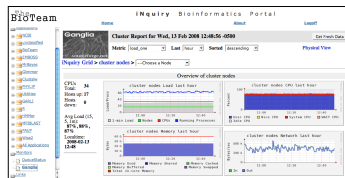


Basic Rules

- For research purpose only. Not a place to put your favorite MP3 or backup your HD.
- Do not overload the systems. It is ok to use ~6 nodes in period of low activities but when it gets busy, limit yourself to only 2-3 nodes if absolutely necessary.
- Always keep track of your jobs. Don't let analyses running unattended for months.
- Use the queue system whenever you can.
- Do not run jobs on the Head node.

Remote Access

- Via SSH or Web Interface



Useful Commands

(Help page available at : http://137.99.46.188/wiki/index.php/Main_Page)

- **qstat** : Shows the current status of the available Grid Engine queues and the jobs associated with the queues.
- **ls** : List directory contents
- **ps** : Display the process status. Allow to get process ID.
 - ps ux : Displays your process only
 - ps aux : displays all the process running on the node
- **du** : display disk usage statistics. Use du -h for a readable output

Useful Commands (cont)

- **mkdir** and **rmdir** : create and remove directories
- **cp** : copy files
- **mv** : moved files (can be used to rename files)
- **rm** : remove files. rm -r to remove files and sub-directories
- **kill** : to kill a running process. Kill -9 'proc_id'

The queue system

"Managing Workload by Managing Resources and Policies"

- **qstat** : Display the queue status.
- **qssh** : Queue remote shell. Automatically select an available node to log on.
- **qsub** : Queue submit. Automatically submit a job to an available node. Used in conjunction with a shell script (see next slide).
- **qdel** : Delete a job running in the queue.
 - qdel - process_ID

How to submit a job using qsub?

A shell script is just a small text file pointing to what you want to run in the queue.

For example, if I want to submit a perl script (phym1.pl), I will create a text file name phym1.sh :

```
#!/bin/sh
cd /Users/nucleus/evolver
perl phym1_trees1.pl
#end of script
```

To submit the shell :
- qsub phym1.sh

Things to be cautious :

-While highly reliable, the cluster might sometimes run into problems and needed to be rebooted. This will cause to loose all the processes that were running at the time. Try to think of ways to break up or save at different stage of your analyses.

-The NFS (Network File System) has temporary amnesia when overwhelmed. The system will forget to write part of the output files. A workaround is to save to the scratch drive of the individual nodes (cd /scratch).

Tricks that I have learned

In Perl, Array of Arrays are useful for grid-like manipulations of data :

```
Seq.txt =
MRRRAIATNQQ
MRLAIISRQD
MRLIISRQD
MRLAIISRQD
MRLAIISRQD

0123456789
MRRRAIATNQQ
1 MRLAIISRQD
2 MRLIISRQD
3 MRLAIISRQD
```

Print \$matrix[2][4]; → S

Retrieving data directly from NCBI using E-tools

- fastacmd -s 49220 -d nr →

You can use E-tools to get the Genbank file for 49220 (http://www.ncbi.nlm.nih.gov/entrez/query/static/entrez_help.html)

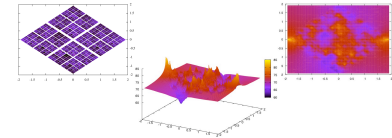
-read-tax.pl :

Database [PubMed]: protein
Query [zanzibar]: 49220
Report [abstract]: genbank

Gnuplot

(<http://www.gnuplot.info>)

- Great at generating plots on the fly using Perl
- Can handle enormous datasets
- Easy to use, very powerful



Gnuplot (cont)

- Only works if you are on the Head node!
- You can install and use it on you personal computer (PC/MAC)

```
#!/usr/bin/perl
open(OUT, ">gnup.out");
print OUT "set term postscript color";
print OUT "set output 'file_%.5p.ps'";
print OUT "set title 'file_%.5p'";
print OUT "set xlabel 'YearMin'";
print OUT "set ylabel 'FracMin'";
print OUT "set xrange [0 to 1]";
print OUT "set yrange [0 to 1]";
print OUT "plot 'gnuplot.gnu.out' using 'x' with lines 'M' title 'M'";
system("gnuplot gnup.out");
...
```

Old Assignments

- 1) Review PSblast [Your questions?](#)
 - 2) Write a 3 sentence outline for your student project [Send me an email on this!](#)
 - 3) Re-read chapter 2 p32 - p34 on control structures and page 142 -146 on for, foreach, and while loops
- For next week:
- Background: @a=(0..50);
 - #assigns numbers from 0 to 50 to an array, so that \$a[0]=0; \$a[1]=1; \$a[50]=50
- 4) Write perlscripts that add all numbers from 1 to 50. Try to do this using at least 2 different control structures.

Control structures: Sum 1..50

```
#!/usr/bin/perl
$sum=0;
$count=0;
while ($count <= 50) {
    $count++; #this is tricky in the last loop $count is 49 and then increased to 50 and added
    $sum += $count;
};
print "$sum\n";

#!/usr/bin/perl/
$sum=0;
$count=0;
for ($count =0; $count <= 50; $count++) {
    $sum=$sum+$count;
    $sum += $count;
};
print "$sum\n";
```

Control structures: Sum 1..50

```
#!/usr/bin/perl/
$sum=0;
@array = (1..50);
foreach (@array) {
    # $sum=$sum+$_;
    $sum += $_;
};
print "$sum\n";

#!/usr/bin/perl/
$sum=0;
$count=0;
while () {
    $sum += $count;
    $count++;
    if ($count >= 50) {last};
};
print "$sum\n";

foreach ( ) { };
while ( ) { };
if ( ) {last};
```

Control structures: Sum 1..50

```
#!/usr/bin/perl
$sum=0;
@array = (1..50);
while (defined($array[$count])) {
    $sum += $array[$count];
    $count ++;
    #print "$array[$count] \t $sum\n";
};
print "$sum\n";

#!/usr/bin/perl -w
$sum=0;
@array = (0..50);
for ($count=1; ($count<=51); $count++){
    $sum += $array[$count];
    # $temp=$array[$count];
    #print "$count=$count sum is $temp\n";
};
print "$sum\n";

for ( , , ) { }
Counting elements of an array
Could have started at 0
```

For Wednesday

- Email your 3 sentence project outline
- Read NCBI info on [geneprot \(here\)](#)
- Try [geneplot](#) comparing your favorite genomes (e.g. [here](#))
- What might be a problem using [geneplot](#)?

For Monday

- Read chapter 3 in Learning Perl
- Write a script that reads in a sequence and prints out the reverse complement.
- Modify your script to that it can handle a sequence that goes over several lines?
- ```
Background: $comp =- tr/ATGC/TACG;
#translates every A in $comp into a T; every T into an A; every G into a C and every C into a G
```



PSIBlast to find transposase homologs

- Download transposase sequence transposase.fa
- Download genome as nucleotide sequence
- Format genome

- formatdb -i Tpet.fna -p F -o T
- blastpgp -i transposase.fa -d nr -I T -h 0.00001 -j 6 -C transposase.chk -a2
- blastall -i transposase.fa -d Tpet.fna -p psitblastn -R transposase.chk -o transposase\_Tpet.tab -a2 -m8 -F F

transposase\_Tpet.tab:

```
seqs 100000000 more transposase_Tpet.tab
#|
#| q115742470[sw][NC_009626.1] 15.54 436 334 9 11 423 46314 46754 3m-101 341
#| q115742470[sw][NC_009626.1] 15.46 354 354 5 10 197 44299 44248 4m-14 72.7
#| q115742470[sw][NC_009626.1] 12.01 424 335 17 5 292 194957 197451 2m-09 52.4
#| q115742470[sw][NC_009626.1] 19.20 125 92 5 249 264 1078742 1080333 4m-08 51.9
#| q115742470[sw][NC_009626.1] 12.19 293 247 12 1 292 469326 469944 2m-07 50.4
#| q115742470[sw][NC_009626.1] 14.41 178 132 6 140 317 1374637 1375151 3m-05 44.6
#| q115742470[sw][NC_009626.1] 10.79 370 345 4 144 244 238363 237663 3m-03 41.3
#| q115742470[sw][NC_009626.1] 10.12 273 199 14 189 26 121811 121762 3m-05 41.4
#| q115742470[sw][NC_009626.1] 12.34 340 291 4 87 343 2023448 202288 4.041 38.0
#| q115742470[sw][NC_009626.1] 11.25 140 125 7 257 389 1942355 194204 4.051 38.0
```