



# Ten Simple Rules

for getting started with  
command-line bioinformatics

Brandies PA, Hogg CJ (2021)  
PLOS Computational Biology 17(2): e1008645

# TEN SIMPLE RULES



01

COMPUTING  
TERMINOLOGY



02

TOOL/PIPELINE  
SELECTION



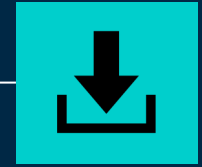
03

ESTIMATING  
RESOURCES



04

SELECTING  
PLATFORMS



05

SOFTWARE  
INSTALLATION



06

SCRIPT  
CURATION



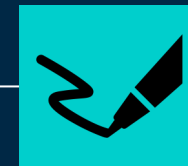
07

MONITORING &  
OPTIMISATION



08

FILE  
MANIPULATION



09

RECORD  
KEEPING



10

PATIENCE!



# COMPUTING TERMINOLOGY

01

# GET FAMILIAR WITH COMPUTER TERMINOLOGY

ALGORITHM

EXECUTABLE

HPC

DEPENDENCY

THREAD

CPU

SCHEDULER

VM

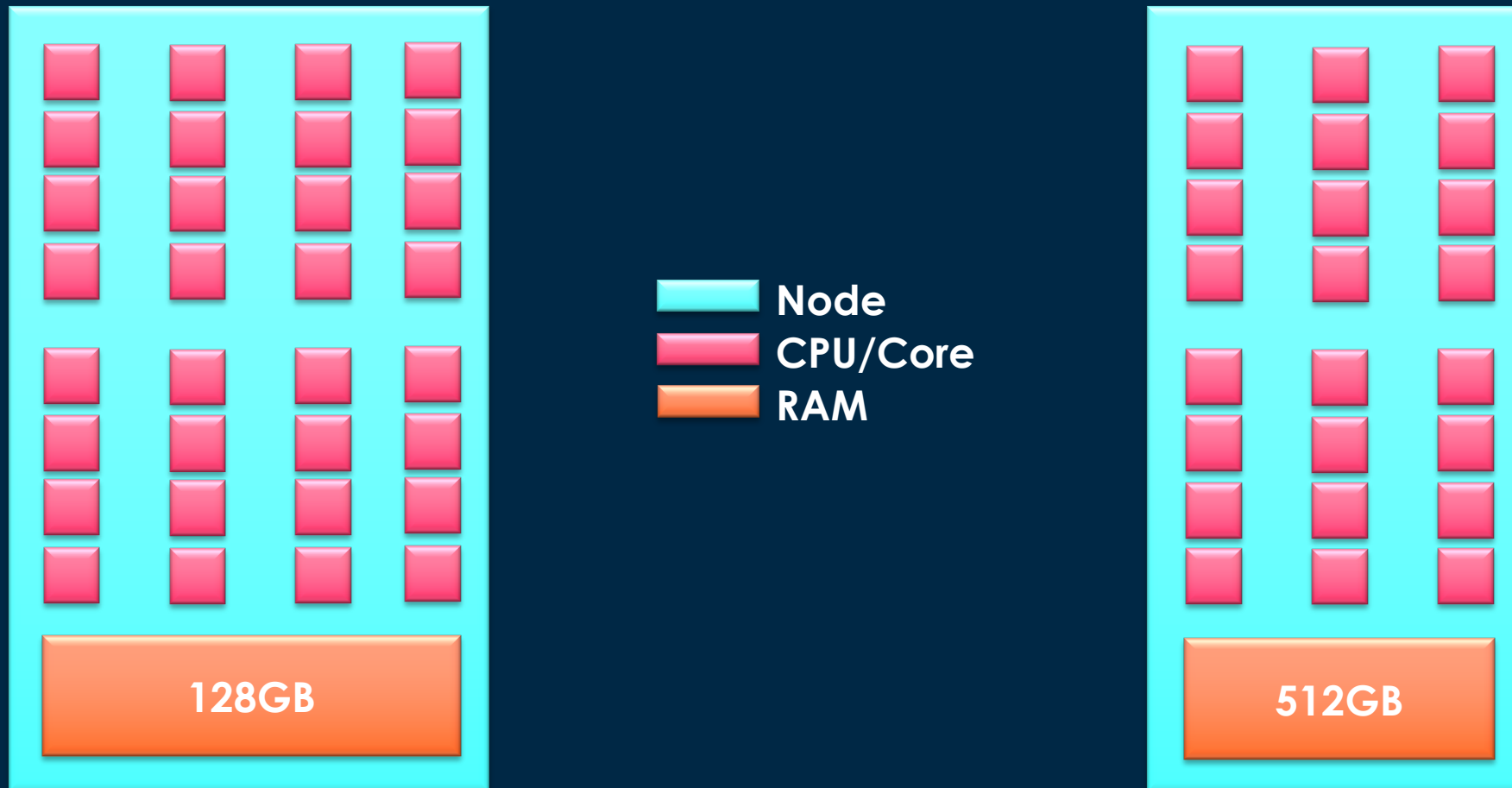
RAM

PIPELINE

WALLTIME

NODE

# GET FAMILIAR WITH COMPUTER TERMINOLOGY



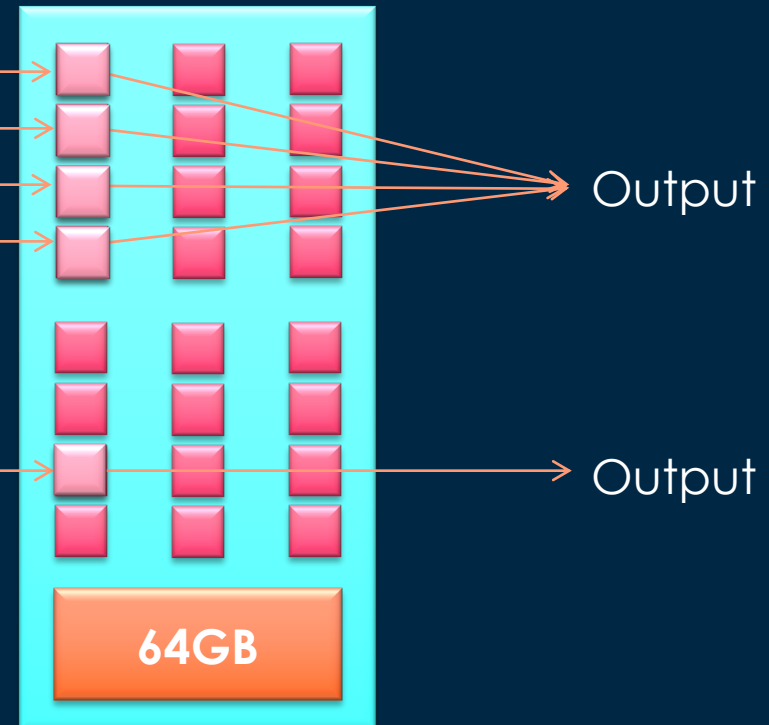
STANDARD

HIGH-MEMORY

# GET FAMILIAR WITH COMPUTER TERMINOLOGY

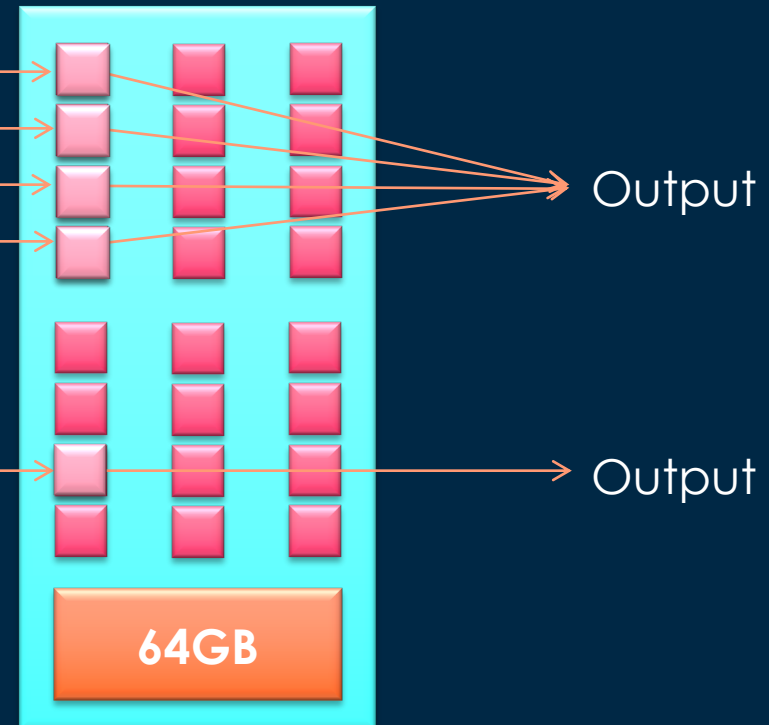
## Multi-Threaded Program:

- Process 1
- Process 2
- Process 3
- Process 4



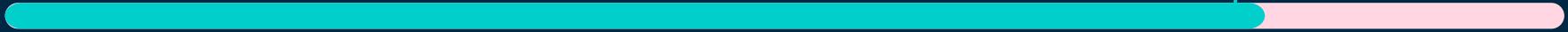
## Single-Threaded Program:

- Process 1 > Process 2 > Process 3 > Process 4



02

# TOOL/PIPELINE SELECTION



# KNOW YOUR DATA AND ASSESS YOUR NEEDS

- TARGET SPECIES/QUALITY OF DATA?
- AVAILABLE COMPUTING RESOURCES?
- AVAILABLE/TESTED TOOLS?







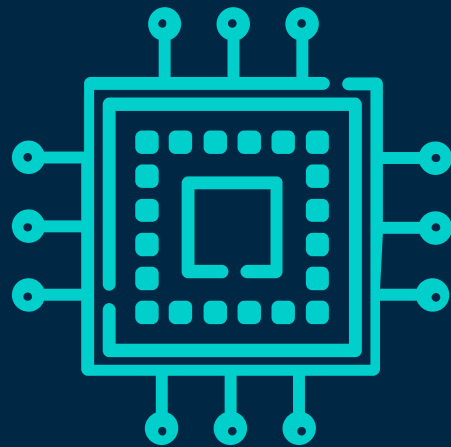


# ESTIMATING RESOURCES

03

# ESTIMATE YOUR COMPUTING REQUIREMENTS

- TOOL DOCUMENTATION + TEST DATASETS
- TOOL PUBLICATION + CITED PAPERS
- ASK COLLEAGES + USE ONLINE Q&A



32 CPUs  
128GB RAM  
1TB Storage

# ESTIMATE YOUR COMPUTING REQUIREMENTS

<https://github.com/AustralianBioCommons>

## IPA on Zeus @ Pawsey Supercomputing Centre

### Accessing workflow

The scripts for using [IPA](#) on Zeus have been made available below in the Quick start tutorial.

### Quick start tutorial

Below is a tutorial for install and annotation of IPA on Zeus.

**Note:** Purge\_Dups is not included below as it requires a manual install. Documentation is currently not available for this, but may be in the future.

### Install the improved phased assembler via bioconda

**Note:**

- the instructions below will install the latest version of miniconda and IPA.
- to [install a specific version](#) of IPA, use the following install script: `conda install pbipa=1.1.2`

It is recommended to use an interactive session on zeus to install IPA. To launch an interactive session use: `>salloc -n 1 -t 1:00:00`

1. Download miniconda: `curl -O https://repo.anaconda.com/miniconda/Miniconda3-latest-Linux-x86_64.sh`
2. Install miniconda: `sh Miniconda3-latest-Linux-x86_64.sh`
  - **Note:** During the setup steps you will be asked where to install Miniconda. DO NOT install under `$HOME`, the suggested directory for miniconda is `/group/<project>/<user>/miniconda3`. Whenever you need to use miniconda from the command line you will need to run `export PATH=$PATH:/group/.../miniconda3/bin`, and in some cases also run `exec bash`



# SELECTING PLATFORMS

04

# EXPLORE DIFFERENT COMPUTING OPTIONS

## CLOUD

Customised computing resources

Unshared resource (live programming)

Complete control of compute environment

Both free and paid commercial options

## SHARED HPC

Fixed computing resources

Shared resource (scheduled jobs)

Compute environment largely controlled by IT

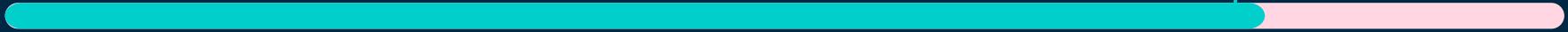
Usually freely available based on institution/merit

<https://ronin.cloud>



# SOFTWARE INSTALLATION

05



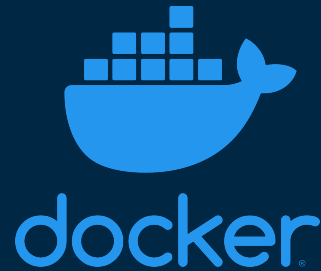
# UNDERSTAND THE BASICS OF SOFTWARE INSTALLATION

## PACKAGE MANAGERS

APT – DEBIAN

YUM – REDHAT

## CONTAINERS



## MANUAL INSTALL

Download +  
Unpack Code

Configure  
Software

Build Software

Install Software

## CONDA





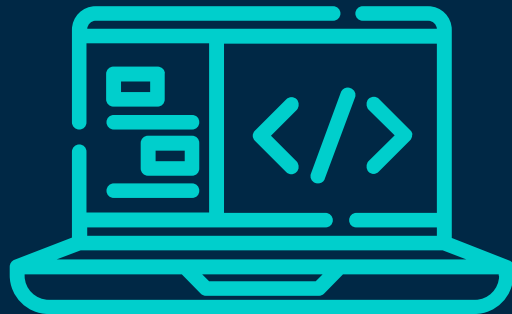


# SCRIPT CURATION

06

# CAREFULLY CURATE AND TEST YOUR SCRIPTS

- ✘ SYNTAX OR MINOR SCRIPTING ERRORS
- ✘ SOFTWARE INSTALLATION ERRORS
- ✘ DEPENDENCY ISSUES
- ✘ SUBOPTIMAL COMPUTE RESOURCES



WHEN WORKING WITH HPC ENVIRONMENTS, INTERACTIVE QUEUES ARE YOUR FRIEND!



# MONITORING & OPTIMISATION

07

# MONITOR AND OPTIMISE YOUR PIPELINES

## CLOUD VM

- HTOP  
Simple real-time monitoring of machine resources
- NETDATA  
Extensive monitoring using hundreds of preconfigured charts

## HPC

- SCHEDULER LOGS  
Summary of maximum and total compute resources used

# MONITOR AND OPTIMISE YOUR PIPELINES

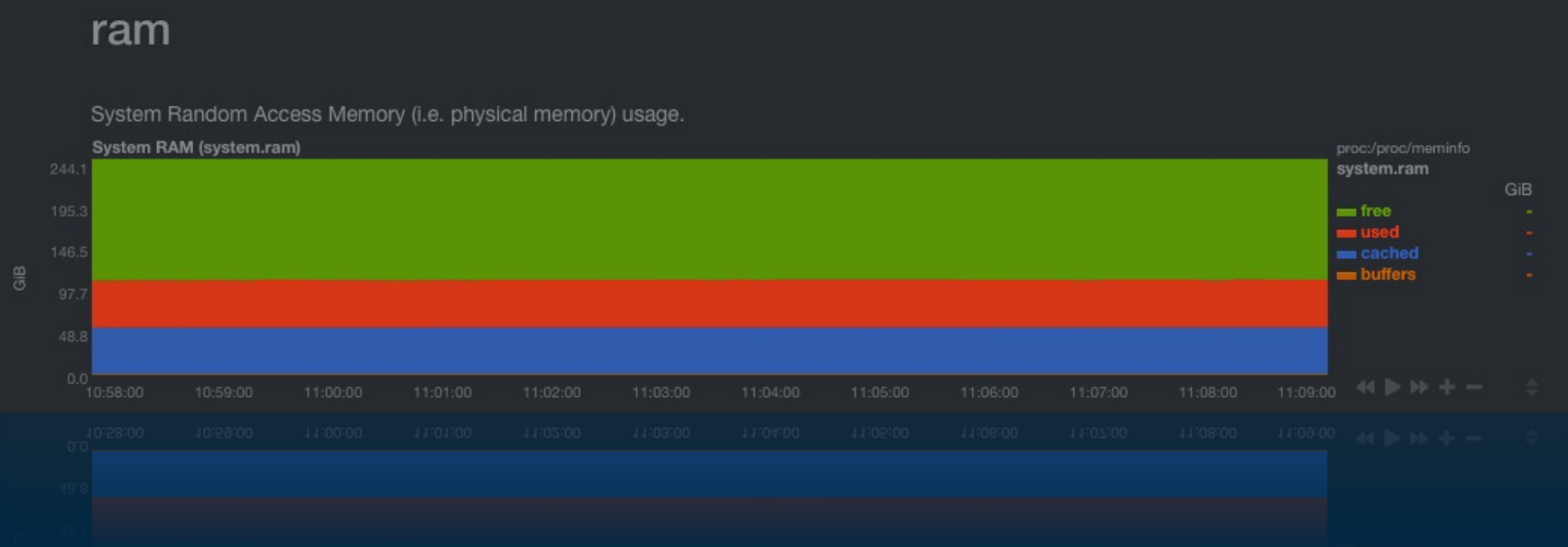
## HTOP

```
1 [|||||] 100.0% 25 [|||||] 100.0% 49 [|||||] 100.0% 73 [|||||] 100.0%
2 [|||||] 100.0% 26 [|||||] 100.0% 50 [|||||] 100.0% 74 [|||||] 100.0%
3 [|||||] 100.0% 27 [|||||] 100.0% 51 [|||||] 100.0% 75 [|||||] 100.0%
4 [|||||] 100.0% 28 [|||||] 100.0% 52 [|||||] 100.0% 76 [|||||] 100.0%
5 [|||||] 100.0% 29 [|||||] 100.0% 53 [|||||] 100.0% 77 [|||||] 100.0%
6 [|||||] 100.0% 30 [|||||] 100.0% 54 [|||||] 100.0% 78 [|||||] 100.0%
7 [|||||] 100.0% 31 [|||||] 100.0% 55 [|||||] 100.0% 79 [|||||] 100.0%
8 [|||||] 100.0% 32 [|||||] 100.0% 56 [|||||] 100.0% 80 [|||||] 100.0%
9 [|||||] 100.0% 33 [|||||] 100.0% 57 [|||||] 100.0% 81 [|||||] 100.0%
10 [|||||] 100.0% 34 [|||||] 100.0% 58 [|||||] 100.0% 82 [|||||] 100.0%
11 [|||||] 100.0% 35 [|||||] 100.0% 59 [|||||] 100.0% 83 [|||||] 100.0%
12 [|||||] 100.0% 36 [|||||] 100.0% 60 [|||||] 100.0% 84 [|||||] 100.0%
13 [|||||] 100.0% 37 [|||||] 100.0% 61 [|||||] 100.0% 85 [|||||] 100.0%
14 [|||||] 100.0% 38 [|||||] 100.0% 62 [|||||] 100.0% 86 [|||||] 100.0%
15 [|||||] 100.0% 39 [|||||] 100.0% 63 [|||||] 99.4% 87 [|||||] 100.0%
16 [|||||] 100.0% 40 [|||||] 100.0% 64 [|||||] 100.0% 88 [|||||] 100.0%
17 [|||||] 100.0% 41 [|||||] 100.0% 65 [|||||] 100.0% 89 [|||||] 99.4%
18 [|||||] 100.0% 42 [|||||] 100.0% 66 [|||||] 100.0% 90 [|||||] 100.0%
19 [|||||] 100.0% 43 [|||||] 100.0% 67 [|||||] 100.0% 91 [|||||] 100.0%
20 [|||||] 100.0% 44 [|||||] 100.0% 68 [|||||] 100.0% 92 [|||||] 100.0%
21 [|||||] 100.0% 45 [|||||] 100.0% 69 [|||||] 100.0% 93 [|||||] 100.0%
22 [|||||] 100.0% 46 [|||||] 100.0% 70 [|||||] 100.0% 94 [|||||] 100.0%
23 [|||||] 100.0% 47 [|||||] 100.0% 71 [|||||] 99.4% 95 [|||||] 100.0%
24 [|||||] 100.0% 48 [|||||] 100.0% 72 [|||||] 100.0% 96 [|||||] 100.0%
Mem [|||||] 71.5G/748G Tasks: 608, 157 thr; 96 running
Swp [|||||] 0K/0K Load average: 75.69 26.71 10.58
Uptime: 00:37:23

PID USER PRI NI VIRT RES SHR S CPU% MEM% TIME+ Command
8072 root 20 0 762M 727M 5180 R 81.3 0.1 0:05.53 /usr/bin/perl /usr/local/RepeatMasker/RepeatMasker /mnt/volume1/maker/bil01gen.dipnr.pri_rnd1.m
7592 root 20 0 752M 717M 5152 R 79.5 0.1 0:04.94 /usr/bin/perl /usr/local/RepeatMasker/RepeatMasker /mnt/volume1/maker/bil01gen.dipnr.pri_rnd1.m
8142 root 20 0 740M 705M 5248 R 78.9 0.1 0:04.40 /usr/bin/perl /usr/local/RepeatMasker/RepeatMasker /mnt/volume1/maker/bil01gen.dipnr.pri_rnd1.m
7725 root 20 0 764M 728M 5192 R 74.0 0.1 0:05.67 /usr/bin/perl /usr/local/RepeatMasker/RepeatMasker /mnt/volume1/maker/bil01gen.dipnr.pri_rnd1.m
7584 root 20 0 767M 731M 5136 R 74.0 0.1 0:05.81 /usr/bin/perl /usr/local/RepeatMasker/RepeatMasker /mnt/volume1/maker/bil01gen.dipnr.pri_rnd1.m
7581 root 20 0 758M 722M 5144 R 74.0 0.1 0:05.34 /usr/bin/perl /usr/local/RepeatMasker/RepeatMasker /mnt/volume1/maker/bil01gen.dipnr.pri_rnd1.m
7594 root 20 0 761M 725M 5096 R 73.4 0.1 0:05.48 /usr/bin/perl /usr/local/RepeatMasker/RepeatMasker /mnt/volume1/maker/bil01gen.dipnr.pri_rnd1.m
7545 root 20 0 761M 725M 5176 R 73.4 0.1 0:05.44 /usr/bin/perl /usr/local/RepeatMasker/RepeatMasker /mnt/volume1/maker/bil01gen.dipnr.pri_rnd1.m
7731 root 20 0 757M 722M 5200 R 70.3 0.1 0:05.26 /usr/bin/perl /usr/local/RepeatMasker/RepeatMasker /mnt/volume1/maker/bil01gen.dipnr.pri_rnd1.m
7829 root 20 0 738M 703M 5124 R 69.7 0.1 0:04.43 /usr/bin/perl /usr/local/RepeatMasker/RepeatMasker /mnt/volume1/maker/bil01gen.dipnr.pri_rnd1.m
7815 root 20 0 736M 701M 5272 R 69.7 0.1 0:04.20 /usr/bin/perl /usr/local/RepeatMasker/RepeatMasker /mnt/volume1/maker/bil01gen.dipnr.pri_rnd1.m
8067 root 20 0 734M 699M 5128 R 69.1 0.1 0:04.14 /usr/bin/perl /usr/local/RepeatMasker/RepeatMasker /mnt/volume1/maker/bil01gen.dipnr.pri_rnd1.m
7873 root 20 0 757M 722M 5084 R 68.5 0.1 0:05.25 /usr/bin/perl /usr/local/RepeatMasker/RepeatMasker /mnt/volume1/maker/bil01gen.dipnr.pri_rnd1.m
7934 root 20 0 746M 710M 5256 R 68.5 0.1 0:04.67 /usr/bin/perl /usr/local/RepeatMasker/RepeatMasker /mnt/volume1/maker/bil01gen.dipnr.pri_rnd1.m
7985 root 20 0 736M 701M 5036 R 67.2 0.1 0:04.21 /usr/bin/perl /usr/local/RepeatMasker/RepeatMasker /mnt/volume1/maker/bil01gen.dipnr.pri_rnd1.m
7656 root 20 0 754M 718M 5148 R 66.6 0.1 0:04.97 /usr/bin/perl /usr/local/RepeatMasker/RepeatMasker /mnt/volume1/maker/bil01gen.dipnr.pri_rnd1.m
F1Help F2Setup F3Search F4Filter F5Tree F6SortBy F7Nice F8Nice +F9Kill F10Quit
```

# MONITOR AND OPTIMISE YOUR PIPELINES

## NETDATA



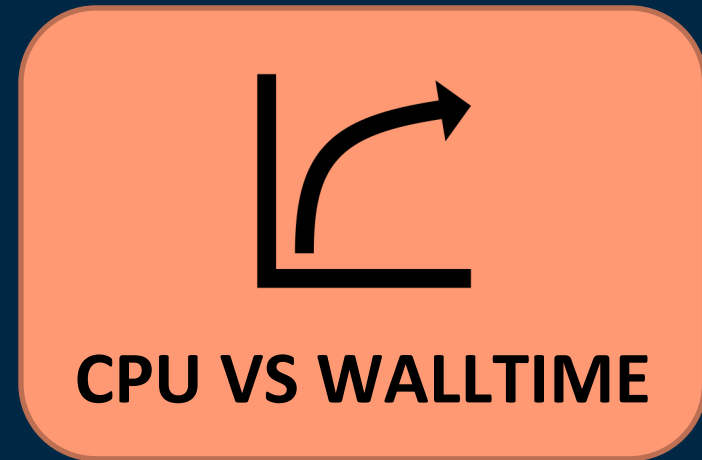
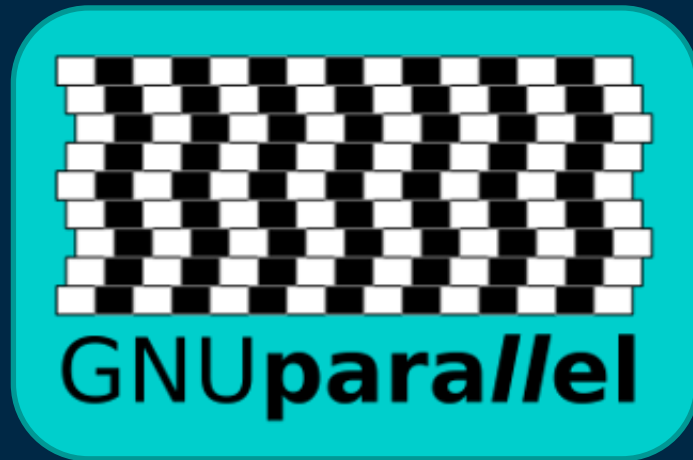
# MONITOR AND OPTIMISE YOUR PIPELINES

## SCHEDULER LOGS

```
Job Id: 2566026.pbserver in queue normal
Job Name: variant_only_vcfs
Project: GENOMICS
Exit Status: 0
Job run as chunks (hpc140:ncpus=12:mem=100663296kb)
Walltime requested: 12:00:00 :
                        :
-- Nodes Summary --
-- node hpc140 summary
   Cpus requested:      12 :
   Cpu Time:          03:29:21 :
   Mem requested:     96.0GB :
                        :
```

# MONITOR AND OPTIMISE YOUR PIPELINES

## OPTIMISATION



CPU VS WALLTIME





# FILE MANIPULATION

08

# GET FAMILIAR WITH BASIC BASH COMMANDS



grep

Pattern  
Matching

```
grep "chromosome 5"
```



awk

Data  
Processing

```
awk '$1 == 5 {print $2, $3}'
```



sed

Find &  
Replace

```
sed 's/sample1/ID7037/g'
```



# RECORD KEEPING

09

WRITE IT DOWN!



**labarchives**  
Research Notebook



**Evernote**



**git**



VISUAL  
STUDIO CODE



ATOM

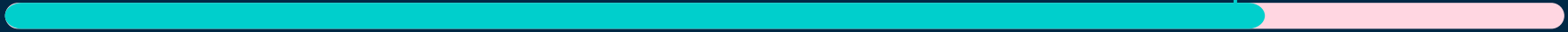


MARKDOWN



PATIENCE

10



PATIENCE IS KEY!



# THANK YOU! 😊



✉ [parice.brandies@sydney.edu.au](mailto:parice.brandies@sydney.edu.au)

🐦 [@PariceBrandies](https://twitter.com/PariceBrandies)