



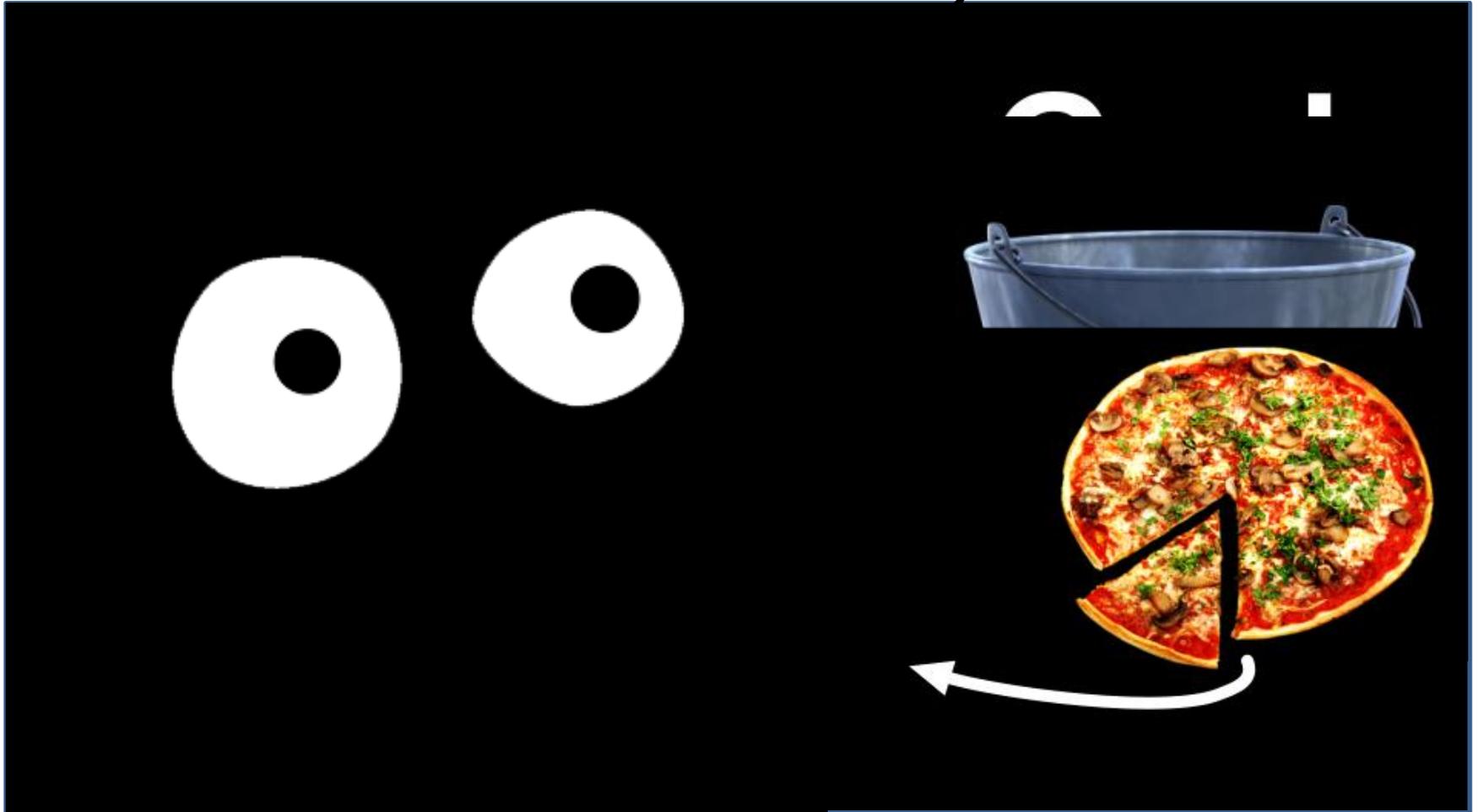
How to build a supercomputer

ON THE CHEAP

G Burton – ICG – Jan-14 – v1.2



Server Guy



HPC Facts

- Chinese Tianhe-2 (MilkyWay -2) - China's National University of Defense Technology
Number 1 in top 500 with 3.12 million Cores(processors).
- 34 petaFLOPS (Floating Point Operations / Second). Sciamia 10 teraFLOPS.
- In Top 500 - USA=213, China=37, UK=29
- 75% of Top 500 use Intel processors.
- Race for first ExaFlop (exescale) .

Computer performance

Name	FLOPS
yottaFLOPS	10^{24}
zettaFLOPS	10^{21}
exaFLOPS	10^{18}
petaFLOPS	10^{15}
teraFLOPS	10^{12}
gigaFLOPS	10^9
megaFLOPS	10^6
kiloFLOPS	10^3

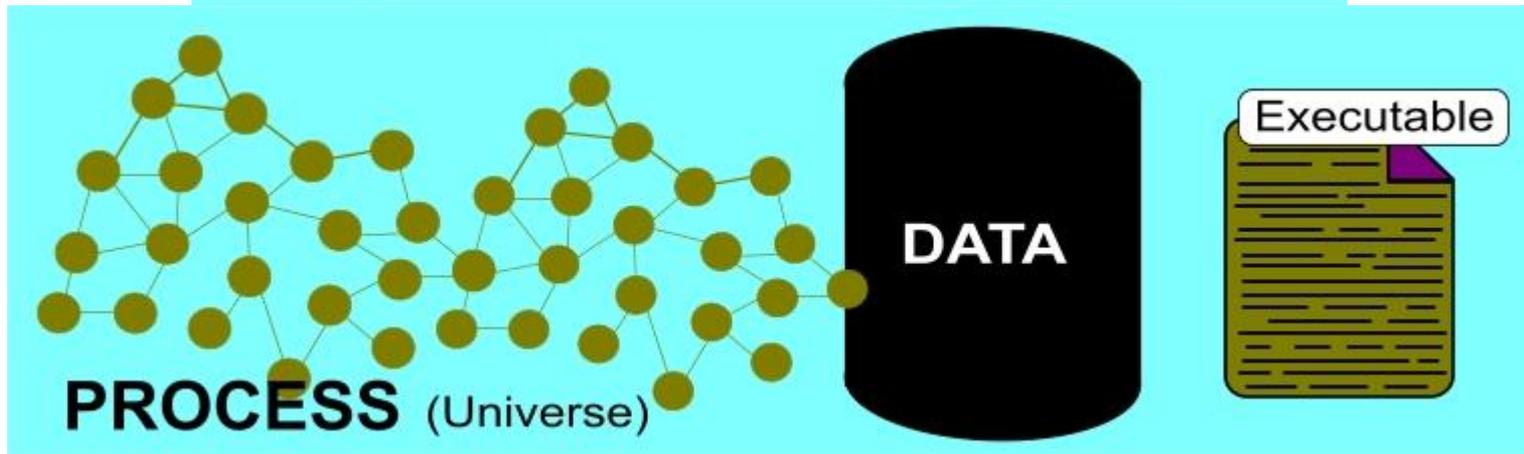
To illustrate how we use an HPC we will set ourselves a problem

What is the answer to the ultimate question of Life, the Universe , and Everything ?

C: Hitch Hikers Guide to the Galaxy

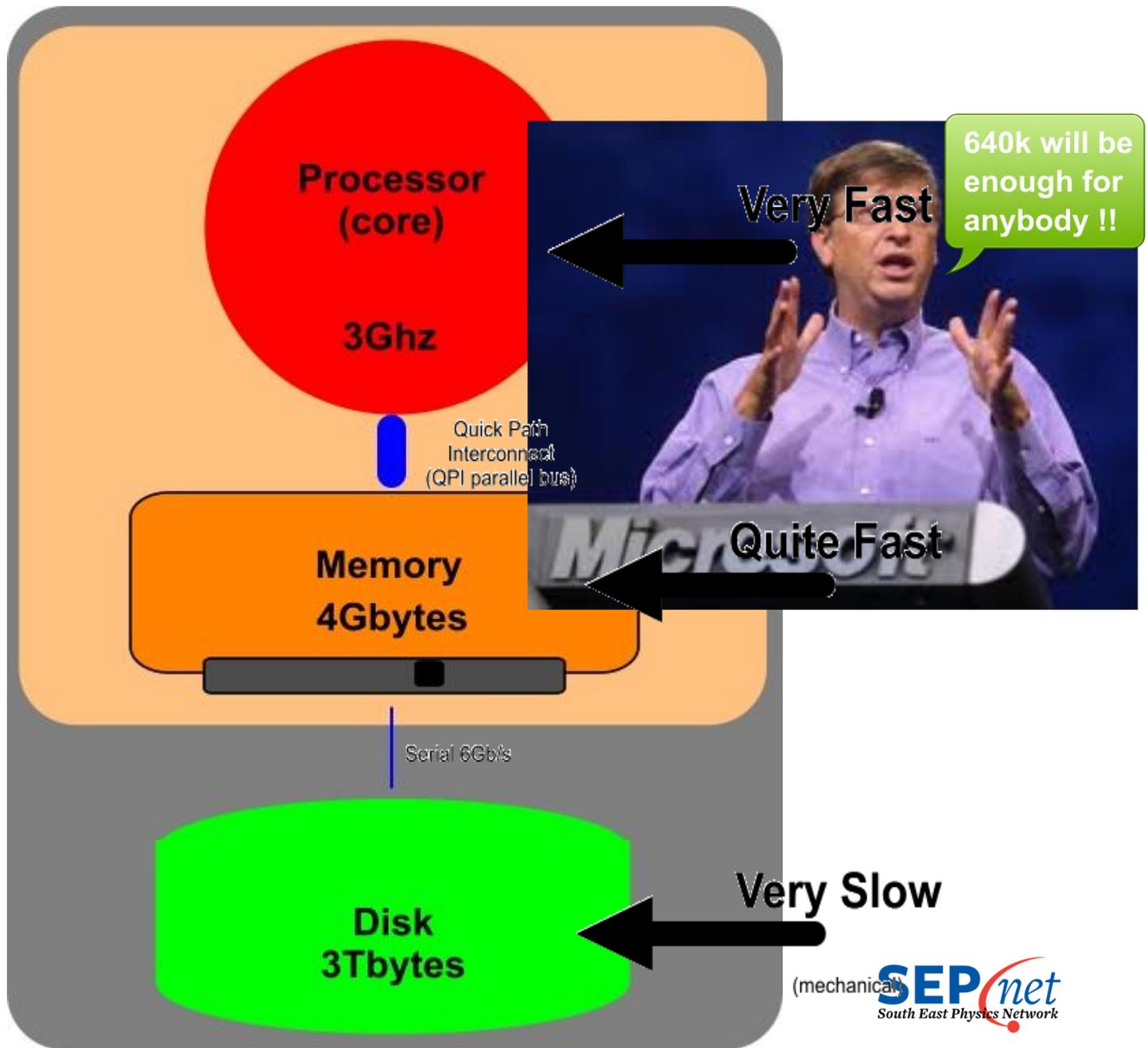


Start with small data, test with larger and then run with full....

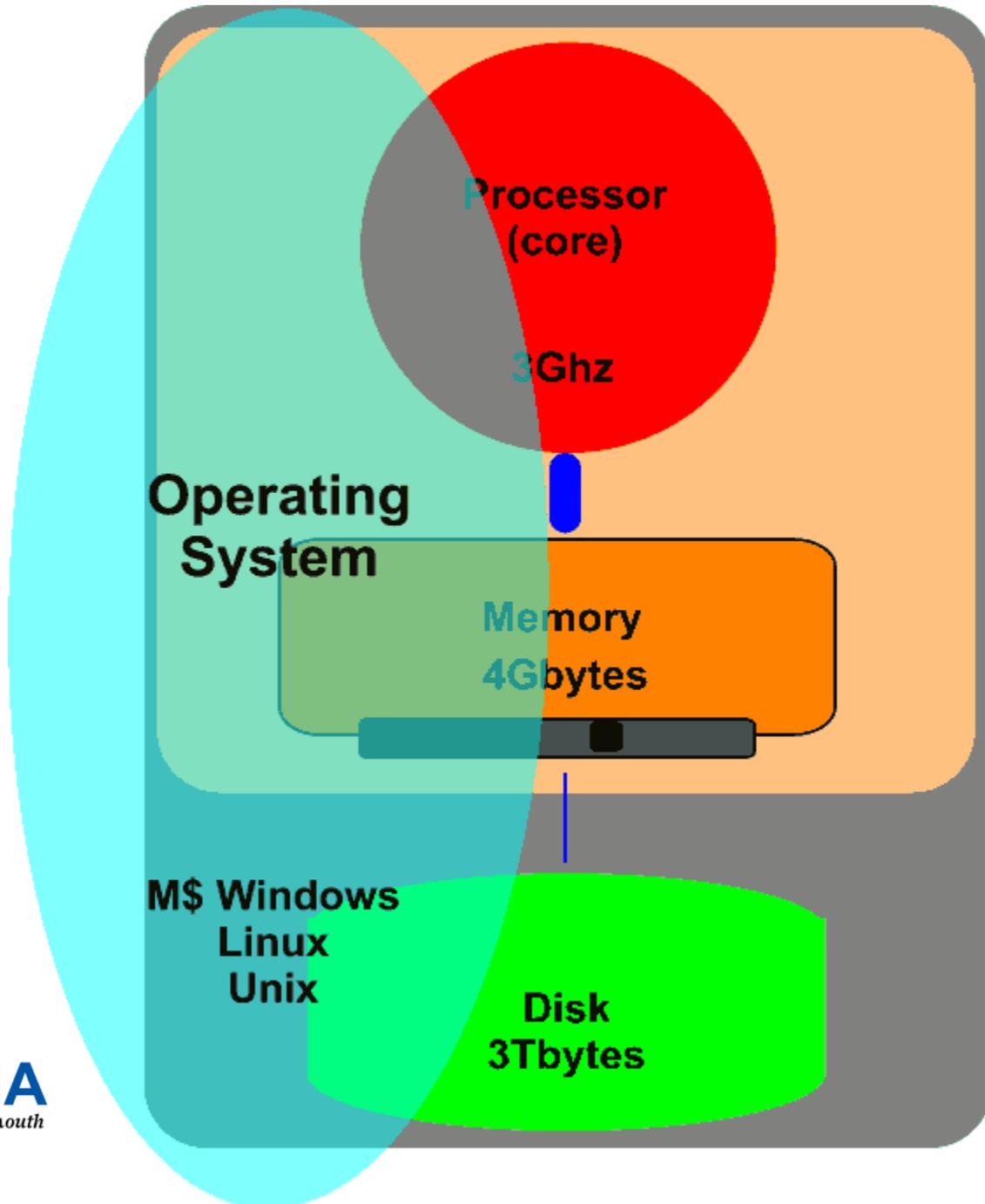


What is the answer to the ultimate question of Life, the ~~Galaxy system~~ , and Everything ?

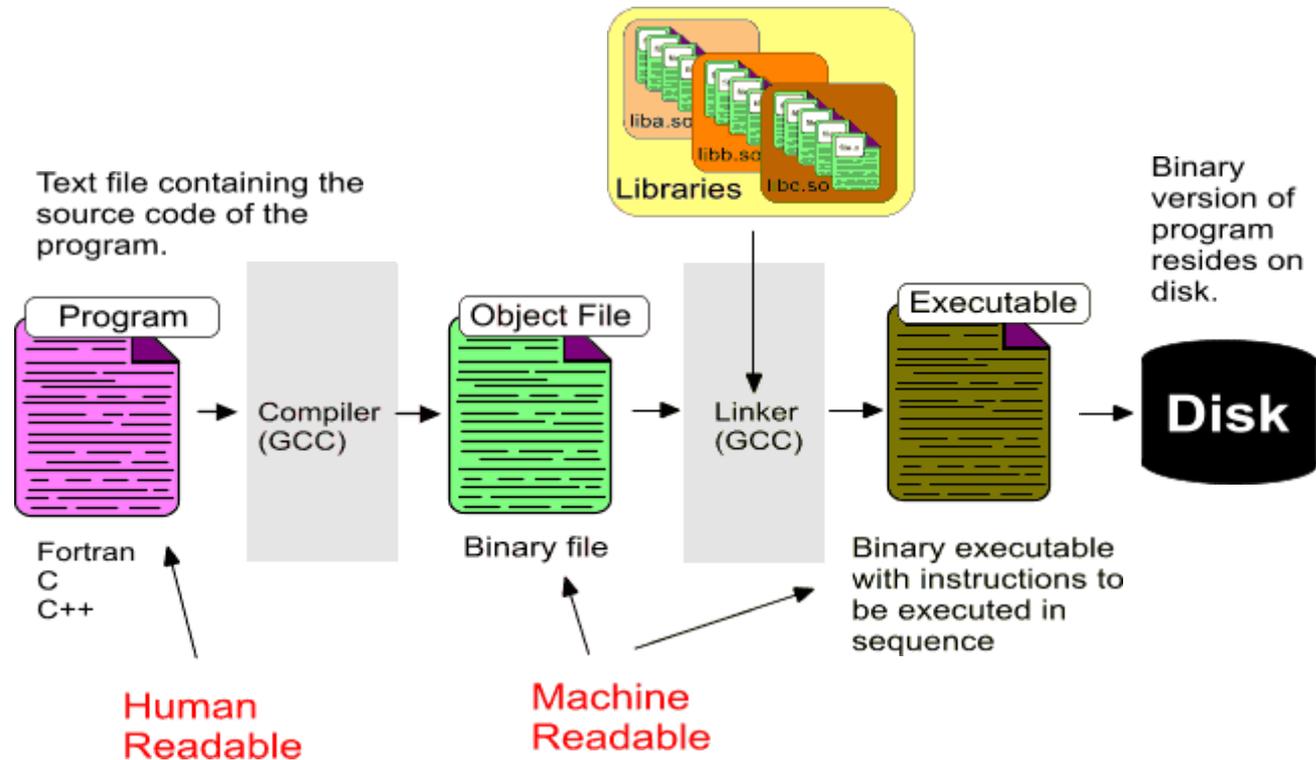
Back to Basics

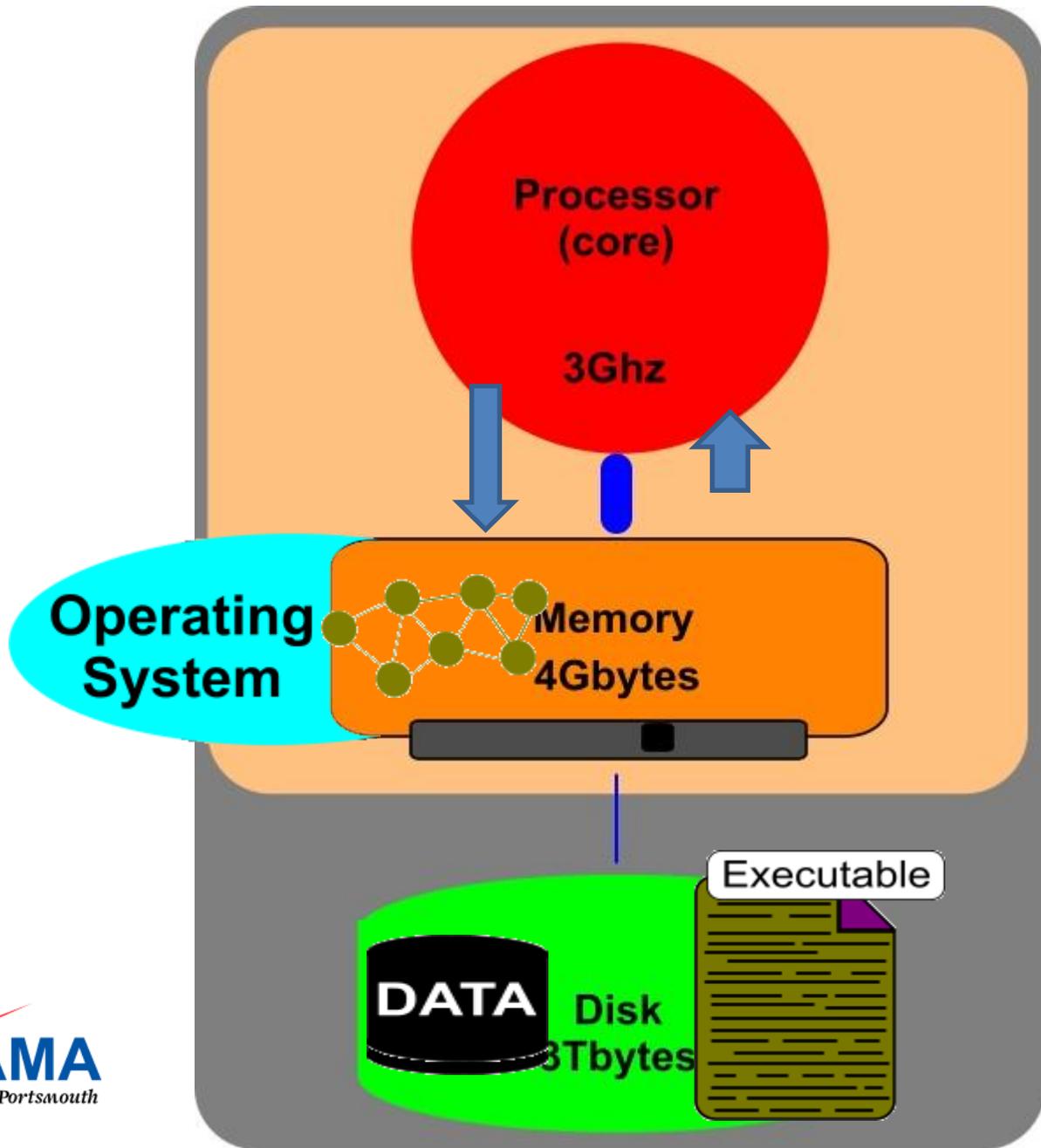


Back to Basics

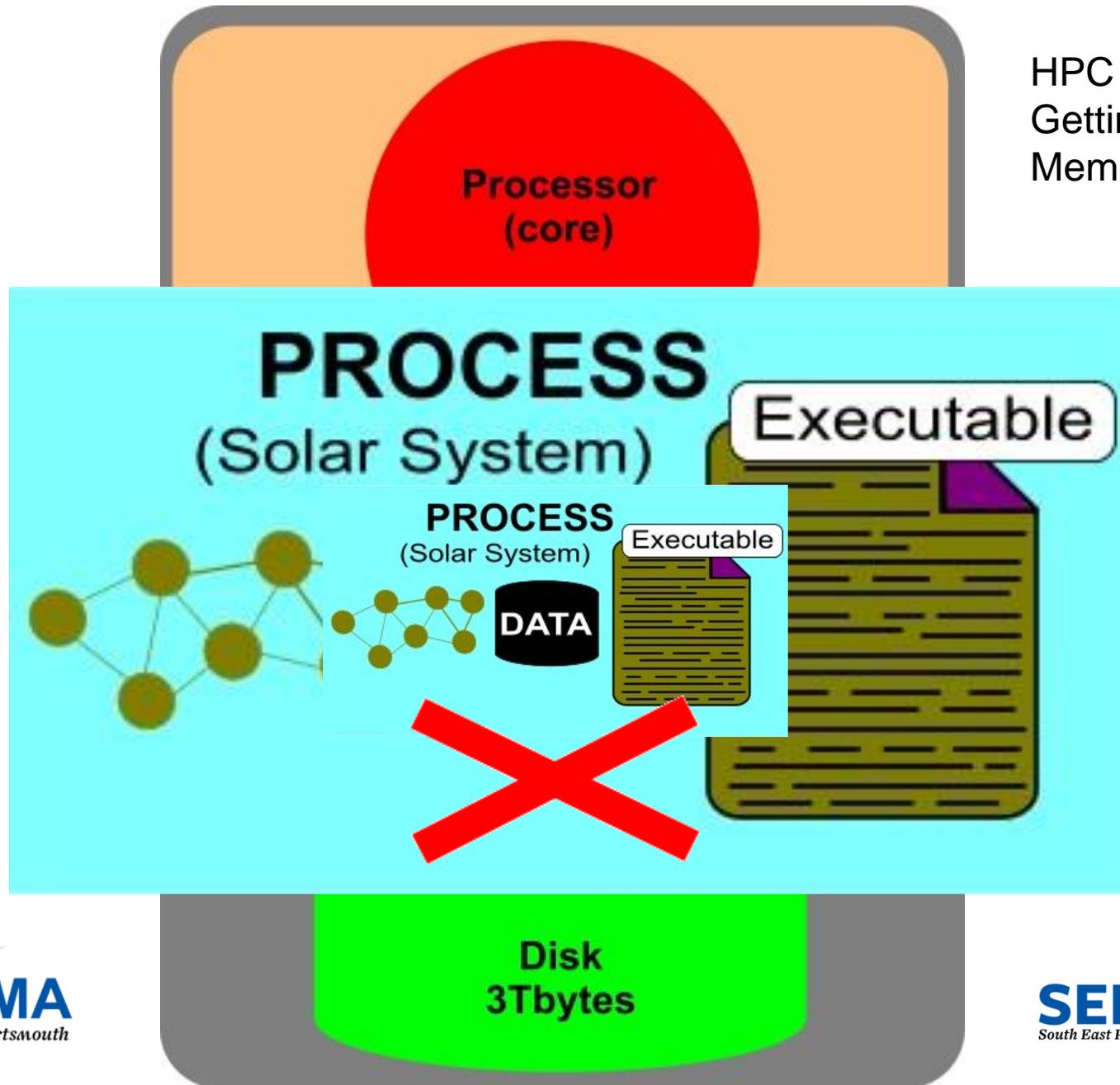


In order to solve our problem we need a “Program” to run.





HPC is all about
Getting access to
Memory.

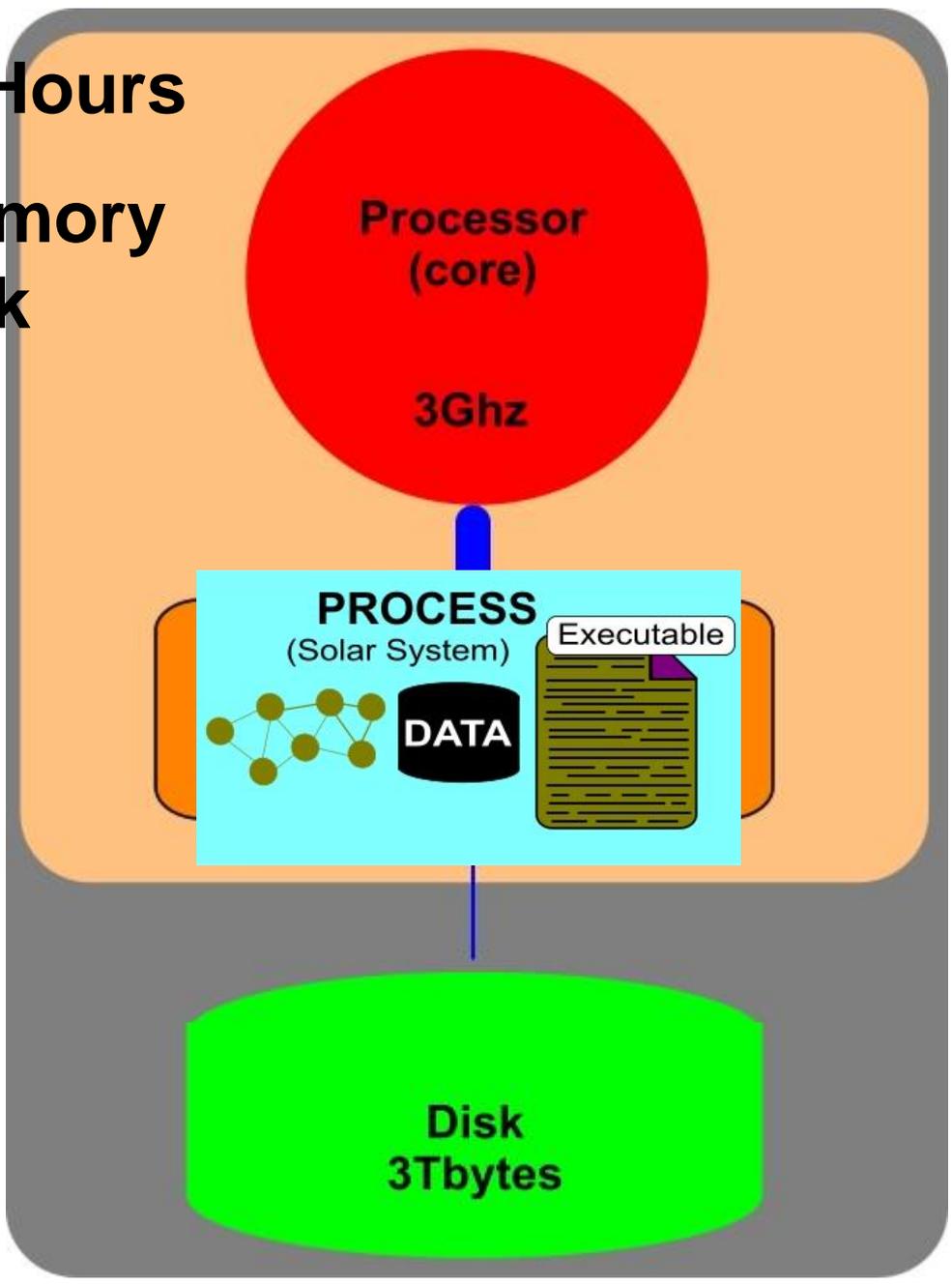


10,000 CPU Hours

(Check Point Restarts)

3 Gbytes Memory

2 Tbytes Disk



Answer

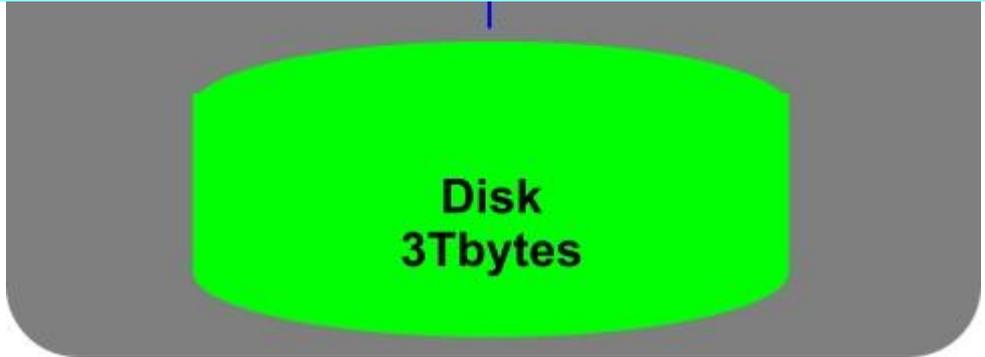
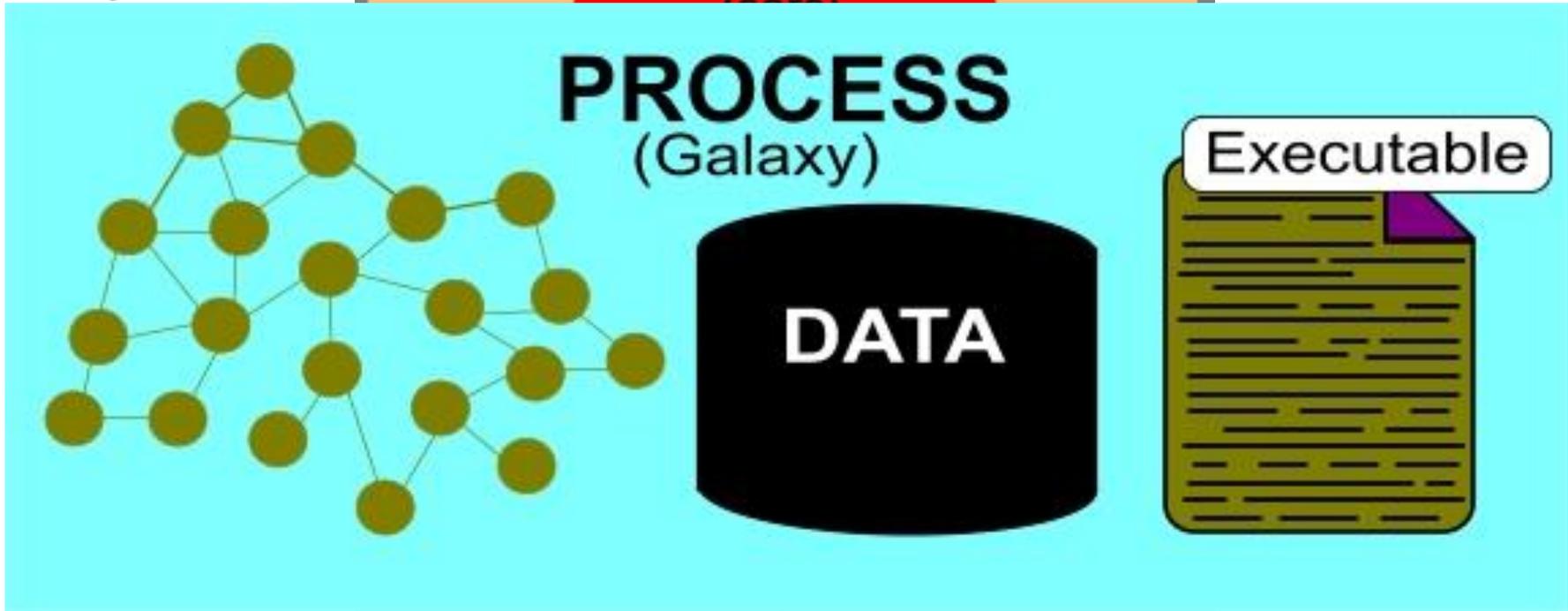
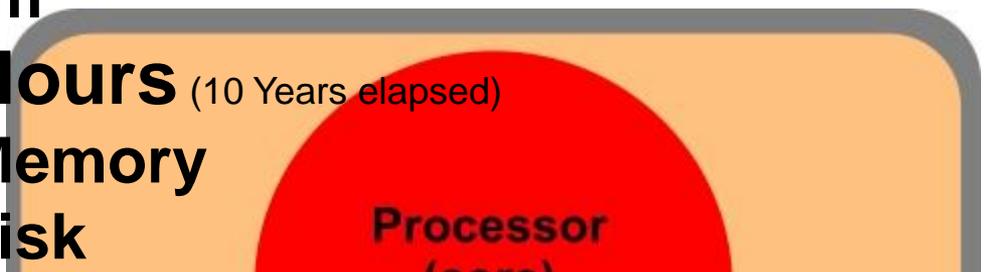
63

Galaxy Problem

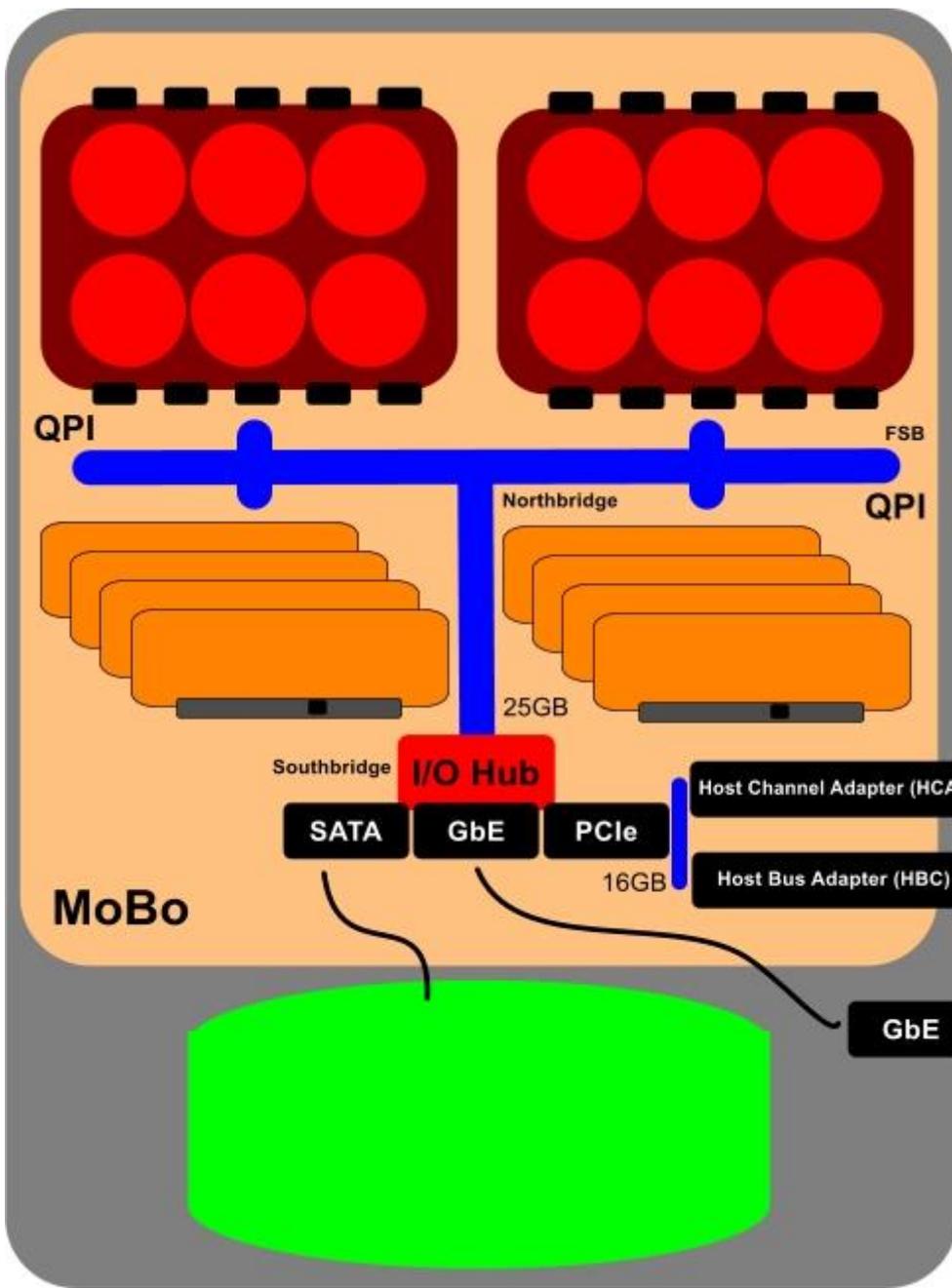
100,000 CPU Hours (10 Years elapsed)

30 Gbytes of Memory

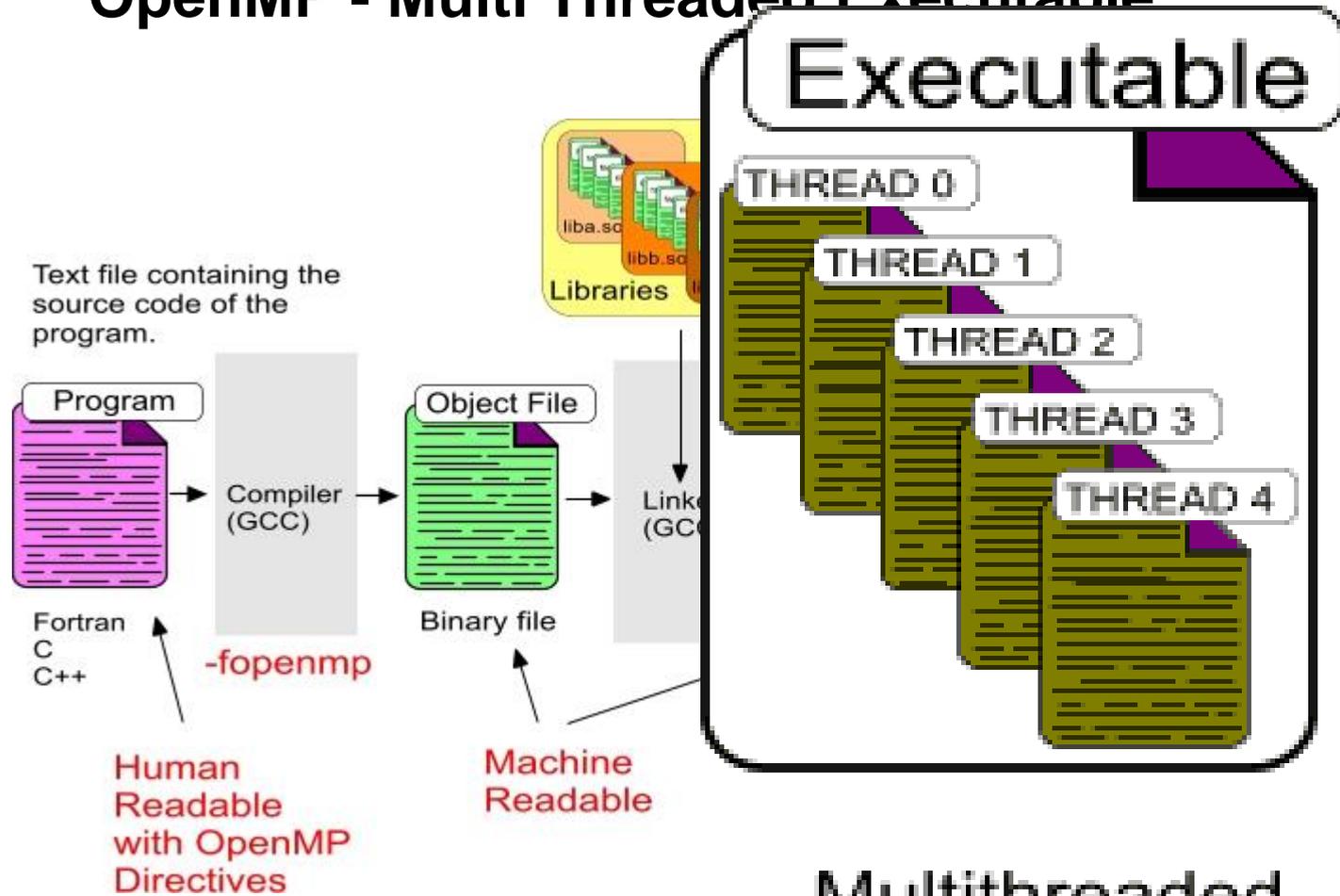
20 Tbytes of Disk



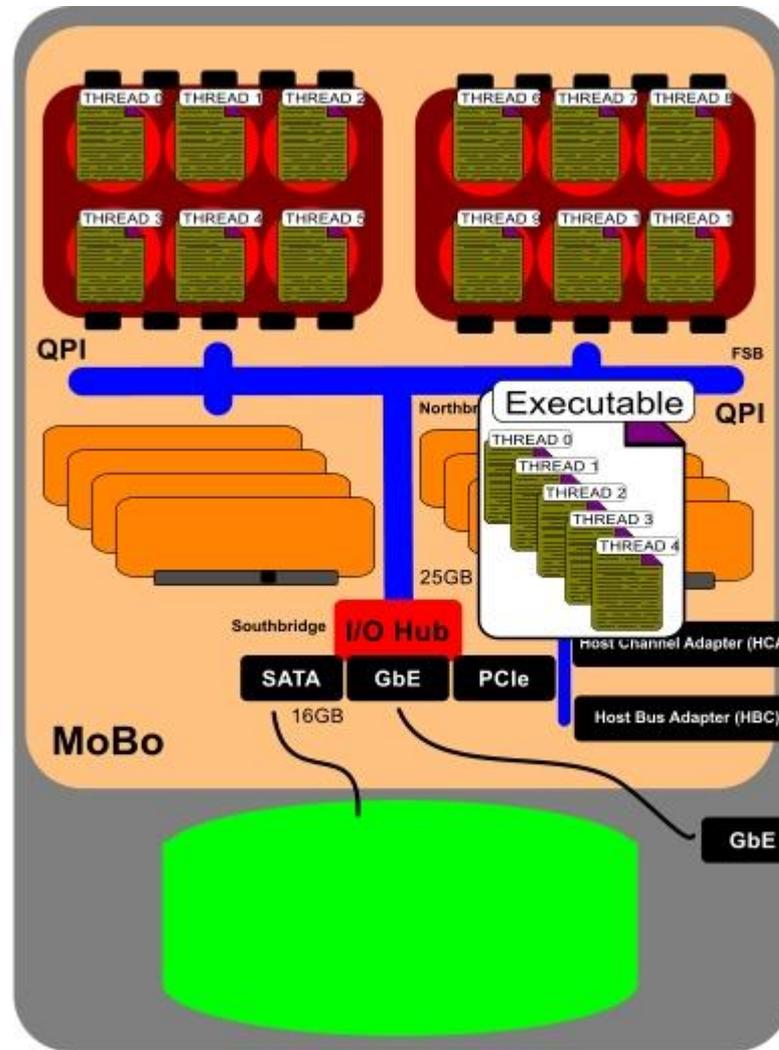
Galaxy Problem
 100,000 CPU Hours
 30 Gbytes of Memory
 20 Tbytes of Disk



OpenMP - Multi Threaded Executable



10 x Faster so
100,000 Cpu Hrs
goes to
10,000 Elapsed Hrs

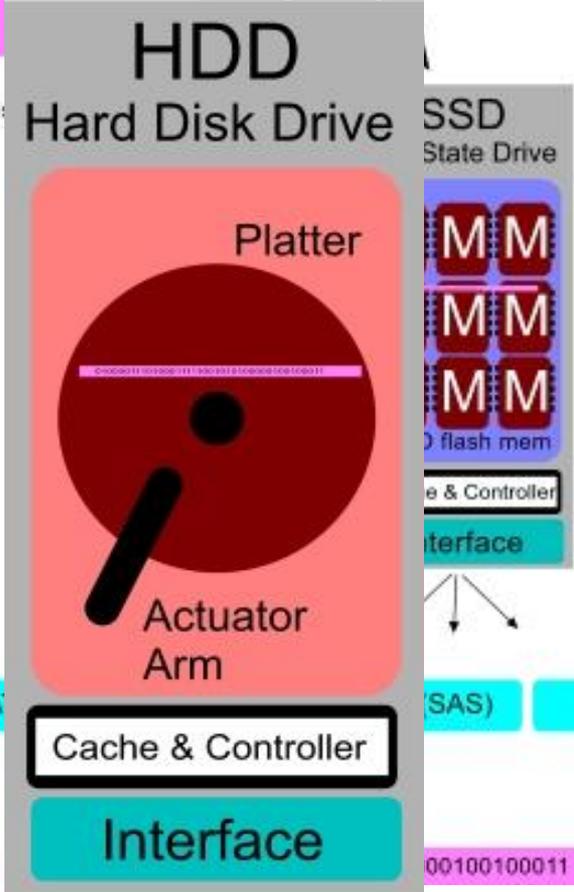


1TB HDD SATA 3Gb/s, 7200rpm, 32MB Cache £43
 1TB HDD SAS 3Gb/s, 7200rpm, 16MB Cache £99
 1TB HDD FATA 7200rpm £480
 480GB SSD SATA 3Gb/s £565

new: Self-encrypting drive:



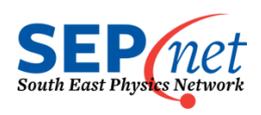
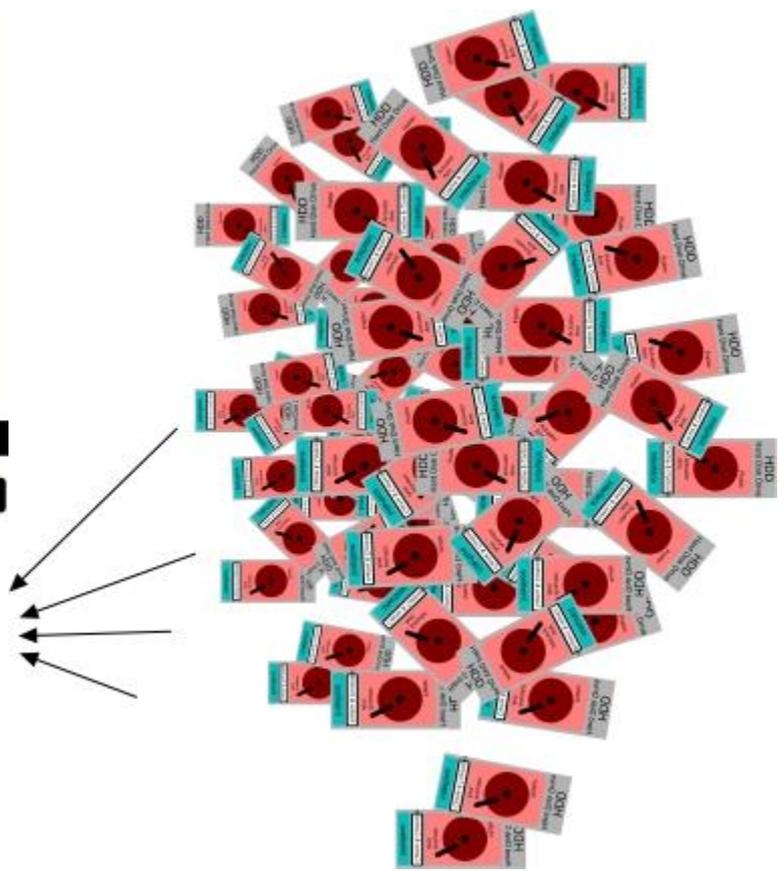
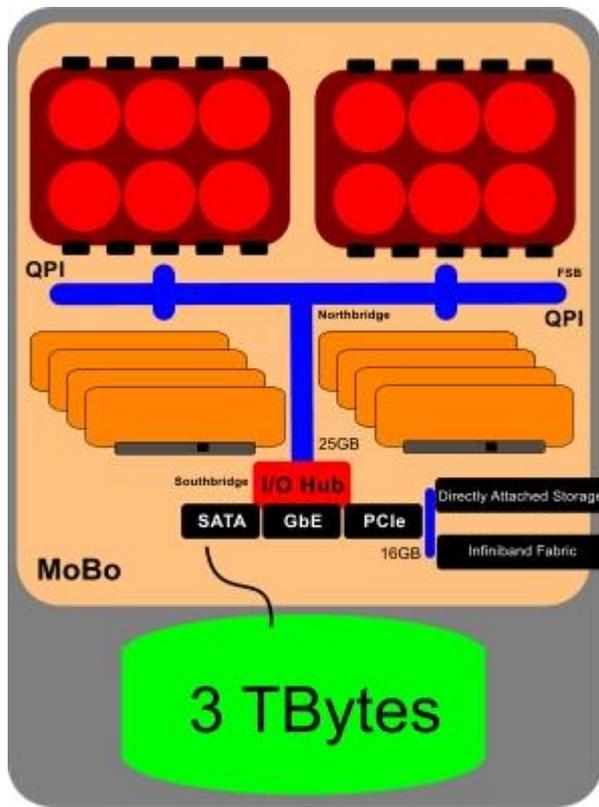
Disks store blocks of data.

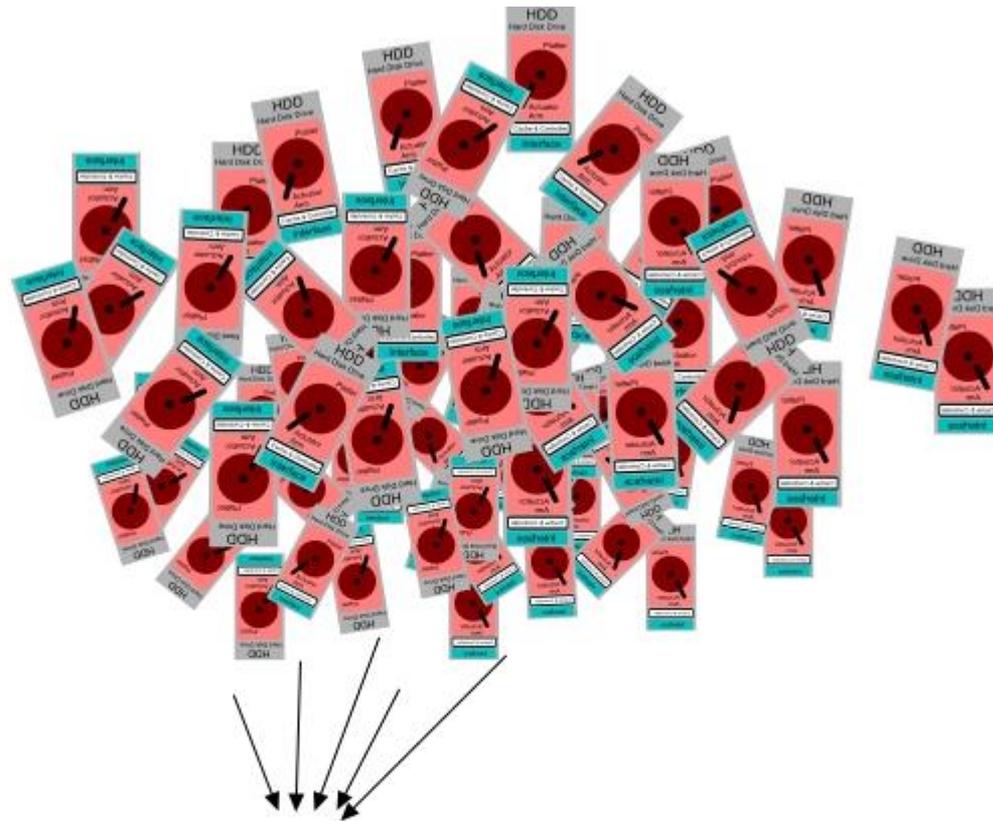


Disks are Dumb

Very Slow
 Quite bulky
 2T-3TBytes
 Cheap





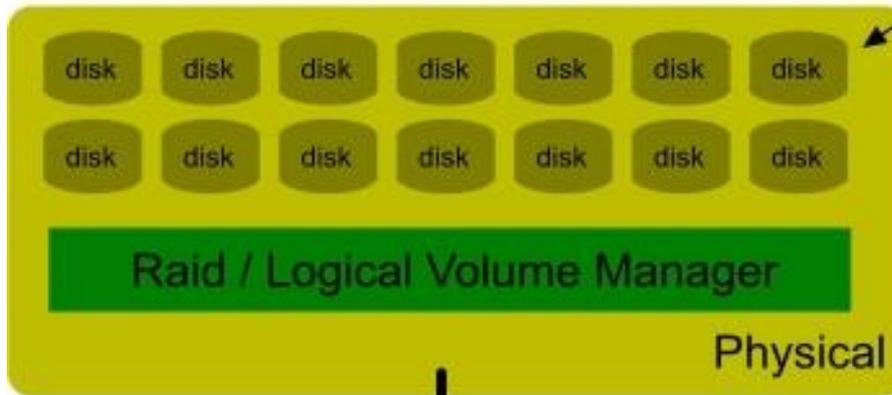


Commodity Disks

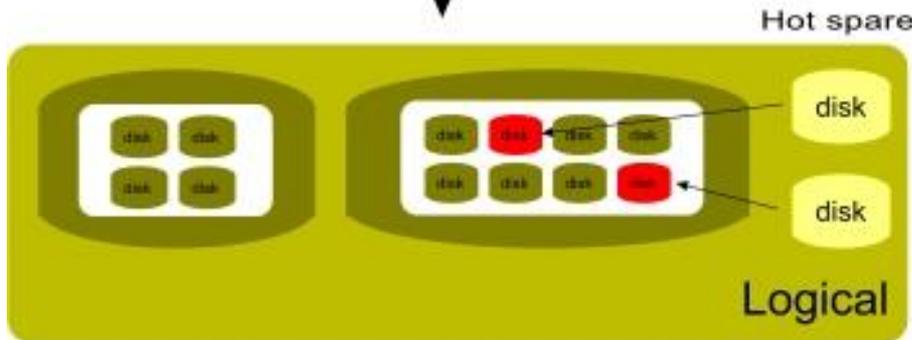


In HPC internal disks usually only contain the Operating System.
Some times two disks are "mirrored" for security.

Single Disk of little use for data :-

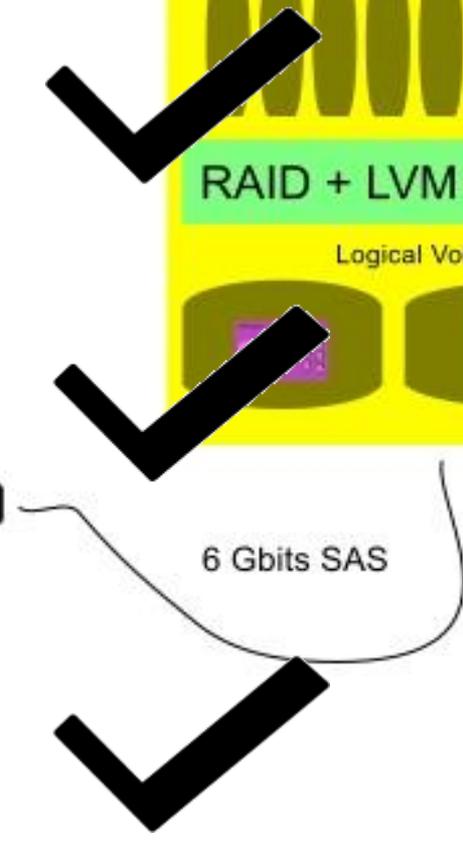
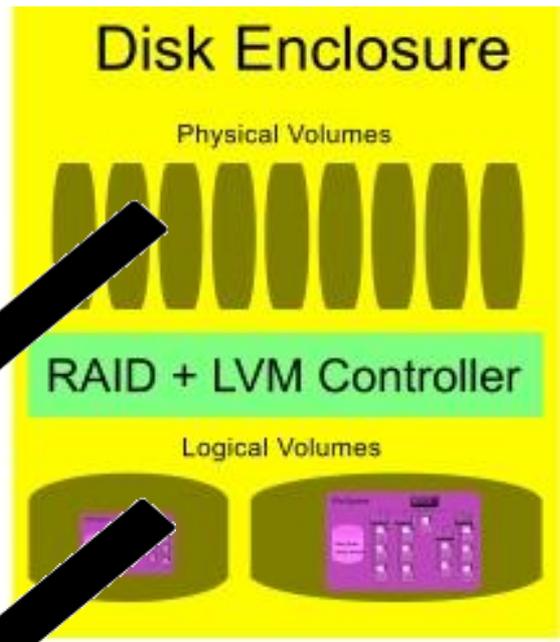
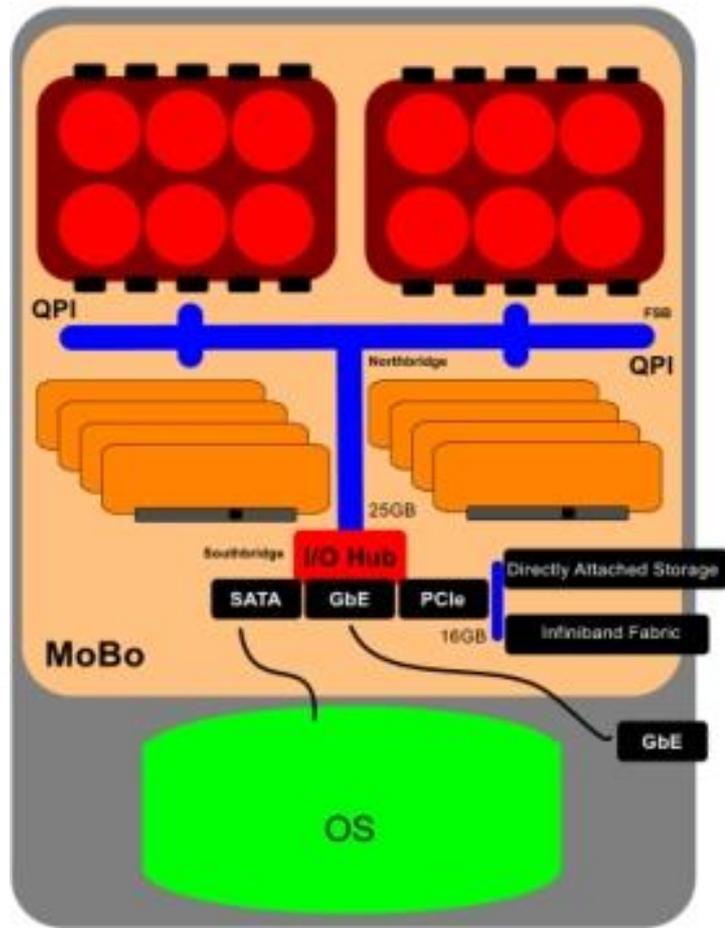


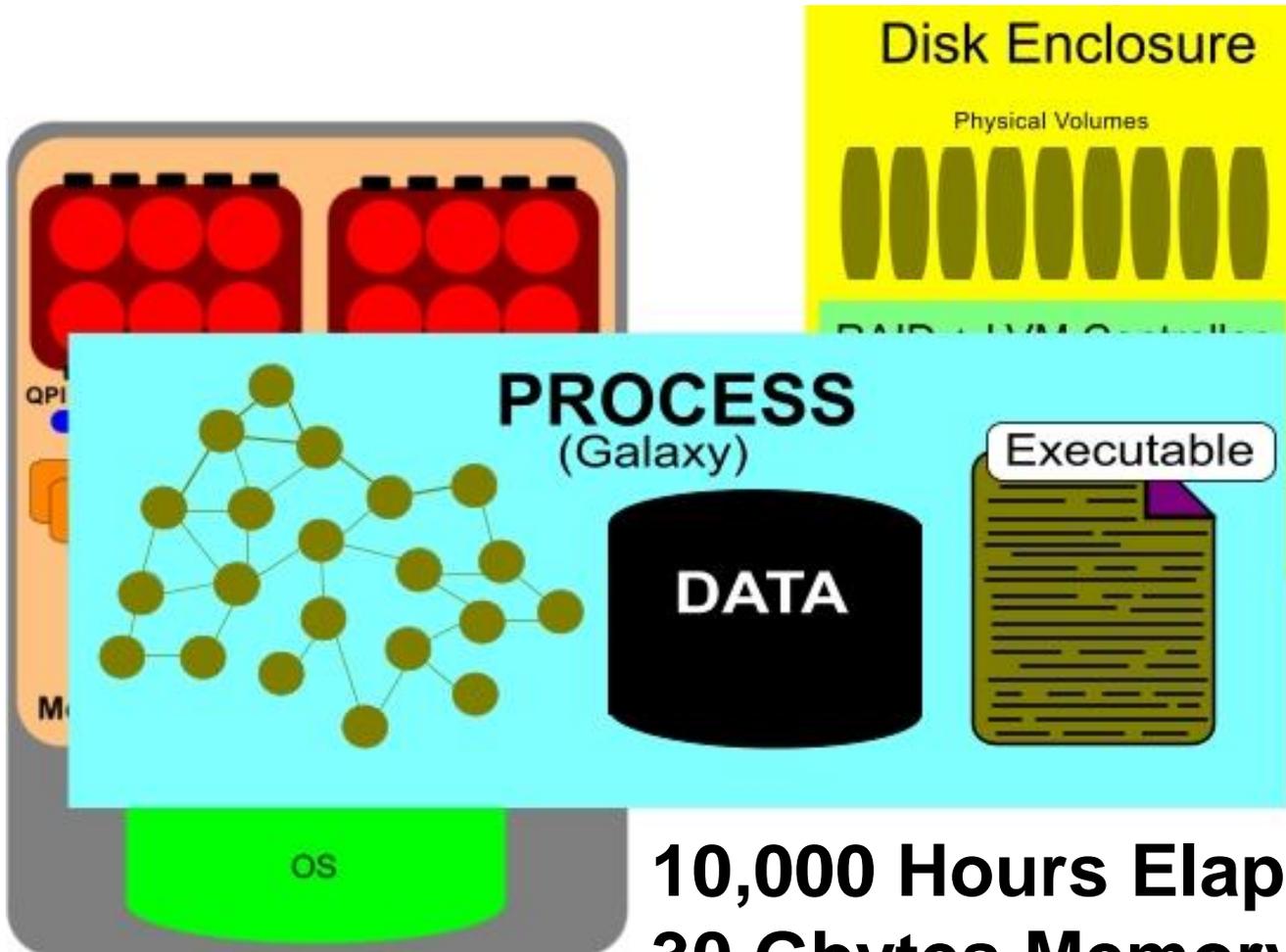
- Limited size.
- Limited performance
- Limited fault tolerance



Logical Volumes

- Increased size
- increased performance through striping.
- Increased resilience through parity and mirroring.



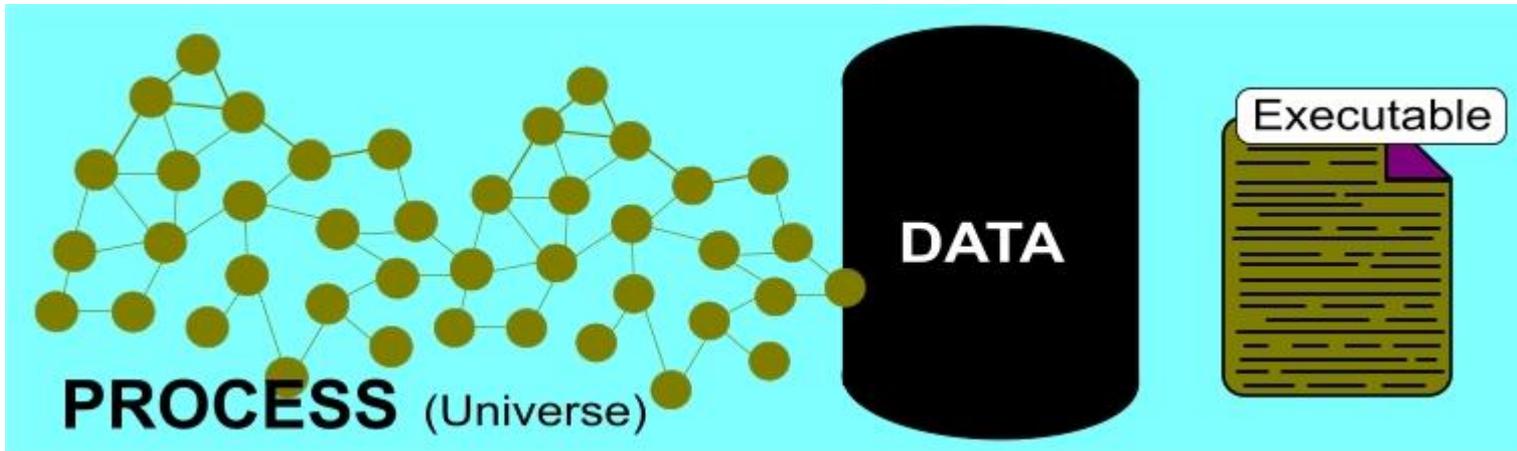


Answer
54

10,000 Hours Elapsed (Parallel Processing)

30 Gbytes Memory (Large memory boards)

20 Tbytes Disk (Directly Attached Storage)



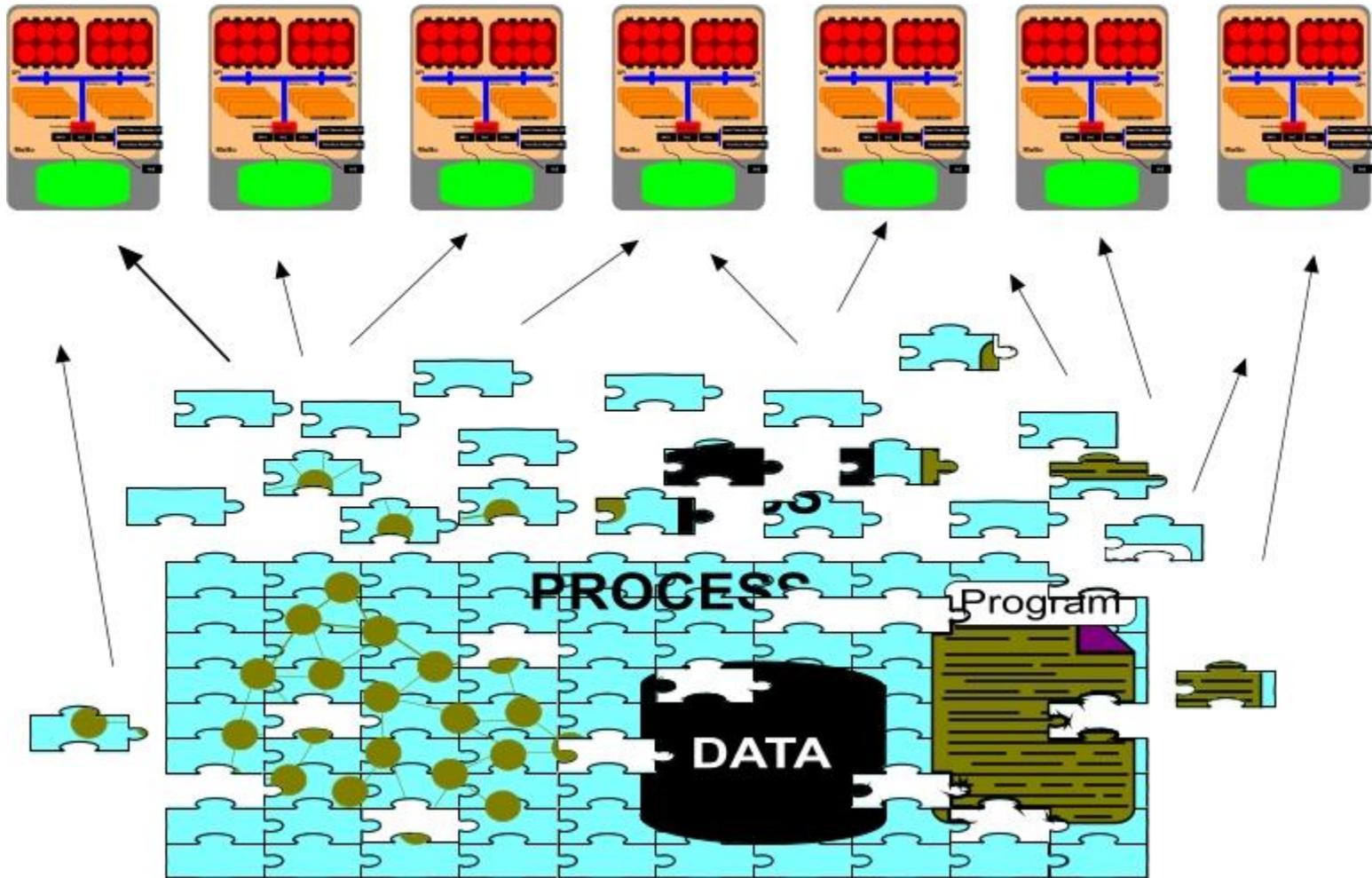
Universe Problem

10,000,000 CPU Hours (1000 years)

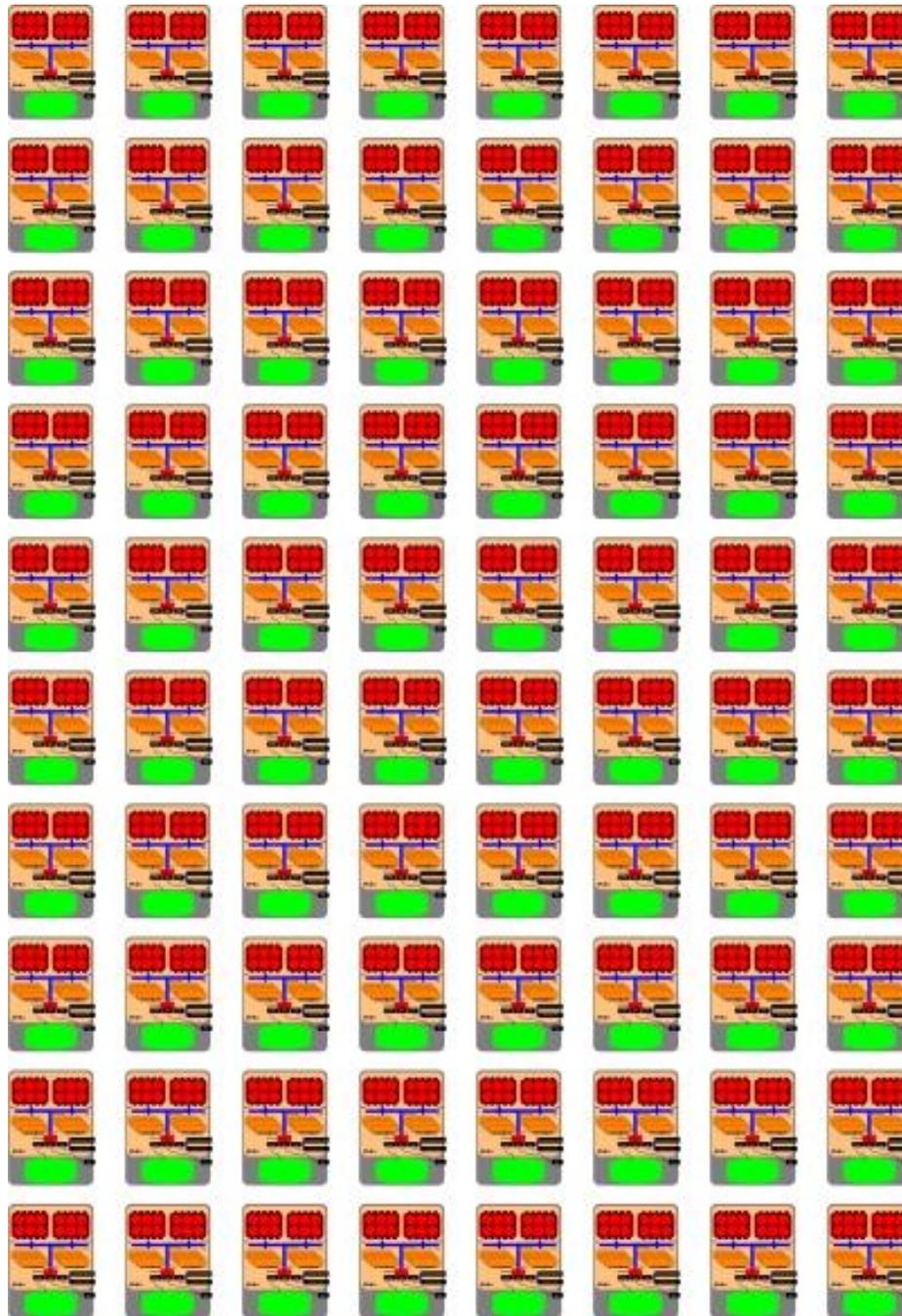
3T Bytes Memory

200 Tbytes Disk

Commodity is Cheap – Specialised is Expensive



Nodes (10 Cores)



Commodity Cluster

Required:-

10,000,000 CPU Hours

3T Bytes Memory

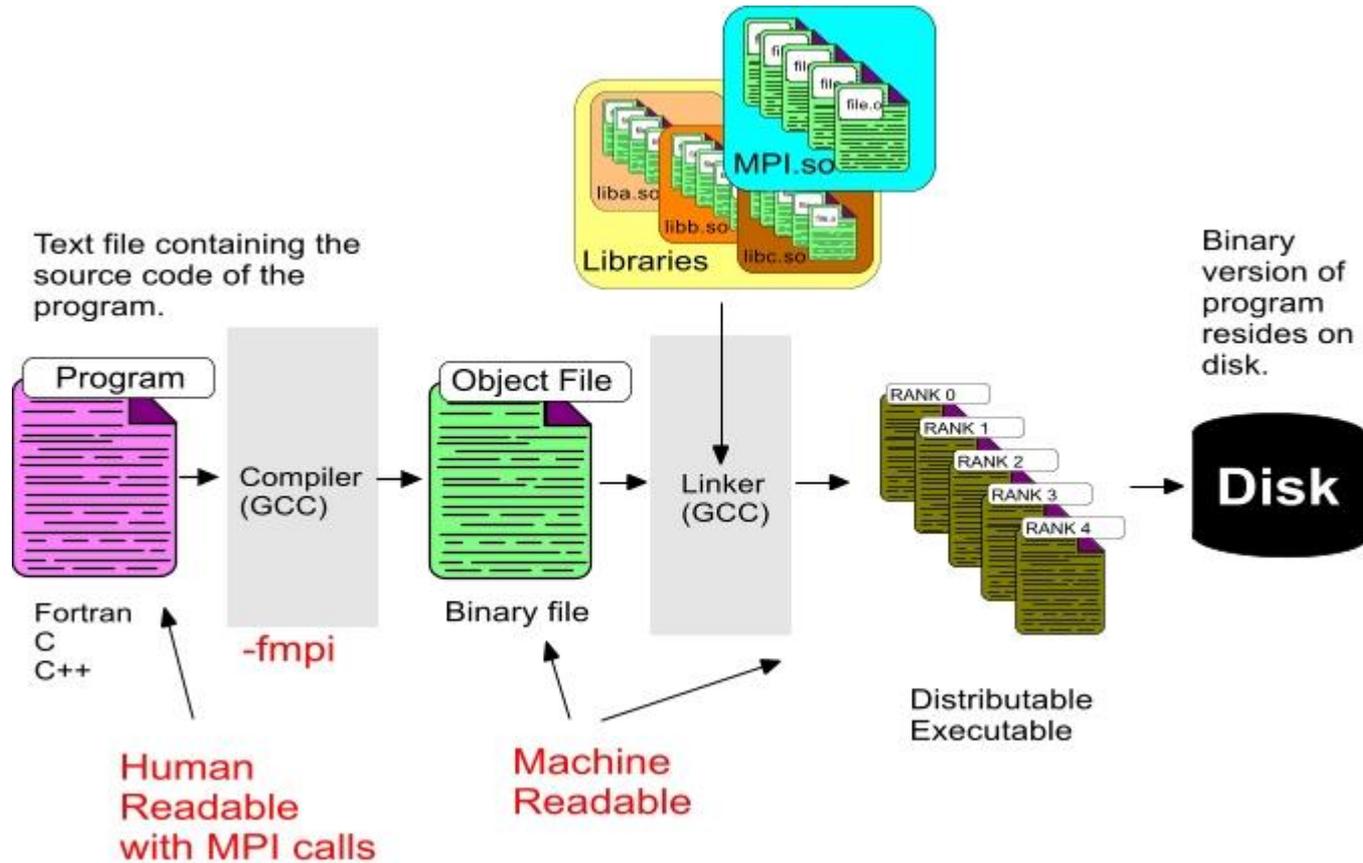
Distributed Over 100
Nodes (100x10 cores):-

10,000 Hours Elapsed

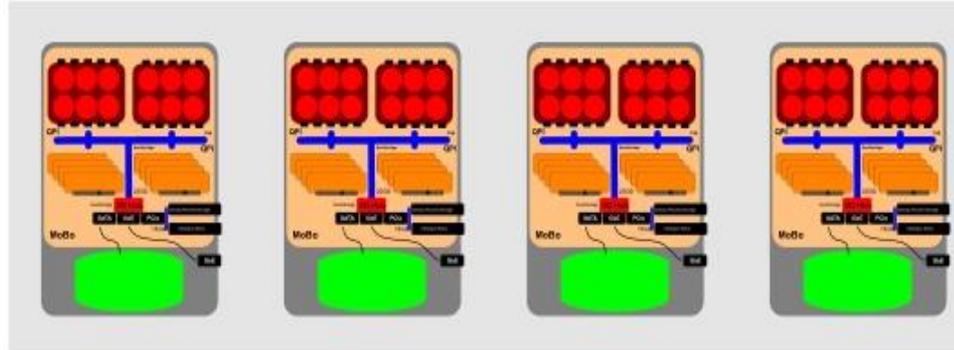
30Gbytes / Node



Message Passing Interface - MPI



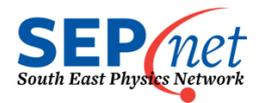
Dell 6100 Front View



Dell 6100 Rear View

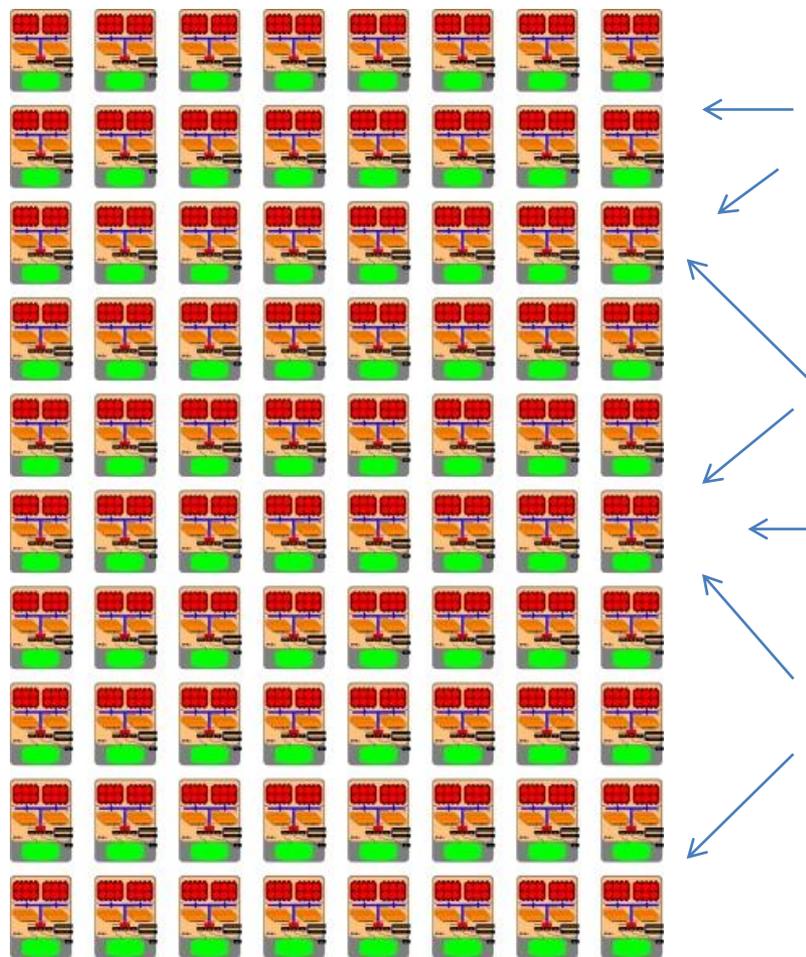
High density - Google / Amazon



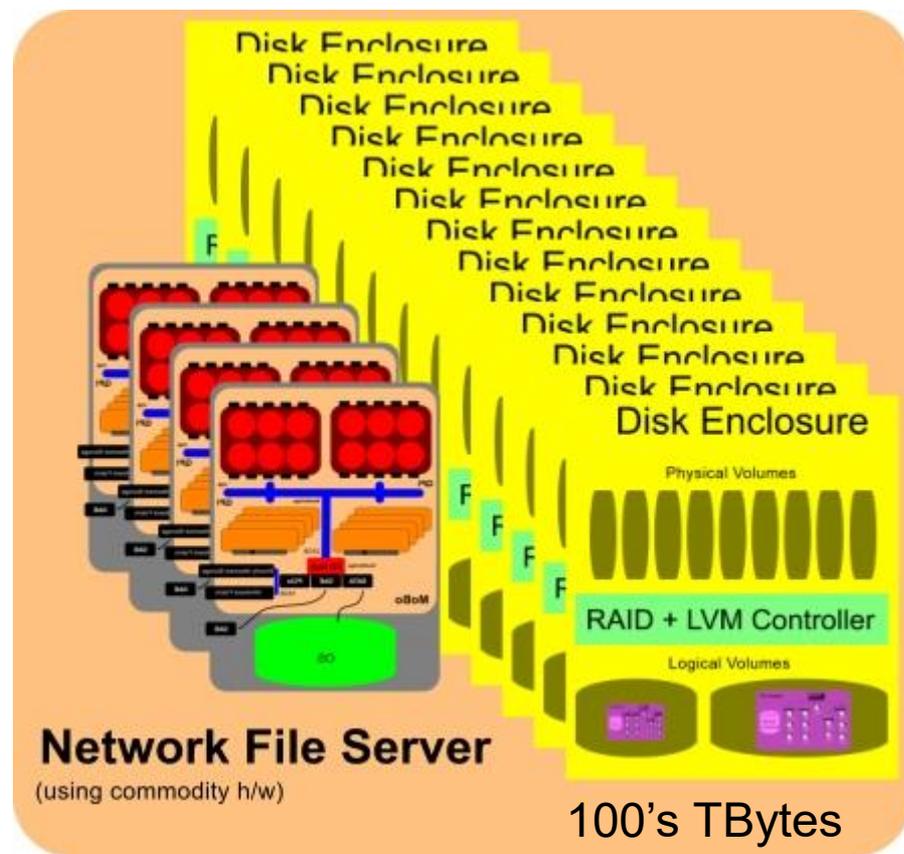


Shared Storage

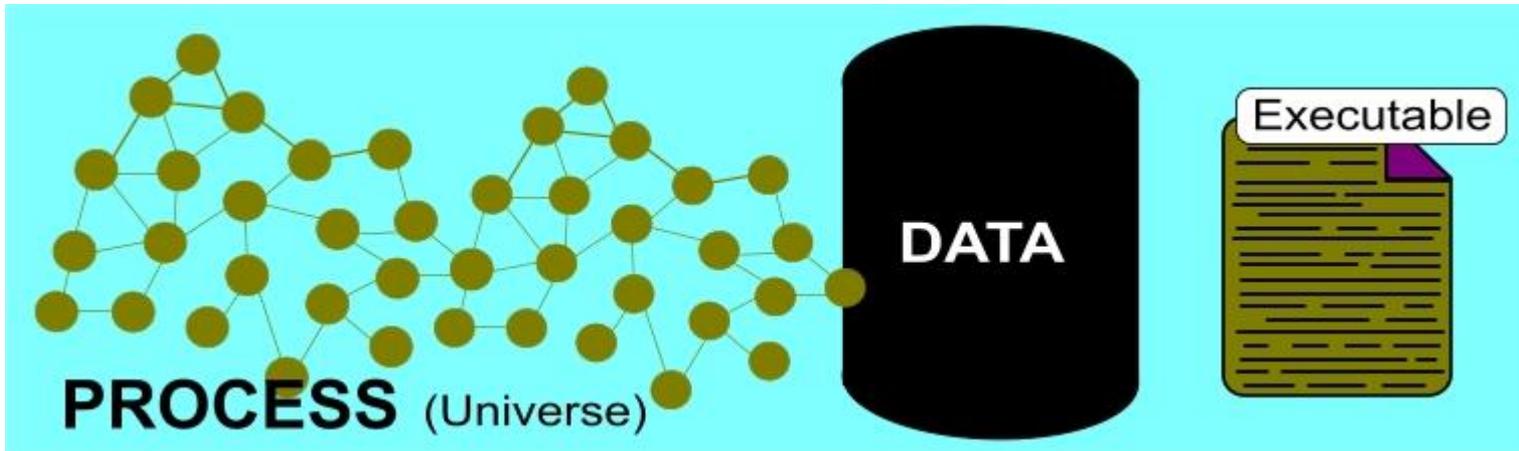
(Require 200TBytes)



Network File System Clients



A change made by one Node will be seen by all others.



Universe Problem

Answer

42

10,000 Hours Elapsed (Distributed Cores)

3 Tbytes Memory (Distributed Memory)

200 Tbytes Disk (Network File Server)

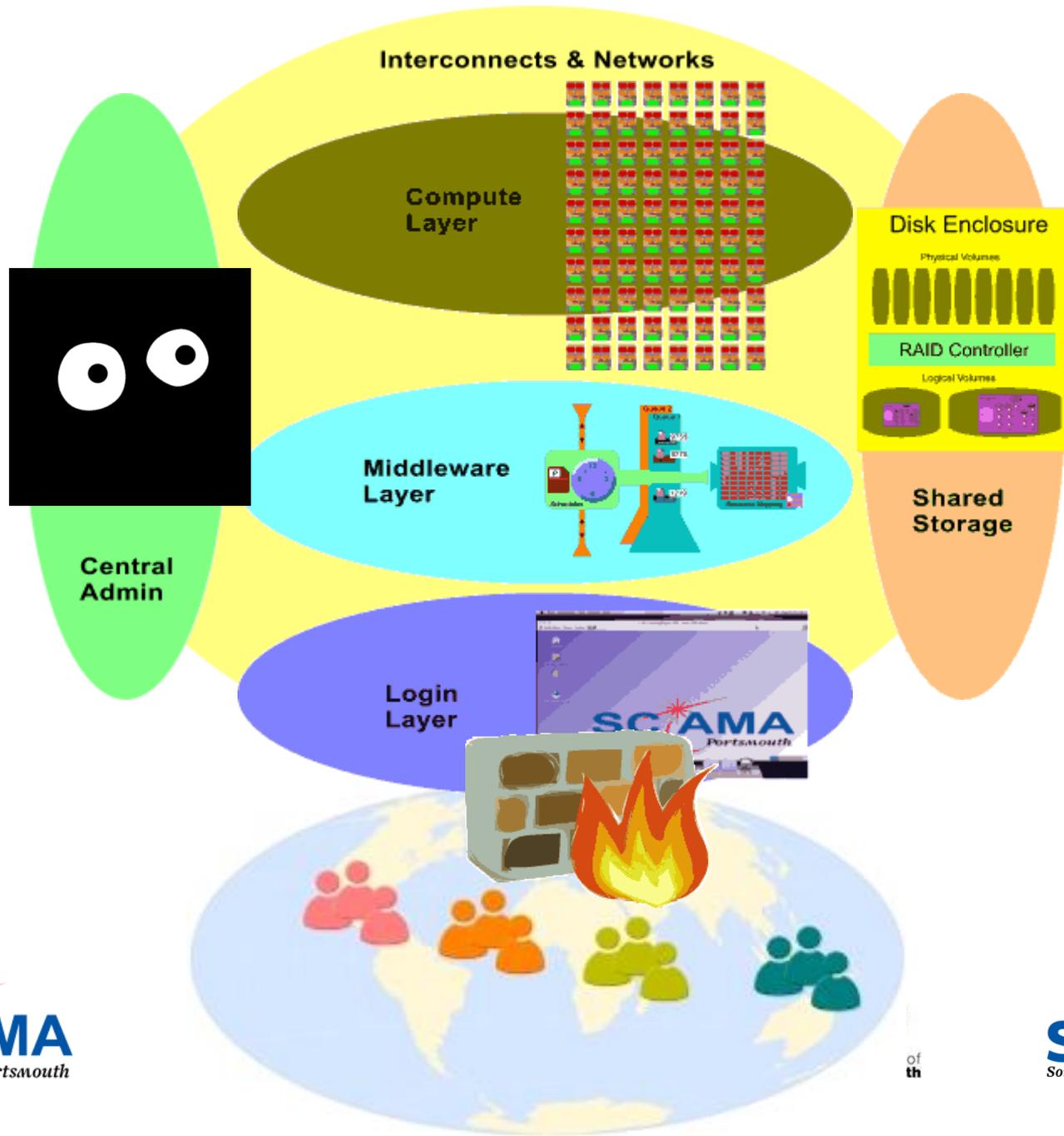
C: Hitch Hikers Guide to the Galaxy



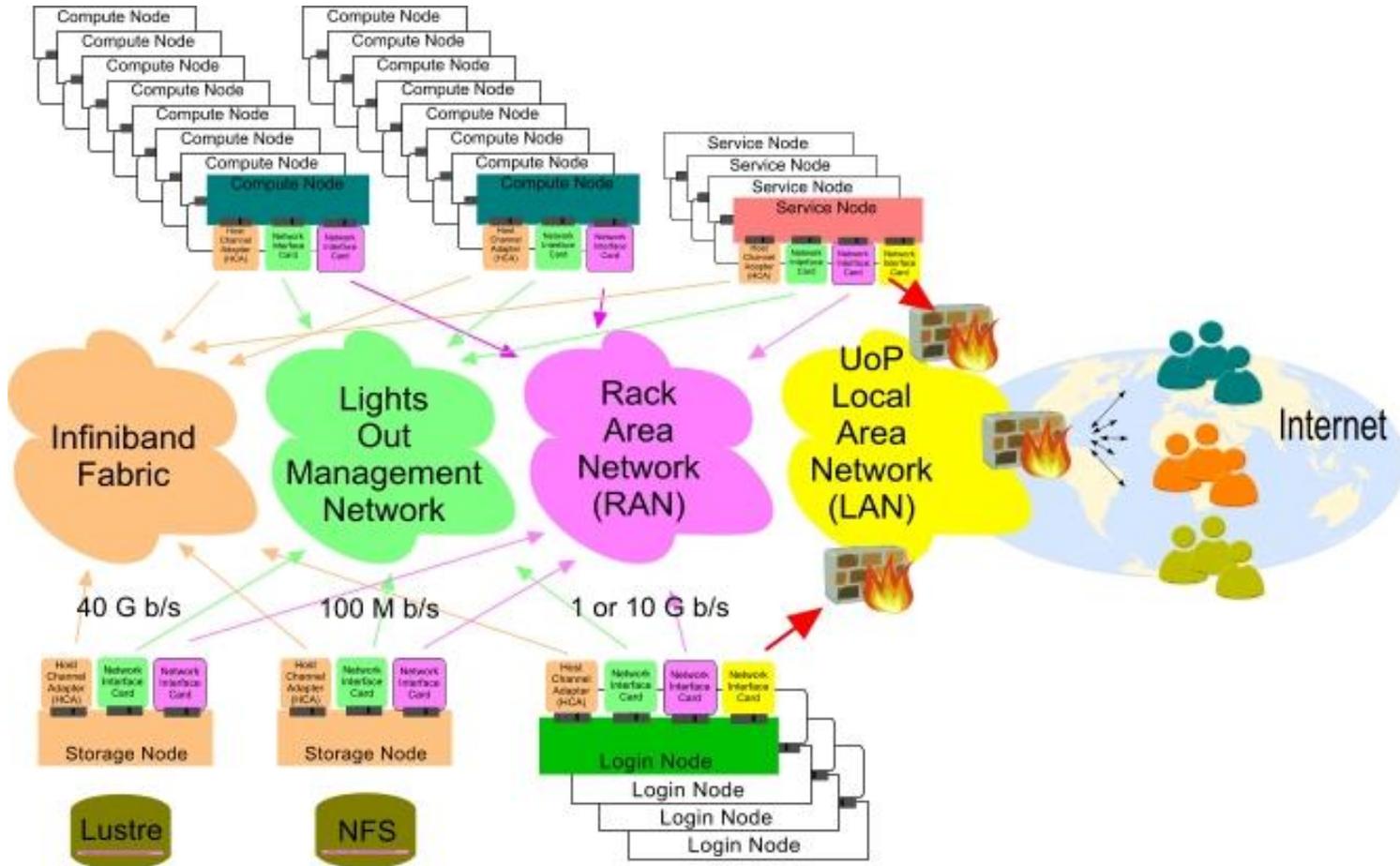
So long and thanks for all the fish !!

C: Hitch Hikers Guide to the Galaxy



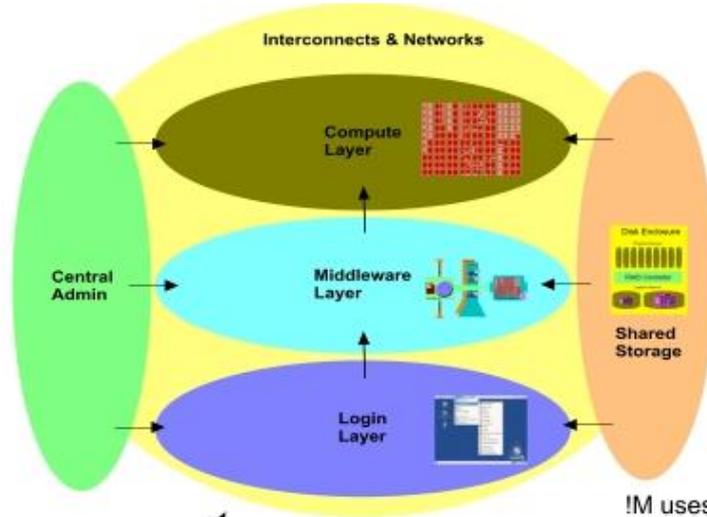


Sciama Network



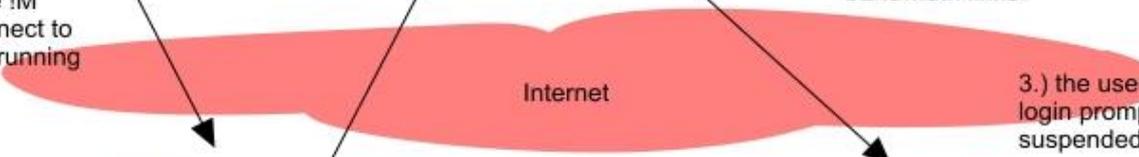
Use of Remote Login Client

1.) User downloads IM client from Internet. Client is available for Linux, Windows and MAC. The client is free.



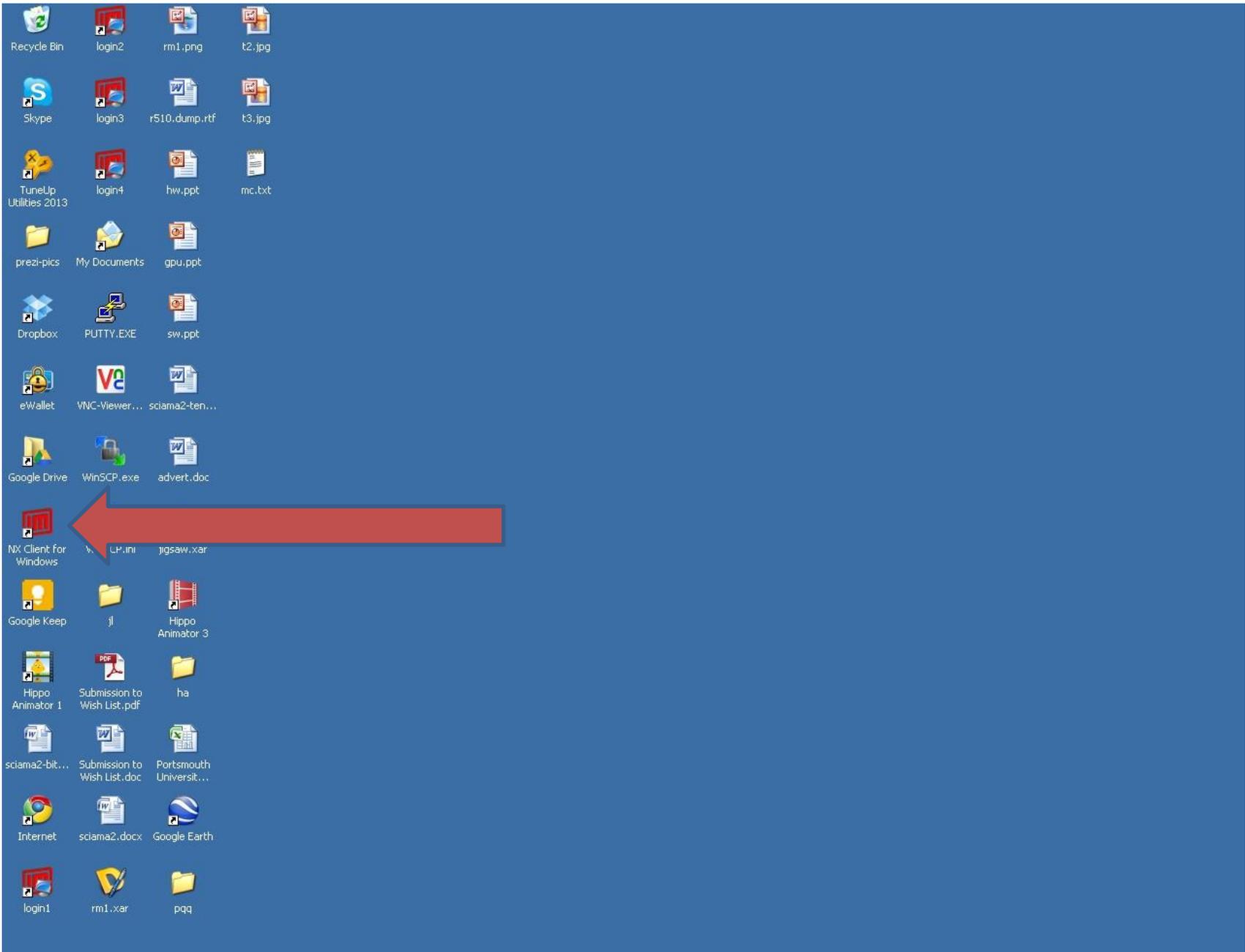
!M uses a very lean protocol for use over high latency low bandwidth links.

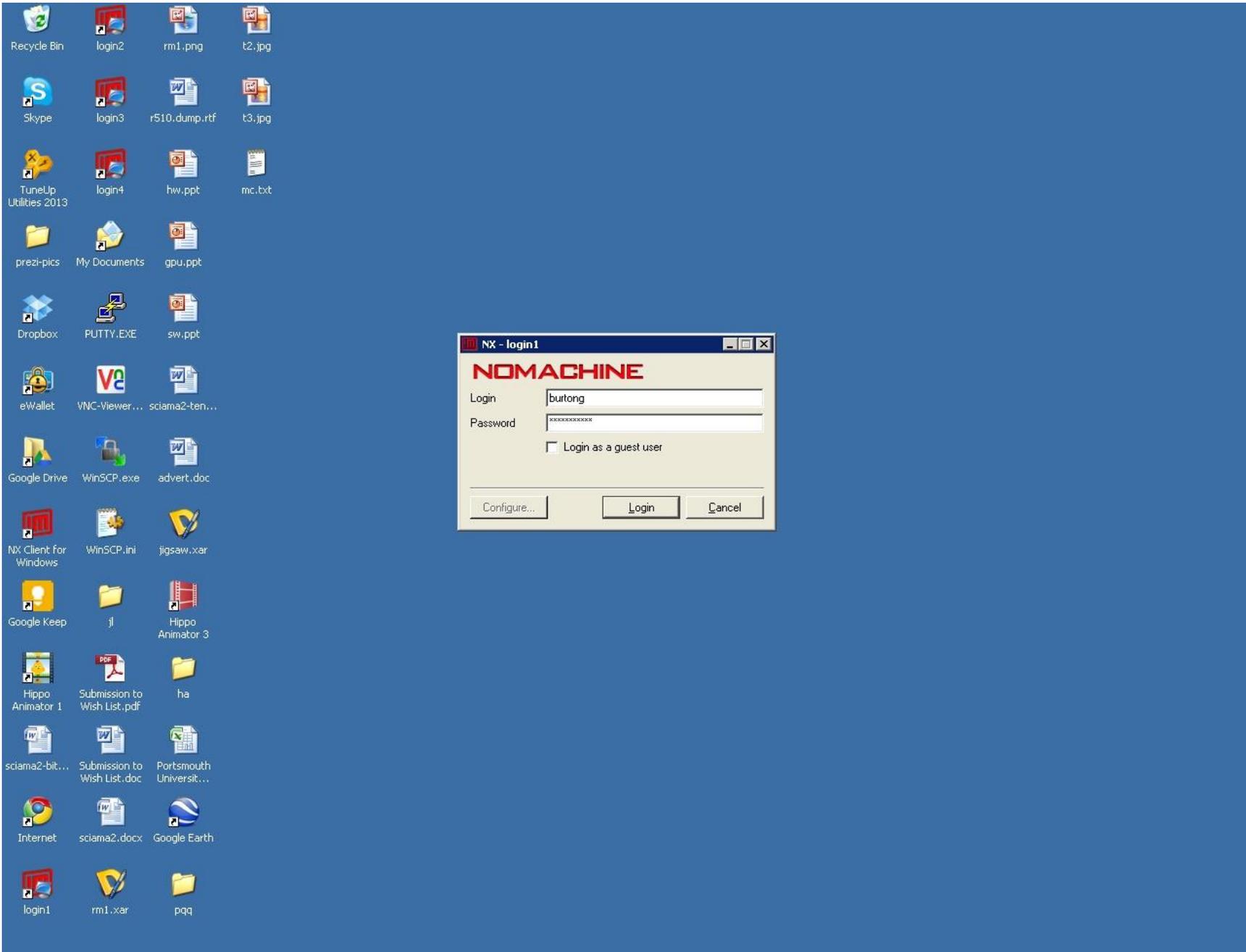
2.) Once installed and configured the IM client can connect to the IM server running in the cluster.



3.) the user is returned a login prompt or a suspended session.







NX - login1

NOMACHINE

Login:

Password:

Login as a guest user

```

File Edit View Terminal Ta
Filesystem 1K-b1
/dev/mapper/system-root
/dev/mapper/system-tmp
/dev/mapper/system-var
/dev/sda1
tmpfs
headnode:/opt/alces
headnode1:/opt/gridware
headnode1:/home
nfs1:/exports/astro1
nfs1:/exports/astro2
nfs1:/exports/astro3
nfs1:/exports/astro4
nfs1:/exports/astro5
mdsl-1b@o21b0:headnode1-ib
[bartong@login1 mnt]$ cd /
[bartong@login1 astro2]$ ls
johnston lofar sudss
[bartong@login1 astro2]$

```

new file

```

root@headnode1:/users/burtong/stats
File Edit View Terminal Tabs Help
Node 1-74 -10 0 cl3 cl3 cl3 cl3 cl3 cl3 cl3 cl3 cl3 cl3 cl3 cl3 --- ---
Node 1-75 -0 0 --- --- --- --- --- --- --- --- --- --- --- ---
Node 1-76 -12 0 ang ang ang ang ang ang ang ang ang ang ang ang
Node 1-77 -0 0 --- --- --- --- --- --- --- --- --- --- --- ---
Node 1-78 -10 0 --- joh cul cul cul cul cul cul --- ste wil joh
Node 1-79 -9 0 lco wil cul --- cul cul cul cul --- --- cul ste

Summary of cores currently not allocated

00 01 02 03 04 05 06 07 08 09
-----
Node 0 12 12 12 12 12 3 2 12 0 0
Node 10 0 0 12 12 12 12 12 12 12 12
Node 20 12 12 12 12 12 12 12 12 12 12
Node 30 12 12 12 12 12 12 12 12 12 0
Node 40 12 12 12 12 12 12 12 12 12 12
Node 50 12 12 12 12 12 12 12 0 3 12
Node 60 12 12 12 2 2 12 12 12 12 4
Node 70 11 2 0 12 2 12 0 12 2 3

node63:ppn=2
node33:ppn=12
You have
[root@h

```

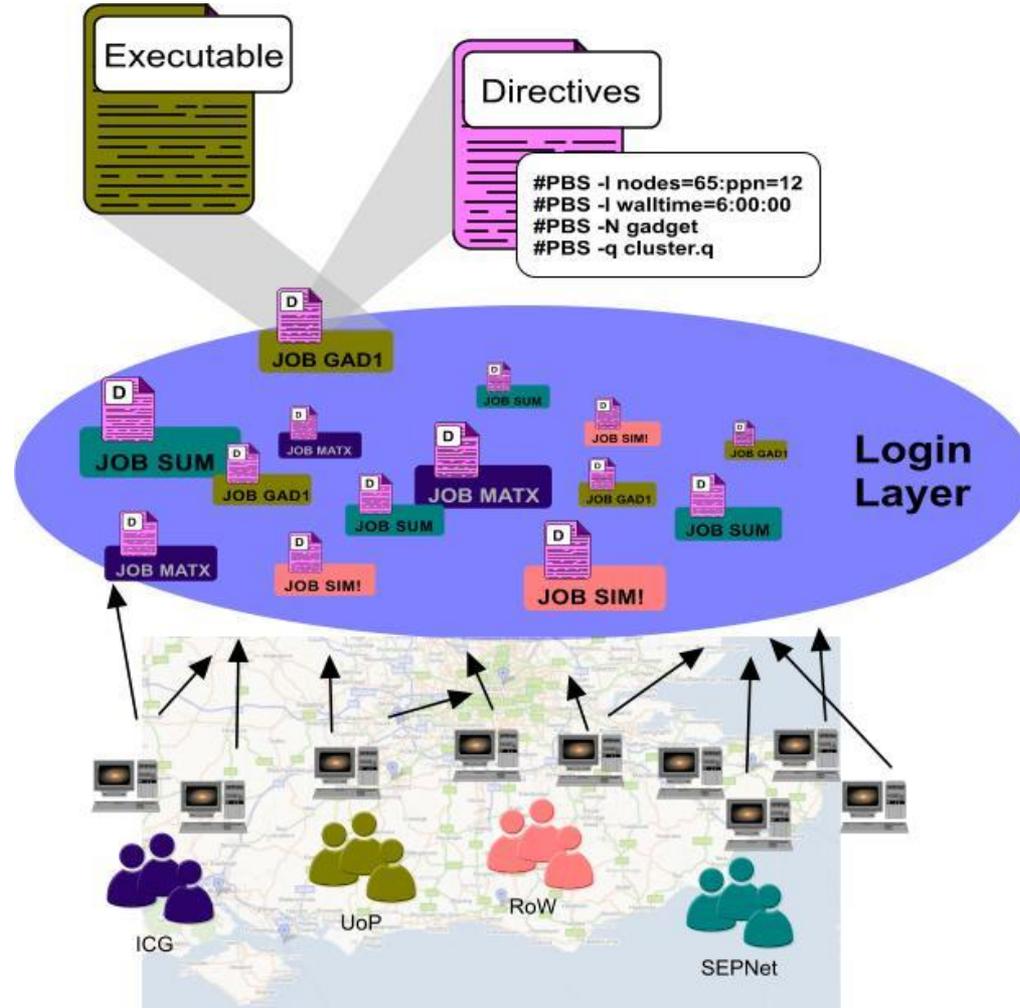
```

campbelh@login1:~/lustre/campbelh/cosmomc_6_1_2014/cosmomc/source
File Edit View Terminal Tabs Help
SDSS_sim_hubble_6_7_2012_new_lc_x1_color_centre_0_0.eps
SDSS_sim_hubble_MASA0.eps
SDSS_sim_hubble_residual_10_7_2012_test.eps
SDSS_sim_hubble_residual_6_7_2012_new_lc_x1_color_centre_0_0.eps
SDSS_sim_hubble_residual_MASA0.eps
SDSS_sim_MASA0_non_Ia_griz_remove_quality_pia_chi_color_x1_circle_centre_0_0_gr_
vs_g_test2_new_0_0.eps
SDSS_sim_MASA0_x1_vs_color_Pia_quality_cut_more_cuts_test2_centre_0_0.eps
SDSS_sim_max_sn_r_10_7_2012_test.eps
SDSS_sim_max_sn_r_6_7_2012_new_lc_x1_color_centre_0_0.eps
SDSS_sim_max_sn_r_MASA0.eps
SDSS_sims_10_7_2012_test
SDSS_sims_20_times_10_7_2012
SDSS_sims_20_times_opt
SDSS_sims_26_6_2012
SDSS_sims_4_7_2012_test
SDSS_sims_alpha_0_199_beta_3_19
SDSS_sims_HC-mods_a_0_22
SDSS_sims_HC-mods_a_0_22_20_times
SDSS_sims_new
SDSS_sims_test_color
SDSS_sim_x1_10_7_2012_test.eps
SDSS_sim_x1_6_7_2012_new_lc_x1_color_centre_0_0.eps
SDSS_sim_x1_color_10_7_2012_test.eps

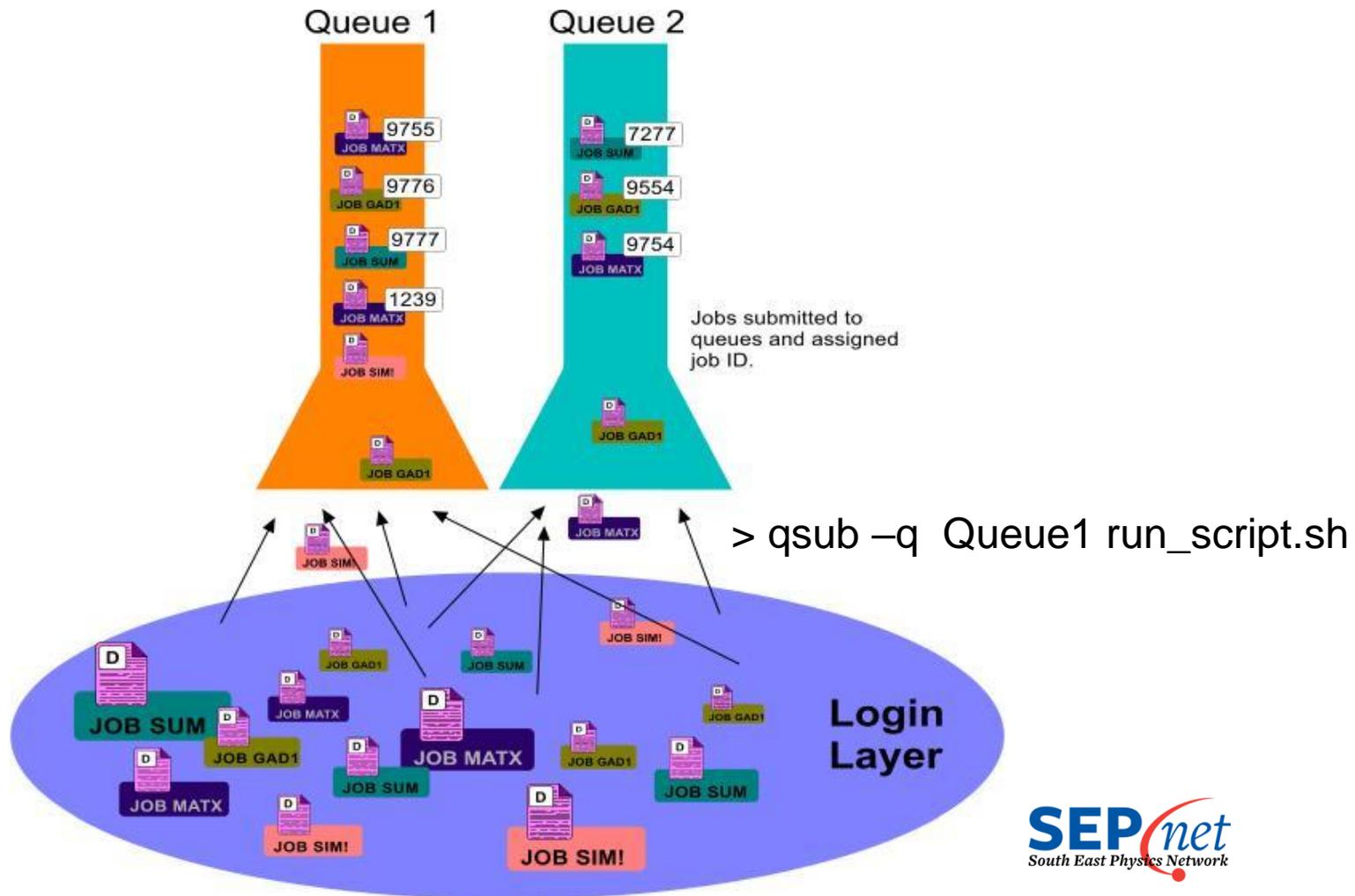
```

MA
mouth

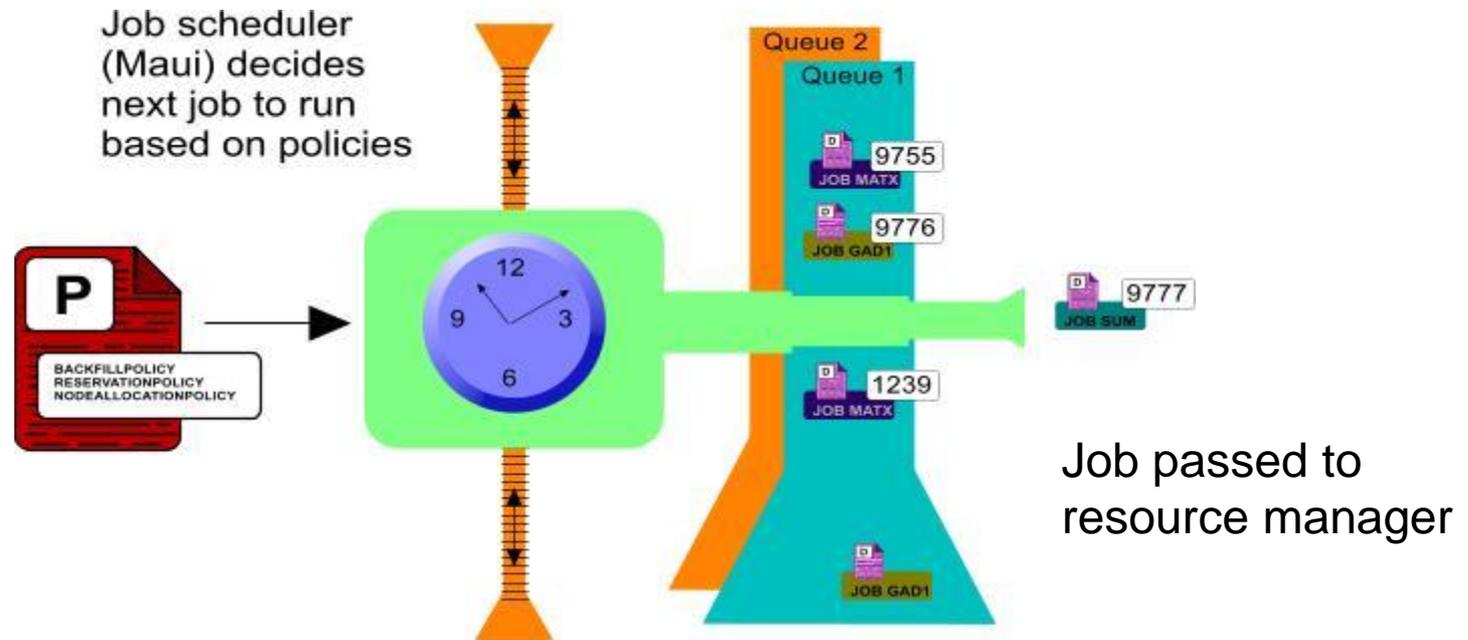
Executable and Jobscript setup in Login Layer



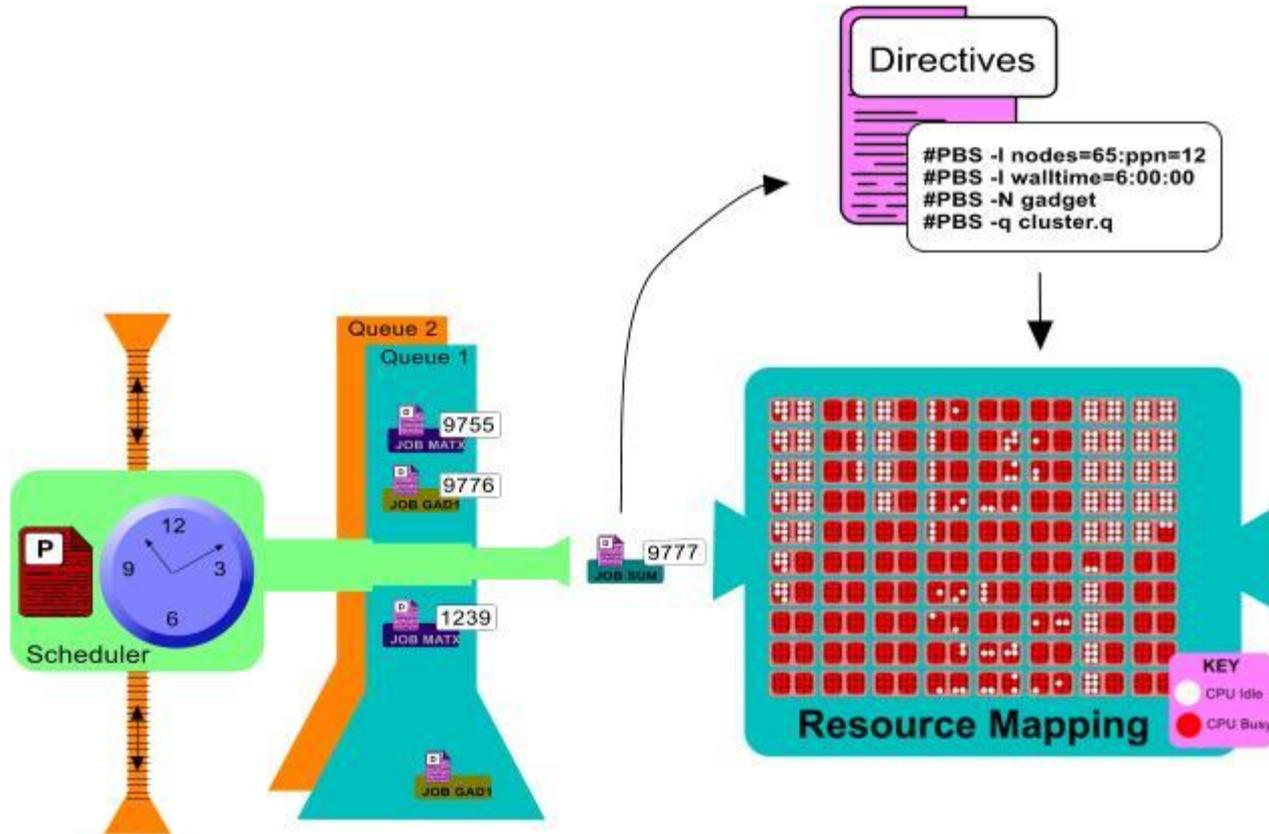
Jobs submitted to the queues



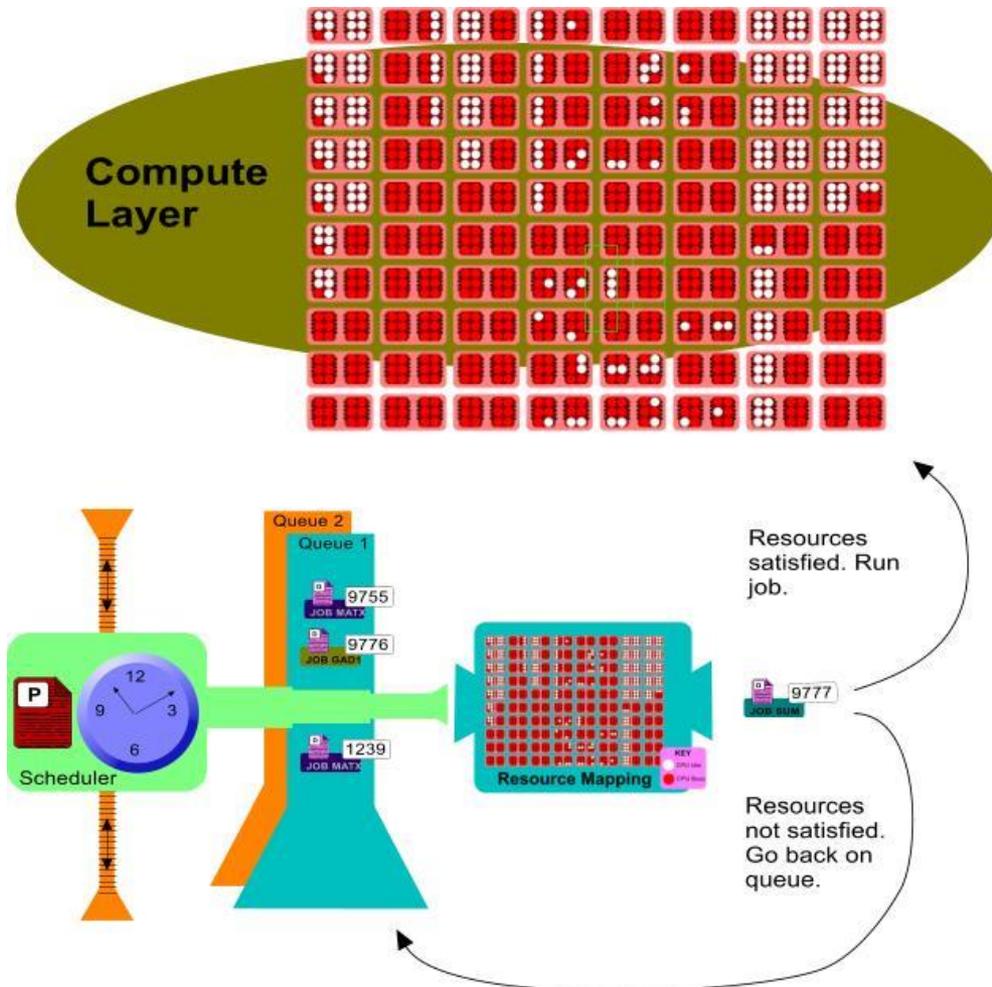
Scheduler prioritises and deems a job ready to run.



Resource manager (Torque) checks for available resources.



Job either runs in the compute pool or returns to queue



So long and thanks for all the ~~fish~~ ^{Coffee} !!

C: Hitch Hikers Guide to the Galaxy

