

**PBS Pro
ALPS
SVN mirror repo
- best practices on MIHIR**

mihir_support@ncmrwf.gov.in

Agenda

PBS

MIHIR queue details

job details

ALPS

basic architecture

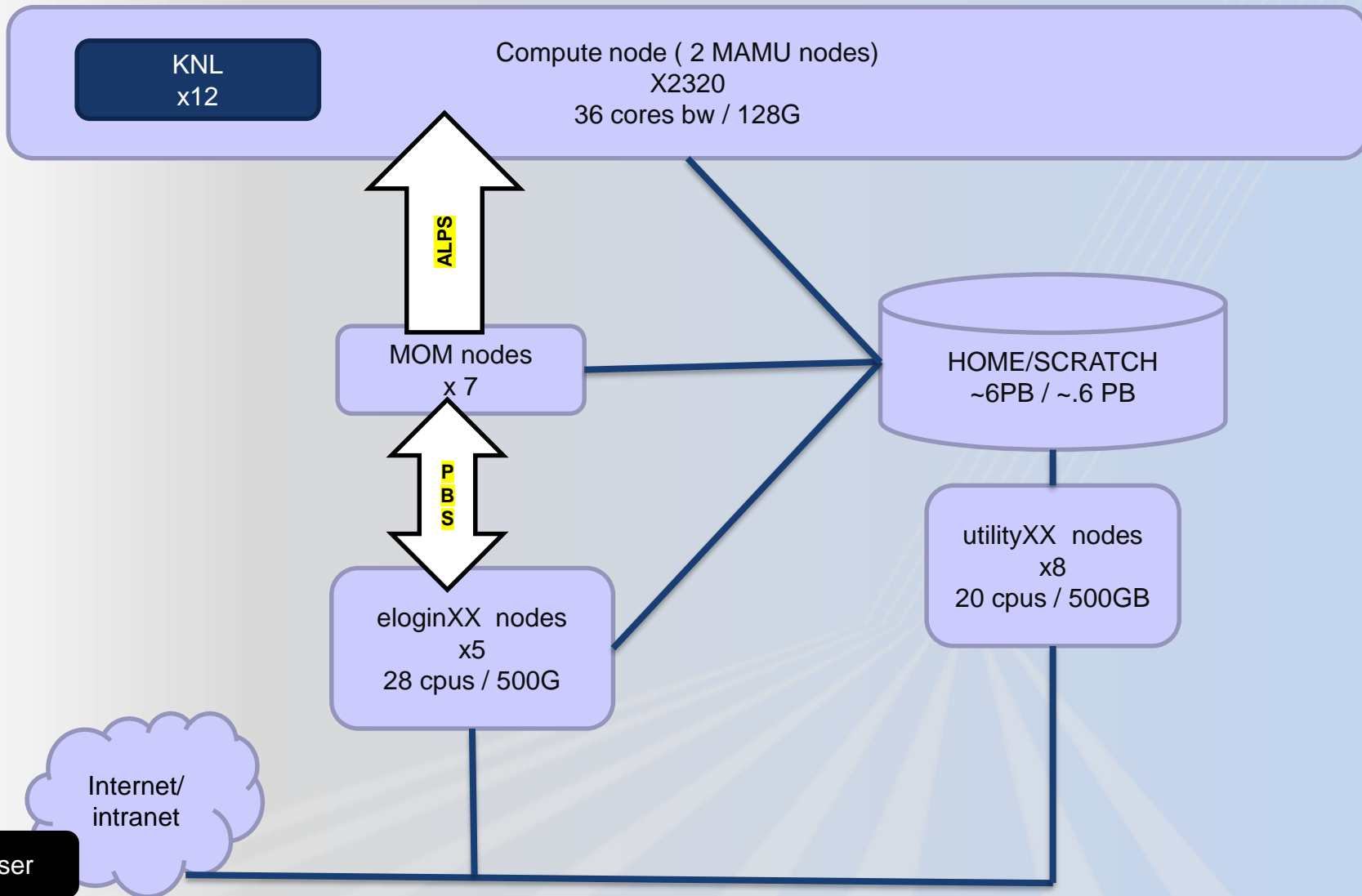
utilizing compute nodes efficiently (flags)

SVN Repo (*current setup*)

Redmine (*Ticketing system*)

Best Practices & Common application issues (Do's Don'ts) – on MIHIR

MIHIR Architecture



Software environment

module avail

module load

module show (which version ? Site installed or cray default)

Site Installed Softwares -

export MODULEPATH=\$MODULEPATH:/home/apps/SiteModules/mostapps

put up custom environment in : ~/.bash_profile instead of .bashrc

Default Environment - PrgEnv-cray/6.0.4

For loading intel environment - module switch PrgEnv-cray/6.0.4 PrgEnv-intel

For loading gnu environment - module switch PrgEnv-cray/6.0.4 PrgEnv-gnu

Note: Only on login/compute (not available on utility nodes)

Compilation

ftn - fortran compiler

cc - c compiler

CC - CXX compiler

Feature	Cray	Intel	GNU
Listing	-hlist=a	-opt-report=2	-fdump-tree-all
Free format (ftn)	-f free	-free	-ffree-form
Floating point optimization	-hfpN, N=0...4	-fp-model [fast fast=2 precise except strict]	-f[no-]fast-math or, -funsafe-math-optimizations
OpenMP recognition	(default)/-hnoomp	-qopenmp	-fopenmp
Variables size (ftn)	-s real64 -s integer64	-real-size 64 -integer-size 64	-freal-4-real-8 -finteger-4-integer-8

PBS

Portable Batch System

module load pbs

QUEUES

Queue	Max	Tot	Ena	Str	Que	Run	Hld	Wat	Trn	Ext	Type
-----	----	----	----	----	----	-----	----	----	----	----	----
NCMRWF	0	0	yes	yes	0	0	0	0	0	0	Exec
COMPROD	0	0	yes	yes	0	0	0	0	0	0	Exec
NCMRWF1	0	0	yes	yes	0	0	0	0	0	0	Exec
INCOIS	0	0	yes	yes	0	0	0	0	0	0	Exec
IMD	0	0	yes	yes	0	0	0	0	0	0	Exec
Serial	0	0	yes	yes	0	0	0	0	0	0	Exec
serial1	0	0	yes	yes	0	0	0	0	0	0	Exec

NCMRWF Queue

acl groups	=	prod1
Excluded users	=	cmprod, erfpod, ocnprod, umrda
Maximum NCPU	=	21600
Maximum Nodes	=	600
Walltime	=	08:00:00
vntype	=	cray_compute
chunk Qlist	=	NCMRWF
Default Chunk Vntype	=	cray_compute

COMPROD Queue

acl users	=	gfsprod & imdges
Maximum NCPU	=	18864
Maximum Nodes	=	524
Walltime	=	06:00:00
vntype	=	cray_compute
chunk Qlist	=	COMPROD
Default chunk vntype	=	cray_compute

NCMRWF1 Queue

acl groups	=	prod1, staff
Excluded users	=	umeps, umfcst, umprod, umreg
Maximum NCPU	=	14400
Maximum Nodes	=	400
Walltime	=	06:00:00
vntype	=	cray_compute
chunk Qlist	=	NCMRWF1
Default chunk vntype	=	cray_compute

INCOIS Queue

acl groups	=	incois
Maximum NCPU	=	28224
Maximum Nodes	=	784
Walltime	=	48:00:00
vntype	=	cray_compute
chunk Qlist	=	INCOIS
Default chunk vntype	=	cray_compute

IMD Queue

acl groups	=	imd
Excluded users	=	imdgefs
Maximum NCPU	=	8064
Maximum Nodes	=	224
Walltime	=	06:00:00
vntype	=	cray_compute
chunk Qlist	=	IMD
Default chunk vntype	=	cray_compute

SERIAL Queue

acl groups	=	prod1
Excluded users	=	imdgefs
Maximum NCPU	=	72
Maximum Nodes	=	2
Default Walltime	=	06:00:00
vntype	=	cray_mamu
Default chunk vntype	=	cray_mamu

SERIAL1 Queue

acl groups	=	prod1 & staff
Excluded users	=	cmprod, ocnprod, umeps, umfcst, umrda, umreg
Maximum NCPU	=	72
Maximum Nodes	=	2
Default Walltime	=	06:00:00
vntype	=	cray_mamu
Default chunk vntype	=	cray_mamu

PBS Commands

- ❖ `qsub` : Submit a job

`qsub myjob.pbs`
- ❖ `qdel` : Delete a batch job

`qdel <job id>`

`qdel 123[5]` for array job
- ❖ `qstat` : Show status of batch jobs (Job states: Q R E F B)

`qstat -r` running jobs

`qstat -u` user job status

- `pbsnodes` : List the status and attributes of all nodes in the cluster.
- ❖ `pbsnodes -Sja|grep -i free`

PBS Job Attributes

Attribute	Values	Description
-l	comma separated list of required resources (e.g. select=1)	Defines the resources that are required by the job and establishes a limit to the amount of resources that can be consumed. If it is not set for a generally available resource, the PBS scheduler uses default values select=1, ncpus=36 & vntype=cray_compute.
-N	name for the job	Declares a name for the job
-o	[hostname:]pathname	Defines the path to be used for the standard output (STDOUT) stream of the batch job.
-e	[hostname:]pathname	Defines the path to be used for the standard error (STDERR) stream of the batch job.
-q	name of destination queue, server, or queue at a server	Defines the destination of the job. The default setting is sufficient for most purposes.

PBS Resources

Attribute	Values	Description
select	positive integer	Declares the node configuration for the job.
walltime	hh:mm:ss	Specifies the estimated maximum wallclock time for the job.
ncpus	positive integer	Declares the number of CPUs requested.

PBS Environment Variables

Variable	Description
PBS_ENVIRONMENT	set to PBS_BATCH to indicate that the job is a batch job; otherwise, set to PBS_INTERACTIVE to indicate that the job is a PBS interactive job
PBS_JOBID	the job identifier assigned to the job by the batch system
PBS_JOBNAME	the job name supplied by the user
PBS_QUEUE	the name of the queue from which the job is executed
PBS_O_HOME	value of the HOME variable in the environment in which qsub was executed
PBS_O_LOGNAME	value of the LOGNAME variable in the environment in which qsub was executed
PBS_O_PATH	value of the PATH variable in the environment in which qsub was executed
PBS_O_WORKDIR	the absolute path of the current working directory of the qsub Command

ALPS

Application Level Placement Scheduler

ALPS

Cray supported mechanism for placing and launching applications on compute nodes.

PBS make policy and scheduling decisions, while ALPS provides a mechanism to place and launch the applications contained within batch jobs.

Function

Assigns application IDs

Launches applications

Delivers signals to applications

Works with the CLE Node Health Software
(NHC)

Architecture

1. ALPS Clients

aprun

apstat

apkill

apmgr

apbasil

Architecture

2. ALPS Daemons

apsys

apinit

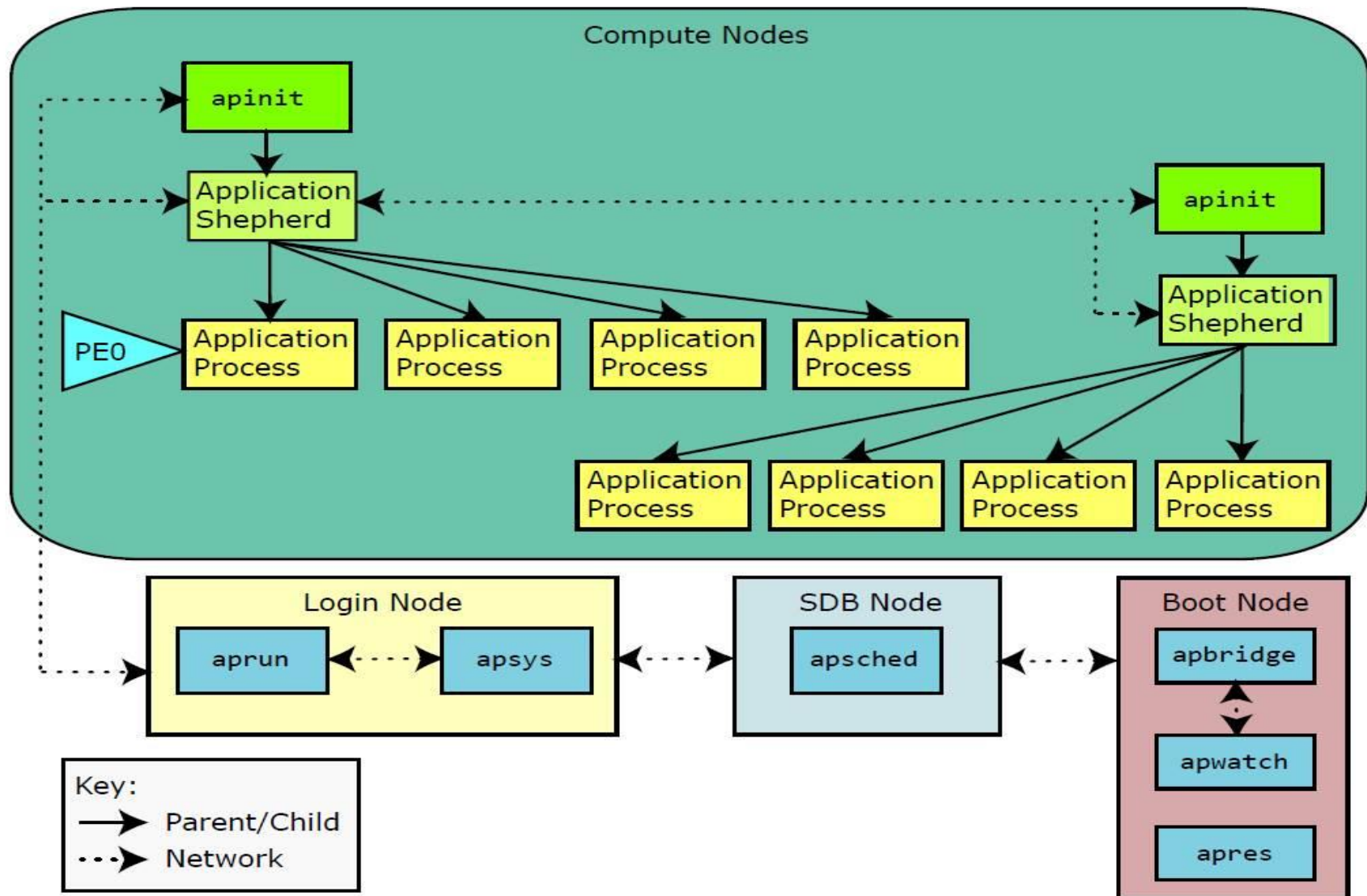
apsched

apbridge

apwatch

apres

ALPS flow



Sample PBS Script (submission)

Script

```
#!/bin/sh
#PBS -N job101
#PBS -l walltime=1:00:00
#PBS -q NCMRWF
#PBS -l select=2:mem=1gb

cd $PBS_O_WORKDIR

aprun -n72 -N36
./my_application
```

Submission/Comments

```
qsub my_script.pbs
qsub -q NCMRWF1 my_script.pbs

qsub -v
INFILE=/tmp/myinfile,INDATA=/tmp/out -q
NCMRWF1 my_script.pbs

qsub -J 10-20:2 my_script.pbs
```

Q: how to avoid running job on mom nodes?

Q: job starts & runs where?

Handling Output and Error redirection

Linux shell I/O redirection symbols ">" and "<" must NOT be used with the job submission command in their job scripts.

Not using redirection symbols leads to creating of the respective files in PBS configuration spaces and fills it up, leading to loss of PBS job submission service. Users must open and close their input and output files through their executables and avoid using these symbols.

On large scale distributed systems like the Cray-XC40 the output and error files generated by your jobs need to be explicitly stated in your jobs scripts. Please use full path names to write these files to your area on the space /home OR /scratch. Using relative paths or just file names in your job scripts pushes these files to PBSPro system spool area and fills that space. As a result the PBSPro scheduler no longer can dispatch jobs for execution and the batch scheduler stalls. This renders the system into a stall mode and every user gets affected and no job gets dispatched.

Handling Output and Error redirection ...

1. Try using specific file descriptors other than stdout and stderr in your executables. These file descriptors can then be set to be written to /home OR /scratch space directly in your source code or through appropriate directives or environment variables of your opensource and licensed codes.
2. The following #PBS directives should not be used/included in the job submission scripts

```
#PBS -o /home/.../output.log  
#PBS -e /home/.../error.log  
#PBS -j oe
```

3. While using the input/output redirection symbols(<,>), please make sure the absolute paths of your input,output and executable file location are specified and pointed to /home OR /scratch space.

Example –

```
aprun -n 24 -N 24 ./a.out &> {out-err-file}  
aprun -n 24 -N 24 ./a.out > {output-file} 2> {error-file}  
e.g. aprun -n 24 -N 24 ./a.out &> out-err.txt
```

aprun command

command for launching an application on compute node(s).

The aprun command can't use more resources (reserved using the qsub command)

Replacement of mpiexec/mpirun – for MPI jobs on compute nodes

Description	Flags
number of PEs to place per node	-N
Number of processing elements (PEs) needed for your application	-n
Specifies the number of PEs to allocate per NUMA node	-S
Strict memory containment per NUMA node.	-ss
Specifies how many CPUs to use per compute unit for an ALPS job.	-j
Specifies the number of CPUs for each PE and its threads	-d

PBS scripts... compute nodes

For launching executable on 2 compute nodes

```
#!/bin/bash
```

```
#PBS -l select=2
```

```
#PBS -l place=scatter
```

```
#PBS -q NCMWRF1
```

```
export OMP_NUM_THREADS=1
```

```
aprun -n1 -N1 ./initialize1 & // serial job
```

```
aprun -n1 -N1 ./initialize2 // serial job
```

```
aprun -n 72 -N 36 ./application // parallel job
```

aprun examples

```
aprun -n72 -N 36 ./a.out ; export OMP_NUM_THREADS=1
                           aprun -n72 -N 36 ./a.out
```

Q: Application is pure MPI? Application is hybrid?

```
export OMP_NUM_THREADS=2 ; export OMP_NUM_THREADS=2;
aprun -n36 -N 18 -d2 ./a.out ; aprun -n36 -N 18 ./a.out
```

```
aprun -n20 -N 10 ./a.out ; aprun -n20 -N 10 -ss ./a.out
aprun -n20 -N 10 -S 5 ./a.out ;
```

PBS scripts... MAMU nodes

```
#PBS -l select=1:ncpus=8:vntype=cray_mamu  
#PBS -q serial1
```

```
./initialize1 &  
./initialize2 &  
module load cray-snplauncher/7.6.3  
mpiexec -n[-np] 6 ./application
```

Shared node !! Performance ?

Do's

Report any unexpected slowdown/ file missing / corrupted related issues to mihir_support@ncmrwf.gov.in asap.

On MAMU nodes for MPI jobs, ensure you have loaded snplauncher (/opt/pbs/default/bin/mpiexec)

As per requirement , please use -ss (suspend/resume is disabled).

Ensure you are using correct vntype

Setting crons/rose jobs on utility/elogin nodes? – inform mihir_support with job start time & duration.

Ensure aprun is prefixed with binaries (example – grads, um-atmos.exe)

Environment settings via .bash_profile

Ensure you are not using old PBS syntax

```
#PBS -l nodes=3:ppn=8
```

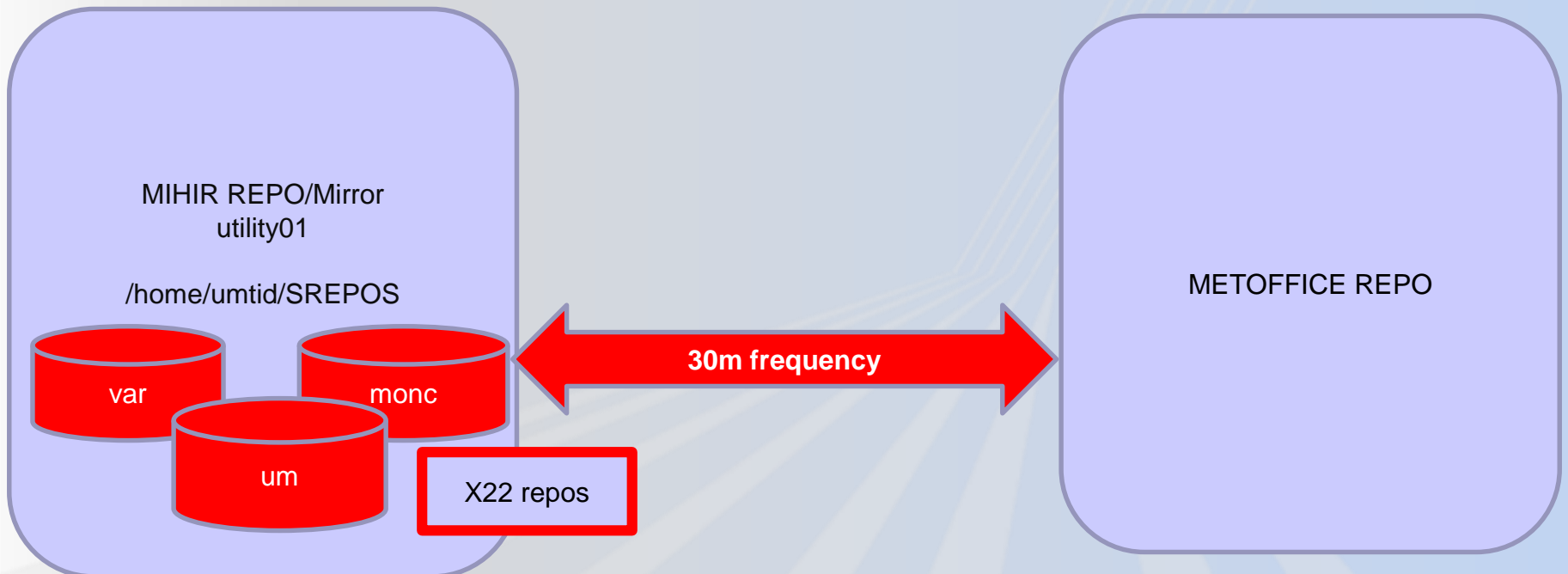

Rose suite

In case of abrupt shutdown of rose suite -

remove ~/cylc-run/erfgc2/.service/contact

kill suite's python process via (ps -ef|grep python|grep \$USER)

SVN Local Repo



Do's

1. Check branch – metoffice website
 2. Raise a ticket
 3. Create your own branch (via utility)
 4. Check-Out from your branch (via local mirror)
 5. Carry out your changes
 6. commit to your branch (from utility node)
1. `fcmls
svn://192.168.0.130/um/main
/trunk@rev`
 2. `--`
 3. `fcmbc -k ticket_number
https://code.met..../sv/repo/tr
unk@rev
into_your_own_branch
(utility)`
 4. `fcmlco your_new_branch`
 5. `touch dev_1/newfile.c; vim
dev_1/newfile.c`
 6. `fcmlcommit -m "physics
changed"`

Don'ts

Avoid checking out from metoffice repo, please use local mirror (max wait should be 30 minutes)

Local mirror is read-only, only checkout – don't commit.

**Avoid using .fcm files, use keyword.cfg file
~/.metomi/fcm/keyword.cfg**

Thank You!

The background features a gradient from white on the left to light blue on the right. At the bottom, there are several wavy, parallel lines in shades of light blue and white that create a sense of depth and movement, resembling a stylized horizon or a path leading into the distance.