

Serveur de calcul « nymphéa »

Guide de l'utilisateur (version 2)

<b>1.</b>	<b>Description de la machine .....</b>	<b>4</b>
1.1.	Architecture générale.....	4
1.2.	Processeurs .....	4
1.3.	Caractéristiques mémoire des nœuds du cluster .....	5
<b>2.</b>	<b>Accès à distance .....</b>	<b>5</b>
2.1.	Cas des utilisateurs extérieurs au réseau IFREMER .....	5
2.1.1.	Accès d'un simple poste de travail.....	5
2.1.2.	Interconnexion entre nymphéa et un serveur .....	6
2.2.	Cas des utilisateurs des sites IFREMER autres que Brest .....	6
<b>3.</b>	<b>Environnement UNIX.....</b>	<b>7</b>
3.1.	Préférences utilisateurs .....	7
3.2.	Travailler en X11 .....	7
3.3.	Ecrire des scripts simples sous UNIX.....	8
3.4.	Editeurs de texte .....	8
3.4.1.	vi.....	8
3.4.2.	nedit.....	8
3.4.3.	emacs.....	9
3.4.4.	Edition par montage NFS.....	10
3.5.	Espace disque utilisateur .....	10
3.5.1.	Zones disque.....	10
3.5.2.	Espace de stockage externe à nymphéa.....	11
3.5.3.	Protocole de sauvegarde automatique .....	11
3.5.4.	Restitution de fichier.....	11
3.6.	Utilitaires .....	12
3.6.1.	Compression de données .....	12
3.6.2.	Transfert de données .....	12
3.6.3.	Imprimer du texte sur une imprimante Postscript .....	12
3.6.4.	Autres.....	13
<b>4.</b>	<b>Les compilateurs.....</b>	<b>14</b>
4.1.	les langages disponibles.....	14
4.1.1.	fortran.....	14
4.1.2.	c.....	14
4.1.3.	c++ .....	14
4.1.4.	java.....	14
4.1.5.	compilateurs du GNU (gcc, g++, g77) .....	14
4.2.	Les bibliothèques de sous-programmes disponibles.....	15
4.2.1.	cxml .....	15
4.2.2.	cxmlp.....	15
4.2.3.	scalapack (ou appellation CS chez HP-COMPAQ) .....	15
4.2.4.	nag .....	15
4.2.5.	netcdf.....	15
<b>5.</b>	<b>Exécuter un programme .....</b>	<b>17</b>
5.1.	En interactif.....	17

## sommaire

5.2.	En batch avec LSF.....	17
5.2.1.	principes du batch .....	17
5.2.2.	les files batchs sur nymphéa.....	17
5.2.3.	Lancer des commandes LSF à partir du WEB INTRANET .....	19
5.3.	Visualiser les jobs en cours sur les machines .....	19
5.3.1.	Jobs actifs .....	19
5.3.2.	Jobs par queue batch.....	20
<b>6.</b>	<b>Le debugger – TOTALVIEW.....</b>	<b>21</b>
<b>7.</b>	<b>Optimiser le code .....</b>	<b>24</b>
7.1.	Mesurer le temps passé dans le code.....	24
7.1.1.	L'outil « prof ».....	24
7.1.2.	Le script cpuProfiler .....	26
7.1.3.	La commande FORTRAN « cpu_time » .....	27
7.1.4.	La commande UNIX « time » .....	27
7.1.5.	Les compte-rendus d'exécution de LSF .....	28
7.2.	La parallélisation .....	28
7.2.1.	Parallélisation OpenMP.....	29
7.2.2.	La parallélisation MPI.....	30
<b>8.</b>	<b>Les progiciels installés.....</b>	<b>32</b>
8.1.	fluent .....	32
8.2.	gcg – Wisconsin package .....	32
8.3.	abaqus .....	32
8.4.	Phylowin.....	32
8.5.	Seaview .....	32
8.6.	Matlab .....	33

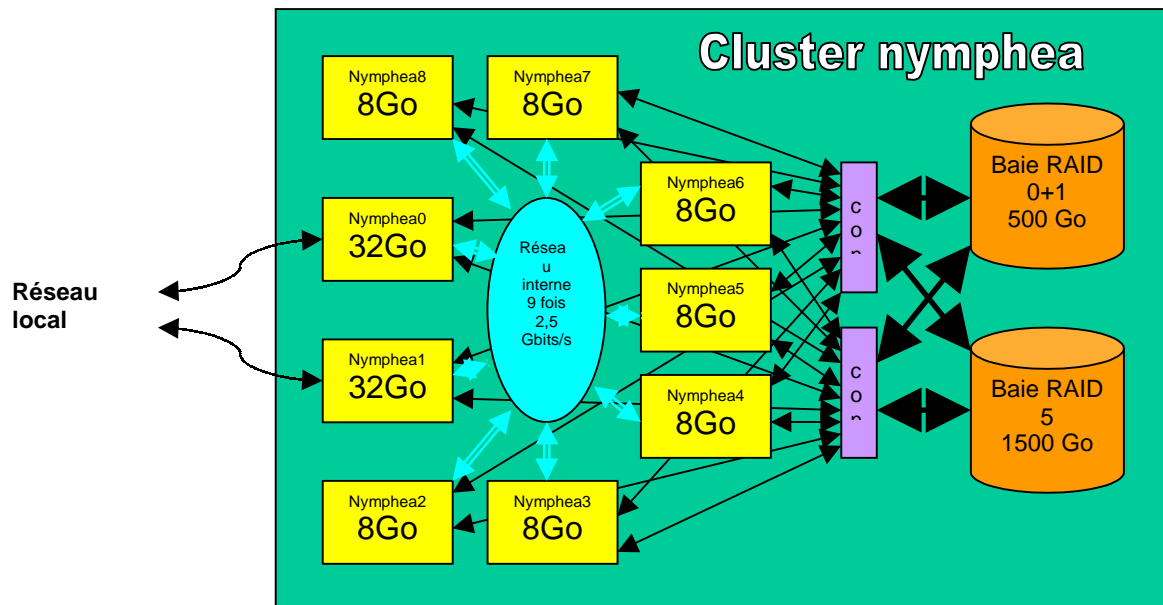
## 1. Description de la machine

La machine est un cluster de 9 nœuds ES45 quadri-processeurs interconnectés par un switch Quadrics. Les processeurs sont des Alpha EV68 à 1 GHz.

2 des nœuds disposent de 32 Go de mémoire et les 7 autres de 8 Go, ce qui fait une mémoire globale de 120 Go. L'espace disque se divise en 2 baies distinctes; la première de 500 Go en RAID 0+1 pour favoriser les performances et la redondance maximales, la seconde de 1500 Go est en RAID 5.

### 1.1. Architecture générale

Le cluster est connecté au réseau IFREMER via une interface Ethernet Gigabit redondée.



### 1.2. Processeurs

Les 36 processeurs sont des Alpha EV68 à 1 GHz, qui comportent 2 niveaux de mémoire cache (L1 et L2). Le tableau qui suit indique les performances comparées des différents niveaux de mémoire.

Zone mémoire	capacité	Temps d'accès	Débit lecture
Cache L1	<b>64 Ko par CPU</b>	<b>2 ns</b>	<b>6,47 Go/s</b>
Cache L2	<b>8 Mo par CPU</b>	<b>19 ns</b>	<b>6,07 Go/s</b>
Mémoire centrale	<b>8 Go par nœud (au moins)</b>	<b>170 ns</b>	<b>2,27 Go/s</b>

### 1.3. Caractéristiques mémoire des nœuds du cluster

- les 2 premiers nœuds (nympha0 et nympha1) ont chacun 32Go
- les autres (nympha2-----nympha8) ont chacun 8Go

## 2. Accès à distance

### 2.1. Cas des utilisateurs extérieurs au réseau IFREMER

#### 2.1.1. Accès d'un simple poste de travail

En faisant une demande auprès de l'assistance l'assistance (tel : 0298224207, e\_mail : [assistance@ifremer.fr](mailto:assistance@ifremer.fr)), un utilisateur peut obtenir un login « extranet » qui lui servira à entrer sur le réseau IFREMER, puis à se connecter sur nympha.

L'accès se fait par le biais d'un navigateur à <https://portail.ifremer.fr/>. Après avoir accepté le certificat, l'écran de login s'affiche.

**Ifremer**  
TMSI/IDM/RIC

**Nokia**  
**Secure Access System**

**Connexion standard**

Nom d'utilisateur :

Mot de passe :

**Connexion avec certificat**

Enabled by Nokia  
Copyright © 2002-2004 Nokia

[English](#) [Français](#)

La connexion sur nympha se fait par le biais d'une session SSH encapsulée dans la session HTTPS.

#### 2.1.2. Interconnexion entre nympha et un serveur

Ce serveur aura nécessairement une adresse IP fixe ; sur la base de cette adresse et d'un protocole IP, un chemin particulier peut être ouvert sur notre firewall pour permettre les échanges automatisés entre nympha et ce serveur.

#### 2.2. Cas des utilisateurs des sites IFREMER autres que Brest

Pour les utilisateurs Ifremer d'un autre site que Brest, il est recommandé de travailler en local sur nympha, et donc d'y avoir un espace disque sur lequel seront les fichiers en entrée comme en sortie, pour que le temps de transfert ne pénalise pas la vitesse de calcul.

Il est toujours assez tôt, après le calcul, pour récupérer les fichiers résultats sur la machine locale.

## 3. Environnement UNIX

### 3.1. Préférences utilisateurs

Pour avoir le bon environnement de travail, nous vous recommandons de vous inspirer ou tout simplement de copier dans vos `.cshrc` et `.login` les modèles que vous trouverez sur le serveur brest, dans le répertoire:

`/home/services/exemples` , sous les noms :

`.cshrc.cpq` et `.login.cpq`

Pour les utilisateurs extérieurs à l'Ifremer, ces fichiers sont automatiquement générés.

Pour les utilisateurs travaillant à la fois sur nymphaea et sur solaris par exemple, il est conseillé d'avoir un `.cshrc` et `.login` « aiguillage » du type :

<pre># .cshrc file switch ( `uname -n` ) case nymphaea*:     source ~/.cshrc.cpq     breaksw default:     source ~/.cshrc.solaris endsw</pre>	<pre># .login file switch ( `uname -n` ) case nymphaea*:     source ~/.login.cpq     breaksw default:     source ~/.login.solaris endsw</pre>
---	---

### 3.2. Travailler en X11

On utilise la norme de multifenêtrage X11 pour ouvrir des sessions interactives sur nymphaea. La fonction "client" X11 est présente dans tous les systèmes UNIX ou LINUX; pour les postes de travail PC MS-Windows, il faut installer un logiciel d'émulation X11 comme EXCEED de Hummingbird, que RIC recommande.

### 3.3. Ecrire des scripts simples sous UNIX

#### 3.4. Editeurs de texte

##### 3.4.1. vi

Cet éditeur conserve une bonne interactivité lorsque l'accès se fait par un réseau distant. Par contre, toutes les actions se passent par des commandes ; il n'y a ni « ascenseur », ni « menu ».

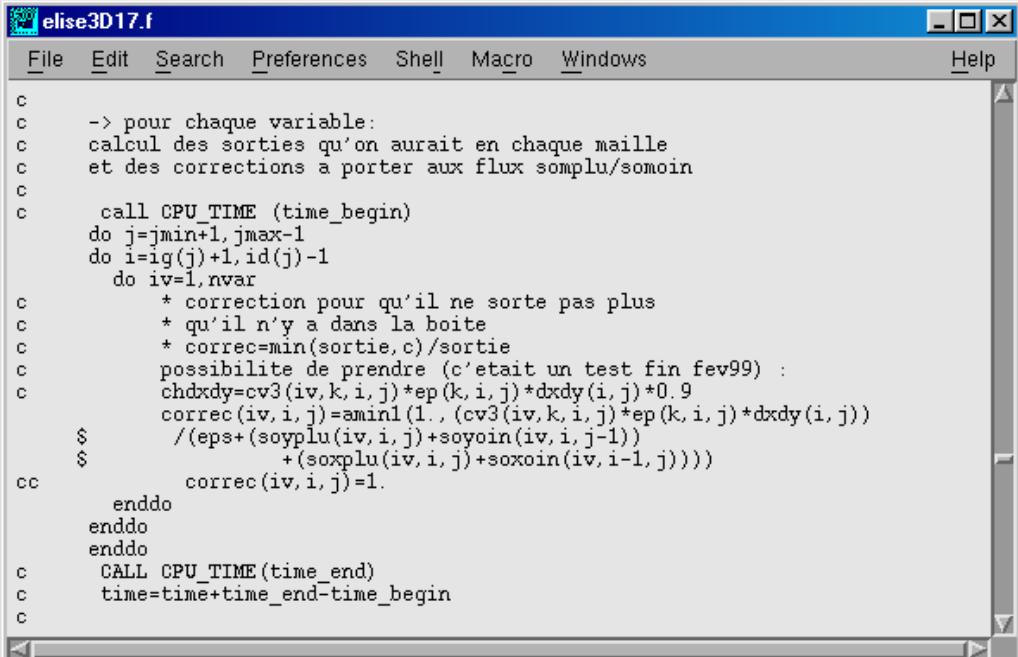
##### 3.4.2. nedit

Cet éditeur conserve une bonne interactivité lorsque l'accès se fait par un réseau distant. Il inclut toutes les fonctions « pleine page » : souris, ascenseurs, menus, ...

Après avoir vérifié que la variable DISPLAY soit correctement positionnée, lancer nedit par la commande :

```
nedit&
```

On obtient :



```

elise3D17.f
File Edit Search Preferences Shell Macro Windows Help
c
c -> pour chaque variable:
c calcul des sorties qu'on aurait en chaque maille
c et des corrections a porter aux flux somplu/somoin
c
c   call CPU_TIME (time_begin)
c do j=jmin+1, jmax-1
c do i=ig(j)+1, id(j)-1
c do iv=1, nvar
c   * correction pour qu'il ne sorte pas plus
c   * qu'il n'y a dans la boite
c   * correc=min(sortie,c)/sortie
c   possibilite de prendre (c'etait un test fin fev99) :
c   chdxxy=cv3(iv,k,i,j)*ep(k,i,j)*dxxy(i,j)+0.9
c   correc(iv,i,j)=amin1(1., (cv3(iv,k,i,j)*ep(k,i,j)*dxxy(i,j))
$   /(eps+(soyplu(iv,i,j)+soyoin(iv,i,j-1))
$   +(soxplu(iv,i,j)+soxoin(iv,i-1,j))))
cc   correc(iv,i,j)=1.
c   enddo
c enddo
c enddo
c   CALL CPU_TIME(time_end)
c   time=time+time_end-time_begin
c

```

### 3.4.3. emacs

Cet éditeur inclut toutes les fonctions « pleine page » : souris, ascenseurs, menus, ...

Après avoir vérifié que la variable DISPLAY soit correctement positionnée, lancer xemacs par la commande :

```
xemacs&
```

On obtient une fenêtre du genre :

```

*****
c
      subroutine flubo3
c
*****
c23456789012345678901234567890123456789012345678901234567890123456789012
c
      SAVE
      include 'zonecommune3Dbio'
c
      dimension soxplu(ndmv, imin:imax, jmin:jmax)
      $          , soxoin(ndmv, imin:imax, jmin:jmax)
      $          , soyplu(ndmv, imin:imax, jmin:jmax)
      $          , soyoin(ndmv, imin:imax, jmin:jmax)
      dimension raph(nijma), transp(ndmv, nijma)
      $          , c2(ndmv, imin:imax, jmin:jmax)
      $          , correc(ndmv, imin:imax, jmin:jmax)
      $          , corepx(imin:imax, jmin:jmax)
      $          , corepy(imin:imax, jmin:jmax)
      $          , eptransx(imin:imax, jmin:jmax)
      $          , eptransy(imin:imax, jmin:jmax)
      real time_end , time_begin
c
c      initialisation
c
c      -----
      data aden24 /.04166667/
      data aden48 /.02083333/
      data aden6 /0.166666666667/
      data aden3 /0.333333333333/
      eps2=2.*eps
c      aden48=1/48
c      aden24=1/24
c      aden6=1/6
c      aden3=1/3
c
      kcpmax=min(kdmax, kcmax+1)

      do i=imin, imax
-----XEmacs: flubo54p (Fundamental)----Top-----
Open a file

```

#### 3.4.4. Edition par montage NFS

Cette technique n'est réellement rapide et interactive que lorsque que l'on est sur le même réseau local que le serveur. Elle consiste en l'utilisation de son éditeur favori sur sa station de travail, qui va lire et écrire sur le disque du serveur « monté » en NFS (Network File System).

### 3.5. Espace disque utilisateur

#### 3.5.1. Zones disque

L'espace disque est divisé en 2 parties :

- 500Go de disque « noble », techniquement conforme au standard RAID 0+1 (c'est l'espace de travail courant)
- 1500Go de disque de stockage de résultats au standard RAID5

La répartition de l'espace disque est la suivante :

#### Espace de travail

home	RIC et Réserve	90 Go
home1	DEL	100 Go
home2	LPO	100 Go
home3	Autres IFREMER	100 Go
home4	ENSIETA	10 Go
home5	Ecole Navale	50 Go
home6	SHOM	50 Go

#### Espace de stockage

home7	LPO	460 Go
home8	DEL (recherche)	340 Go
home9	Autres	360 Go
home10	RIC et Réserve	200 Go
home11	DEL (labos côtiers)	80 Go

### Espace temporaire

Cet espace, d'une taille de 120 Go est un espace d'appoint sur la machine nympha. Pour l'utiliser, il vous suffit de demander à l'assistance la création d'un répertoire à votre nom sur cet espace /home12/nympha, qui est sauvegardé normalement.

La durée de vie des fichiers sur cet espace est limitée à 5 jours ; au delà, ils sont automatiquement détruits.

#### 3.5.2. Espace de stockage externe à nympha

Le serveur « » héberge des disques de grandes capacité avec des performances moyennes. Il est recommandé de faire en sorte que les jobs tournant sur nympha écrivent sur des disques locaux, décrits au §3.5.1

#### 3.5.3. Protocole de sauvegarde automatique

L'ensemble de l'espace disque (travail et stockage) est sauvegardé quotidiennement. Les sauvegardes restent en ligne pendant 2 mois, après quoi on ne conserve qu'une image trimestrielle des espaces disque. Le système de sauvegarde est le même que sur l'ensemble du réseau IFREMER.

#### 3.5.4. Restitution de fichier

L'interface de restitution est disponible sur le WEB à l'adresse :

<http://w3.ifremer.fr:580/intranet/etc/doroweb/htdocs/>

Pour pouvoir exprimer « à la mode UNIX » les noms de fichiers, il faut cliquer sur « Recherche avancée »

## Recherche à une date



**Répertoire de recherche:**  
 Vous pouvez restreindre la recherche en précisant un répertoire  
 Ex: Q:\internet\netscape\...

Q:\

Q\  N\ [Recherche avancée](#)

[Recherche entre deux dates](#)

**Date de recherche(\*):**  
 16 Mar 2005

(\*) : Positionnez-vous 1 jour après la date de la version voulue.  
 Ex : sélectionnez le 17 pour avoir la version du 16 au soir.

### 3.6. Utilitaires

#### 3.6.1. Compression de données

« compress » et « gzip » sont disponibles. Pour plus de détails taper la commande « man compress » ou « man gzip ».

#### 3.6.2. Transfert de données

Le protocole ftp est opérationnel sur le serveur nymphaea.

#### 3.6.3. Imprimer du texte sur une imprimante Postscript

L'utilitaire GNU *enscript* est installé sur nymphaea. Pour imprimer un fichier texte « toto » sur une imprimante « imprim », taper la commande :

```
enscript -d imprim toto
```

Pour plus d'options, faire :

```
enscript --help
```

#### 3.6.4. Autres

dos2unix et unix2dos servent à convertir des fichiers texte du mode PC au mode UNIX et inversement.

## 4. Les compilateurs

### 4.1. les langages disponibles

#### 4.1.1. fortran

Le compilateur peut être appelé indifféremment par les commandes f77, f90 ou f95. Exemple très simple d'appel pour un programme source appelé toto.f dont le binaire s'appellera toto:

```
f95 - o toto toto.f
```

Le compilateur supporte de très nombreuses options décrites dans le WEB de référence

<http://h18009.www1.hp.com/fortran/docs/unix-um/dfumperf.htm>

#### 4.1.2. c

Le compilateur langage C est accessible par la commande UNIX « cc » et la documentation complète peut être visualisée par la commande :

```
man cc
```

#### 4.1.3. c++

Le compilateur langage C++ est accessible par la commande UNIX « cxx » et la documentation complète peut être visualisée par la commande :

```
man cxx
```

#### 4.1.4. java

L'interpréteur JAVA est appelé par la commande UNIX « java ». Le JDK est présent sur la machine. La documentation existe sous forme de man ou de pages WEB.

#### 4.1.5. compilateurs du GNU (gcc, g++, g77)

Vous trouverez aussi les compilateurs du GNU dans :

```
/home/nympea/services/logiciels/compilateurs/gcc-3.3/bin
```

## 4.2. Les bibliothèques de sous-programmes disponibles

### 4.2.1. cxml

Ces deux bibliothèques correspondent aux bibliothèques BLAS, LAPACK et scaLAPACK du domaine public. Elles sont optimisées pour l'architecture de nymphaea.

Pour compiler un programme qui appelle des routines de cxml, il faut ajouter un paramètre à la commande de compilation.

```
f90 my_prog.f -lcxml
```

### 4.2.2. cxmnp

Cxmnp est la version parallélisée de cxml. Certaines routines exploitent plusieurs processeurs en parallèle.

Mode d'emploi :

```
f90 my_prog.f -lcxmnp
```

### 4.2.3. scalapack (ou appellation CS chez HP-COMPAQ)

Scalapack est une bibliothèque du domaine public contenant la version parallélisée de certaines routines de Lapack. Elle est optimisée pour l'architecture de nymphaea.

Appel :

```
f90 my_prog.f -lscalapack
```

### 4.2.4. nag

La bibliothèque NAG FORTRAN SMP est installée sur nymphaea. NAG regroupe un certain nombre de sous-programme mathématique optimisés sur une architecture parallèle à mémoire partagée.

### 4.2.5. netcdf

NETworkCommonDataForm est un format de fichier autodescriptif, associé à une librairie (FORTRAN ou C) d'entrée/sortie d'informations de ces fichiers. Il peut contenir des variables, des caractéristiques sur ces variables (nom, dimensions...), des axes (noms, tailles...). Ce format permet entre autres de structurer et de simplifier l'écriture et la lecture des valeurs de variables stockées dans un de ces fichiers.

---

La version disponible sur nympha est la 3.5.0 de “ University Corporation for Atmospheric Research/Unidata “

Pour utiliser la bibliothèque, spécifier l’option “**-lnetcdf**” dans la ligne de commande

Les différents utilitaires associés (NetcdfOperator version 2.2.0) se trouvent dans : /home/nympha/services/bibli/netcdf/bin

## 5. Exécuter un programme

### 5.1. En interactif

Tout programme peut être lancé en « interactif », comme une commande UNIX. Cependant, ce type d'exécution ne permet pas un usage optimal de la puissance des processeurs, car il est soumis à une mutualisation liée au nombre d'utilisateurs. Une limite est donc fixée pour les programmes en interactif : ils ne peuvent dépasser 1 heure de temps CPU.

### 5.2. En batch avec LSF

#### 5.2.1. principes du batch

Le progiciel LSF de Platform Computing est utilisé pour la soumission de batchs sur la machine. Le "batch", c'est le lancement de programme sur la machine sans interaction avec l'utilisateur. L'ensemble du programme doit être autonome : les données d'entrée et de sortie doivent être en fichiers. L'intérêt du batch réside dans l'usage optimal des ressources de la machine, y compris en heures non-ouvrables. LSF permet de gérer les files d'attente, les priorités en utilisant au mieux les ressources disponibles et chaque programme se voit attribuer de manière exclusive le nombre de processeur dont il a besoin pendant toute la durée de son exécution..

#### 5.2.2. les files batchs sur nymphaea

##### **Les différentes queues batch sont:**

- sequentiel (queue implicite)
- long
- parallel4
- parallel16 ( parallelt)
- bigmem (bigfast)
- bigmems

##### **Jobs nécessitants moins de 8Go de mémoire**

*job sequentiel < 48h CPU*

`bsub -N -o output programme`

*job sequentiel < 120h CPU*

`bsub -N -o output -q long programme`

*(attention, cette queue n'est prioritaire que la nuit et les week-end et peut donc être suspendue par les autres jobs dans la journée)*

**job openmp (exemple pour 4 CPU)**

```
setenv OMP_NUM_THREADS 4
```

```
bsub -N -o output -n 4 -R "span[ptile=4]" -q parallel4 programme
```

**job MPI de 2 a 4 CPU (exemple pour 4 CPU)**

```
bsub -N -o output -n 4 - q parallel4 prun programme
```

**job MPI de 5 a 16 CPU (exemple pour 12 CPU)**

```
bsub -N -o output -n 12 - q parallel16 prun programme
```

**Jobs nécessitants plus de 8Go de mémoire (un seul nœud disponible)****job sequentiel**

```
bsub -N -o output -q bigmems programme
```

**job openmp (exemple pour 4 CPU)**

```
setenv OMP_NUM_THREADS 4
```

```
bsub -N -o output -n 4 -q bigmem programme
```

**job MPI de 2 a 4 CPU (exemple pour 4 CPU)**

```
bsub -N -o output -n 4 - q bigmem prun programme
```

**Queues batch de test ( queues prioritaires)**

**parallelt** (*jobs MPI de 5 à 16 process < 5 minutes elapse*)

**bigfast** (*gros jobs (> 8Go) MPI ou OPENMP < 60 minutes elapse*)

**Paramètres généraux**

- le temps UC max est de 48h par UC (sauf pour la queue "long")
- le nœud interactif est mis a disposition du batch la nuit et les week-end

**Les commandes utiles à connaître de LSF sont : bsub, bjobs et bkill**

- **BSUB** : pour soumettre un job (sous la forme d'un shell\_script) dans une queue batch

```
ex : bsub -N -o fichier_sortie [ -n nombre ] [ -R "reservation d'une ressource" ]  
      -q queue_batch [ prun ] shell_script
```

options : **-N** : pour recevoir un mail à la fin du job

**-o** : précise le fichier de sortie (équivalent de l'écran en interactif)

**-n** : précise le nombre de process (pour les jobs **MPI**)  
ou de processeurs (pour les jobs **OPENMP**)

à ne pas utiliser pour les jobs **séquentiels**

**-R** : (pour les jobs **OPENMP**) ex: -R "span[ptile=4]"

**-q** : précise la queue batch

**prun** : à utiliser uniquement (et obligatoirement ) pour les jobs **MPI**

- **BJOBS** : pour surveiller l'état des queues batchs

ex : **bjobs -u all** (pour tous les users)  
**bjobs -u user** (pour un user particulier)

- **BKILL** : pour tuer un job dans une queue batch

ex : **kill -r jobid** (le jobid étant donné par la commande **bjobs**)

### 5.2.3. Lancer des commandes LSF à partir du WEB INTRANET

Il est désormais possible de lancer des jobs LSF sur nympha à partir du serveur WEB INTRANET « ilestudy ». Voici comment procéder :

## 5.3. Visualiser les jobs en cours sur les machines

### 5.3.1. Jobs actifs

Taper la commande « rinfo »

```
nympha0>85% rinfo
MACHINE CONFIGURATION
nympha ifremer
```

PARTITION	CPUS	STATUS	TIME	TIMELIMIT	NODES
root	38				nympha[0-8] nympheams
normal	25/36	running	2:03:50:35		nympha[0-8]

RESOURCE	CPUS	STATUS	TIME	USERNAME	NODES
normal.28410	1	allocated	1:23:51:47	balessan	nympha5
normal.28413	1	allocated	1:23:31:17	balessan	nympha8
normal.28438	1	allocated	1:21:21:07	baraille	nympha7
normal.28443	1	allocated	1:21:12:27	thuck	nympha7
normal.28484	4	allocated	1:03:31:24	vgarnier	nympha1
normal.28489	1	allocated	1:03:10:14	pcugier	nympha5
normal.28536	4	allocated	06:17:00	ardhuin	nympha3
normal.28541	1	allocated	03:56:59	blevier	nympha7
normal.28548	1	allocated	01:52:19	oartzel	nympha5
normal.28552	1	allocated	01:31:39	priou	nympha5
normal.28554	4	allocated	01:26:49	ardhuin	nympha4
normal.28563	1	allocated	30:59	avanhout	nympha8
normal.28567	4	allocated	02:09	bpre	nympha2

JOB	CPUS	STATUS	TIME	USERNAME	NODES
normal.28915	1	running	1:23:51:46	balessan	nympha5
normal.28918	1	running	1:23:31:17	balessan	nympha8
normal.28944	1	running	1:21:21:06	baraille	nympha7
normal.28949	1	running	1:21:12:26	thuck	nympha7
normal.28992	1	running	1:03:31:23	vgarnier	nympha1
normal.28997	1	running	1:03:10:11	pcugier	nympha5
normal.29045	1	running	06:16:59	ardhuin	nympha3
normal.29050	1	running	03:56:58	blevier	nympha7
normal.29057	1	running	01:52:19	oartzel	nympha5
normal.29061	1	running	01:31:38	priou	nympha5
normal.29063	1	running	01:26:49	ardhuin	nympha4
normal.29072	1	running	30:58	avanhout	nympha8
normal.29074	4	running	02:09	bpre	nympha2

nymphhea0>86%

- dans la première partie du tableau (**RESOURCE**) la colonne CPUS donne le nombre de CPUS réservés par le job
- dans la seconde partie du tableau (**JOB**) la colonne CPUS donne le nombre de PROCESS tournant pour le job
  - les lignes en **noir** sont donc des jobs **séquentiels**
  - les lignes en **bleus** sont donc des jobs **OPENMP**
  - la ligne en **rouge** est donc un job **MPI**

### 5.3.2. Jobs par queue batch

Taper la commande « wo » pour obtenir l’affichage suivant :

```

vgarnier pts/9      eos.ifremer.fr  14:19      29  23:49      -csh
pgarreau pts/10     iphis.ifremer.fr 09:56      1:53  2:37      ftp ftp.ifrem
baraille pts/11     pacte-sv.cnes.fr 13:31      1:02  1:08      -csh
esimon   pts/12     neree.shom.fr   14:38      3    31:02     2:37 test_gs
yannf    pts/13     hyades.univ-bres 13:53      45                    -csh
bpre     pts/16     louet.ifremer.fr 15:26      46    2          -sh
mjouan   pts/17     nemesis.ifremer. 15:22      8days 8:35      -csh
bpre     pts/20     louet.ifremer.fr 15:28      56                    -csh
yannf    pts/21     hyades.univ-bres 14:35      45                    -csh
baraille pts/22     spike.cst.cnes.f 15:43      57                    vi fort.9
pgarreau pts/23     iphis.ifremer.fr 14:26      1:34 108:42    2 xlsbatch
mjouan   pts/24     nemesis.ifremer. 14:56      8days 14:08     -csh
bwaelles pts/26     corto.ifremer.fr 11:15      2days 136:11    -csh
lhayina  pts/33     br158-145.ifreme 11:43      2days                    -csh
avanhout pts/35     soledad.ifremer. 15:27      9days                    -csh
rthiery  br166-193. br166-193.ifreme 15:20      -                      -
rthiery  15:20                                          /usr/sbin/get

```

```

*****
*                JOBS BATCH RUNNING                *
*****

```

JOBID	USER	STAT	QUEUE	FROM_HOST	EXEC_HOST	JOB_NAME	SUBMIT_TIME
60031	oarzel	RUN	sequentiel	nymphhea7	nymphhea7	*2s16.bsub	Jul 15 21:21
60112	bthouve	RUN	sequentiel	nymphhea0	nymphhea7	batch_g1	Jul 17 13:32
60119	gcayrol	RUN	sequentiel	nymphhea0	nymphhea7	modessai	Jul 17 14:10
60122	oarzel	RUN	sequentiel	nymphhea0	nymphhea7	*2s19.bsub	Jul 17 14:16
60099	barail	RUN	sequentiel	nymphhea0	nymphhea4	main	Jul 17 10:43
60130	blevier	RUN	sequentiel	nymphhea0	nymphhea4	Biscay	Jul 17 15:28
60131	vgarnie	RUN	parallel4	nymphhea0	nymphhea5	*tch_menor	Jul 17 16:02

▲

## 6. Le debugger – TOTALVIEW

Une fois passée l'étape de la compilation, il se peut que votre programme ne se comporte pas comme attendu : plantage, résultats incohérents. Le debugger « totalview » peut vous permettre de détecter les erreurs, en contrôlant l'exécution par des points d'arrêt, et en permettant à chaque arrêt de visualiser le contenu des variables.

Pour pouvoir l'utiliser, il faut préalablement compiler le programme avec l'option « -g » pour que toutes les variables soient connues de totalview.

```
f95 -g -o toto toto.f
```

Lancement :

```
totalview toto&
```

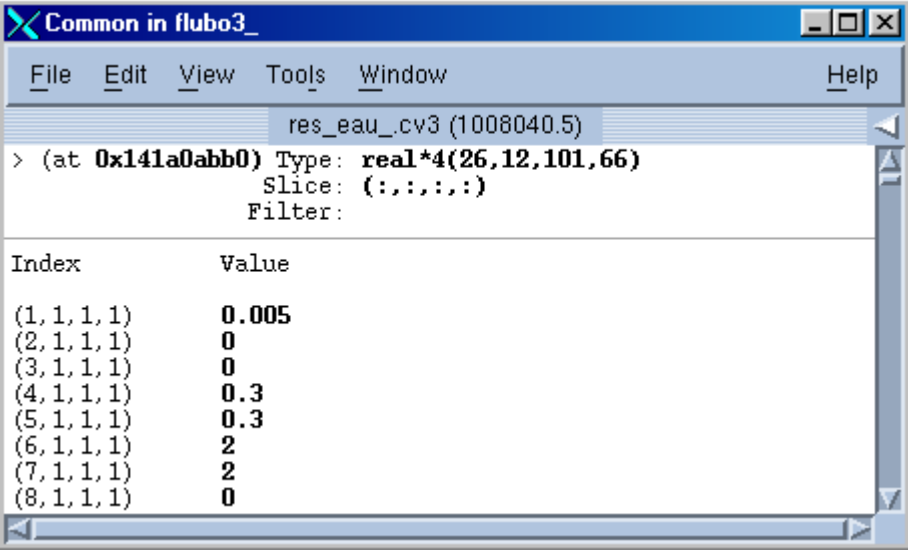
Plusieurs fenêtres s'affichent. Fenêtre principale :

The screenshot displays the elise33prof debugger window. The title bar reads "elise33prof". The menu bar includes File, Edit, View, Group, Process, Thread, Action Point, Tools, Window, and Help. Below the menu bar are buttons for Group, Control, Go, Halt, Next, Step, Out, Run To, Next, Step, and navigation buttons (P-, P+, T-, T+). The status bar indicates "Process 1008040: elise33prof (Stopped)" and "Thread 1008040.5: elise33prof (Stopped) <Trace Trap>".

The main window is divided into several panes:

- Stack Trace:** A table showing the call stack. The top frame is "flubo3\_" at FP=11ffffb810. Below it are "condif3\_" (FP=11ffffba80), "main\$elise33\_" (FP=11ffffbcf0), "main" (FP=11ffffbff0), and "\_start" (FP=11ffffc000).
- Stack Frame:** A table showing the current frame's variables: raph: (Array), soyoin: (Array), soyplu: (Array), soxoin: (Array), soxplu: (Array). Below this are "Common blocks:" with entries like opt\_termfrot\_: (Common), couplage\_: (Common), melange\_: (Common), maill3d\_: (Common), res\_topo\_: (Common), and res\_sau\_: (Common).
- Function flubo3\_ in elise33.f:** A code editor showing the source code. Line 6279 is highlighted in red and marked with a "STOP" icon. The code includes a loop for k from kmin to kmax, and an inner loop for j from jmin to jmax, calculating raph(i) based on epxn and corepx values.
- Thread (1):** A table showing thread information. The first entry is "5 T in flubo3\_".
- Action Points:** A table showing action points. The first entry is "STOP 1 line 6279 at flubo3\_+0x1a40 in".

Exemple d'autre fenêtre visualisant le contenu d'un tableau :



Index	Value
(1, 1, 1, 1)	<b>0.005</b>
(2, 1, 1, 1)	<b>0</b>
(3, 1, 1, 1)	<b>0</b>
(4, 1, 1, 1)	<b>0.3</b>
(5, 1, 1, 1)	<b>0.3</b>
(6, 1, 1, 1)	<b>2</b>
(7, 1, 1, 1)	<b>2</b>
(8, 1, 1, 1)	<b>0</b>

A noter qu'en cas de plantage du programme, totalview signale précisément la ligne de code ayant provoqué le plantage et l'état de toutes les variables à cet instant.

## 7. Optimiser le code

### 7.1. Mesurer le temps passé dans le code

Sur une machine partagée avec d'autres utilisateurs, les temps d'exécution des programmes peuvent varier d'une fois sur l'autre en fonction de la charge de la machine. Les outils de mesure du temps sont donc à considérer en gardant à l'esprit cet état de fait.

#### 7.1.1. L'outil « prof »

L'outil prof permet de lister le temps passé à l'exécution en échantillonnant par subroutine ou même par instruction.

##### 7.1.1.1. Approche globale par subroutine

Compiler avec les options `-p -g3` :

```
f95 -p -g3 -o toto toto.f
```

Exécuter :

```
Toto
```

(ça crée un fichier nommé par défaut `mon.out`)

Lancer prof :

```
Prof toto mon.out > toto.prof
```

Le fichier texte toto.prof contient les informations de temps d'exécution comme suit :

```

profile listing generated Thu Apr  3 10:30:08 2003 with:
  prof elise3D.e mon.out

```

```

-----
* -p[rocedures] using pc-sampling;                               *
* sorted in descending order by total time spent in each procedure; *
* unexecuted procedures excluded                                 *
-----

```

Each sample covers 8.00 byte(s) for 0,0001% of 938,2021 seconds

%time	seconds	cum %	cum sec	procedure (file)
67.3	631,5664	67,3	631,57	flubo3_ (elise3D,f)
8.4	78,5156	75,7	710,08	condif3_ (elise3D,f)
5.0	46,7080	80,7	756,79	hydro3_ (elise3D,f)
4.5	42,6416	85,2	799,43	hydroy_ (elise3D,f)
4.4	41,2109	89,6	840,64	interboites_ (elise3D,f)
2.9	26,7939	92,5	867,44	hydrox_ (elise3D,f)
1.9	18,1846	94,4	885,62	intraboiteseau_ (elise3D,f)
1.9	17,6084	96,3	903,23	t_bilmat_ (elise3D,f)
1.0	9,3203	97,3	912,55	main\$elise3d_ (elise3D,f)
1.0	9,1719	98,2	921,72	s_depeff_ (elise3D,f)

#### 7.1.1.2. Pour échantillonner par instruction

Compiler avec les options `-g3` :

```
f95 -g3 -o toto toto.f
```

Lancer prof :

```
Prof toto mon.out > toto.prof
```

Le fichier texte toto.prof contient les informations de temps d'exécution comme suit. Le temps passé dans chaque ligne de code est évalué.

^Profile listing generated Thu Jun 19 17:14:13 2003 with:  
prof -heavy elise33prof mon.out

```
* -h[eavy] using pc-sampling; *
* sorted in descending order by time spent in each source line; *
* unexecuted lines excluded *
```

Each sample covers 8.00 byte(s) for 0.00033% of 300.3047 seconds

procedure (file)	line	bytes	millisec	%	cum %
flubo3_ (elise33.f)	6574	1144	11174.3	3.72	3.72
flubo3_ (elise33.f)	6573	1564	10597.7	3.53	7.25
flubo3_ (elise33.f)	6344	1924	9960.4	3.32	10.57
flubo3_ (elise33.f)	6609	1444	9714.4	3.23	13.80
flubo3_ (elise33.f)	6575	1184	9490.2	3.16	16.96
flubo3_ (elise33.f)	6612	1424	8862.8	2.95	19.91
flubo3_ (elise33.f)	6486	1140	8496.6	2.83	22.74
flubo3_ (elise33.f)	6611	1168	8418.0	2.80	25.55
flubo3_ (elise33.f)	6614	988	6586.9	2.19	27.74
flubo3_ (elise33.f)	6345	1012	5921.4	1.97	29.71
flubo3_ (elise33.f)	6487	984	5495.6	1.83	31.54
flubo3_ (elise33.f)	6410	584	5071.8	1.69	33.23
flubo3_ (elise33.f)	6346	1324	4681.6	1.56	34.79
flubo3_ (elise33.f)	6610	716	3979.5	1.33	36.11
flubo3_ (elise33.f)	6411	760	3914.6	1.30	37.42
flubo3_ (elise33.f)	6490	1008	3914.1	1.30	38.72
flubo3_ (elise33.f)	6408	520	3912.1	1.30	40.02
flubo3_ (elise33.f)	6347	644	3773.4	1.26	41.28
t_bilmat_ (elise33.f)	7706	508	3710.4	1.24	42.52

### 7.1.2. Le script cpuProfiler

Pour rendre meilleure la lisibilité des temps passés échantillonnés par instruction, il est possible d'installer le script « cpuProfiler ».

Installation : il faut copier le script à partir de /home/nymphaea/services/bin dans le répertoire où se trouve le source du programme à analyser.

Usage : se positionner dans ce même répertoire et lancer la commande

```
./cpuProfiler exemple.f lignedebut lignefin
```

exemple.f est le programme source

lignedebut est le numéro de ligne où démarre l'analyse

lignefin est le numéro de ligne de fin de l'analyse

Résultat produit :

Un répertoire « resultatProfiling » est créé dans le répertoire courant contenant les fichiers « exemple.part.profile » et « exemple.profile ». Ce sont les fichiers source modifiés comme indiqué sur l'image qui suit

```

0      c
0      c      couches intermediaires
0      c      .....
435145196          do 7 k=kmi(i,j)+1,kma(i,j)-1
295977272          aa(k)=a1(k)+fa(k)
271110468          bb(k)=b1(k)+fb(k)
273832940          cc(k)=c1(k)+fc(k)
805185552          dd(k)=(ep(k,i,j)*(cv3(iv,k,i,j)+ddd(iv,k,i,j)*dt)
437807702          $          +(-flux(iv,k,i,j)+flux(iv,k,i-1,j)
381747867          $          -fluy(iv,k,i,j)+fluy(iv,k,i,j-1)
156132348          $          )/(dxdy(i,j)+eps))
57489060          $          *dt3i
0      7      continue
n      r

```

La colonne de gauche a été ajoutée : elle contient le nombre de cycles machine pour l'instruction correspondant à cette ligne. A noter que ce résultat est obtenu sans aucune option d'optimisation sur le compilateur. Par ailleurs, il faut éviter les lignes blanches (sans même le « c » du commentaire) dans le code pour prévenir les décalages de ligne.

#### 7.1.3. La commande FORTRAN « cpu\_time »

Cette instruction FORTRAN permet de faire des mesures sur une section de code seulement. Le résultat subit la relativité mentionnée au §7.1.

Exemple :

```

time=0.
....
call CPU_TIME (time_begin)
.....
call CPU_TIME(time_end)
time=time+time_end-time_begin
....

```

A noter que l'usage de ces instructions retarde sensiblement l'exécution du programme.

#### 7.1.4. La commande UNIX « time »

Il suffit de lancer (pour le programme toto):

```
time toto
```

A la fin de l'exécution une ligne s'affichera comme suit :

```
353.39u 4.07s 5:57 99% 0+1257k 34+36788io 0pf+0w
```

Le nombre 353.39u indique que la partie utilisateur du programme a consommé 353,39 secondes. Ensuite le nombre 4.07s indique 4,07 secondes de consommation pour la partie système du programme.

### 7.1.5. Les compte-rendus d'exécution de LSF

A chaque exécution, LSF renvoie un message à l'utilisateur qui contient des informations de temps passé. Le temps « ellapsed » est la différence entre la date de lancement du job et sa date de fin.

Le temps CPU est la consommation mesurée de CPU. La différence entre ces 2 temps est notamment due aux entrées-sorties, qui peuvent être ralenties par la concurrence avec d'autres travaux de la machine.

```

Sujet: Job 58352: <elise3D17> Done
De : LSF <lsfadmin@ifremer.fr>
Date : 10/06/03 18:29
A : cotty@ifremer.fr

```

```

Job <elise3D17> was submitted from host <nymphaea0> by user <cotty>.
Job was executed on host(s) <nymphaea4>, in queue <sequentiel>, as user <cotty>.
</home1/beniguet/tmsiric/cotty> was used as the home directory.
</export/home/bench/sources/sam/elise/batch> was used as the working directory.

```

```
Started at Tue Jun 10 18:23:03 2003
```

```
Results reported at Tue Jun 10 18:29:27 2003
```

différence fin-début =  
temps ellapsed

Your job looked like:

```

-----
# LSBATCH: User input
elise3D17
-----

```

Successfully completed.

Resource usage summary:

```

CPU time : 353.00 sec.
Max Memory : 126 MB
Max Processes : 2

```

temps CPU

## 7.2. La parallélisation

Cette technique est basée sur le principe que, dans un programme, des instructions qui se suivent ne sont pas nécessairement dépendantes l'une de l'autre, et qu'on peut donc envisager de les exécuter simultanément sur des processeurs différents pour gagner du temps. 2 types de parallélisation sont possibles : l'une fait travailler des processeurs partageant la même mémoire (directives openMP) et l'autre peut faire travailler des processeurs ne partageant pas la même mémoire (directives MPI).

OpenMP est assez facile à mettre en œuvre, mais le gain en temps ne peut dépasser un facteur 4 sur « nympha ».

MPI nécessite une refonte totale de l'algorithme pour qu'il puisse admettre d'être divisé en plusieurs parties pouvant s'exécuter sur différents processeurs.

### 7.2.1. Parallélisation OpenMP

#### 7.2.1.1. Préparer le code

La parallélisation est une fonctionnalité incluse dans le compilateur. On introduit à certains endroits dans le code source des « directives » demandant à celui-ci de répartir la charge d'exécution sur plusieurs processeurs. Si le code concerné par la directive s'y prête, on peut obtenir (théoriquement) des temps d'exécution divisés par le nombre de processeurs travaillant en parallèle.

Voir le WEB :

[http://w3.ifremer.fr/intrarc/Assistance/Guides/nympha/Formation-Compaq-Web\\_fichiers/v3\\_document.htm](http://w3.ifremer.fr/intrarc/Assistance/Guides/nympha/Formation-Compaq-Web_fichiers/v3_document.htm)

à partir du transparent 284.

#### 7.2.1.2. Compilation

L'option du compilateur demandant la mise en œuvre de la parallélisation est « -omp ». Exemple  
f95 -omp -o toto toto.f

#### 7.2.1.3. Execution

Pour pouvoir utiliser plusieurs processeurs, il faut positionner une variable d'environnement comme suit :  
setenv OMP\_NUM\_THREADS=4

Pour lancer le programme « toto » en batch, la commande BSUB de LSF requiert quelques options comme suit :

```
bsub -N -n 4 -R "span[ptile=4]" -q parallel4 toto
```

#### 7.2.1.4. Quelques grands principes pour la parallélisation openMP

##### 7.2.1.4.1 Variables « private »

Dans la directive de parallélisation openMP, on indique au compilateur les variables de la section à paralléliser qu'il doit considérer comme « private », c'est-à-dire que les différents processeurs impliqués dans

l'exécution n'ont pas besoin de se tenir au courant des modifications concernant ces variables.

Il est donc utile de veiller à ce qu'il y ait dans la section le moins possible d'appels à des variables autres que « private ».

#### 7.2.1.4..2 Cas des boucles DO

Il est souvent intéressant de paralléliser des boucles « DO ». Pour obtenir de bonnes performances, il faut rendre indépendante la boucle de rang « i » par rapport à celle de rang « i-1 » ou les précédentes. Plusieurs artifices de programmation peuvent y aider. Voir les transparents du cours :

#### 7.2.1.4..3 Instructions IF

Les instructions « IF » sont malvenues dans une section à paralléliser. Dans pas mal de cas, elles peuvent être remplacées par des calculs faits systématiquement et des instructions de type « min », « max », « dim », ...

Sur un processeur, l'exécution est plus longue, mais sur 4, ça devient globalement plus intéressant.

### 7.2.2. La parallélisation MPI

#### 7.2.2.1. Principe

Le programme est lancé simultanément en plusieurs exemplaires sur des processeurs différents. Grâce à des instructions spécifiques MPI, les différents exemplaires (process) du programme s'identifient (donc se différencient) puis peuvent s'échanger des données par « messages ». Sans modification du code et sans échange de message le résultat est que chaque process exécute strictement la même chose et les performances sont les mêmes qu'en monoprocesseur.

Toute la technique MPI consiste à modifier le programme de telle sorte que chaque process ne se charge que d'une partie du travail, et que le résultat final, issu de tous les process, soit rassemblé à la fin.

#### 7.2.2.2. Compilation

L'appel à MPI est matérialisé par l'option `-lmpi`

```
f90 -lmpi -o myprog myprog.f
```

#### 7.2.2.3. Exécution

Dans LSF, un programme MPI se lance de la manière suivante :

```
bsub -N -o myprog.out -n 4 -q parallel4 prun myprog
```

---

-n détermine le nombre de process choisi  
-q introduit la classe de jobs choisie (nécessairement en rapport avec le nombre de process)  
la directive prun est nécessaire en MPI

Si le programme est très bien parallélisé, on peut choisir le nombre de process au dernier moment, c'est-à-dire lors du lancement de l'exécution.

## 8. Les progiciels installés

### 8.1. fluent

Le binaire du progiciel FLUENT est installé sur nymphéa, mais aucun ticket de licence n'y existe. Il faut donc, dans l'environnement utilisateur (.cshrc), indiquer au logiciel le serveur externe où il pourra trouver des tickets de licence.

Insérer une commande du genre :

```
setenv LM_LICENSE_FILE 7241@serveurexterne
```

### 8.2. gcg – Wisconsin package

Ce logiciel inclut la plupart des fonctionnalités nécessaires au biologiste moléculaire pour analyser ses séquences issues du travail de laboratoire, les confronter avec les principales bases publiques de séquence (EMBL, SWISSPROT, ...). Deux modes d'accès sont disponibles : l'un en mode commande UNIX, l'autre en graphique via une interface X11.

Documentation : <http://www.infobiogen.fr/doc/GCGdoc/Doc.html>

### 8.3. abaqus

La suite logicielle ABAQUS, spécialisée en calcul de structures par la méthode des éléments finis, est installée sur nymphéa. Aucun pool de licence n'y est cependant disponible, si bien que chaque utilisateur doit préciser dans son environnement le nom de son serveur de licence externe par une commande du genre :

```
setenv LM_LICENSE_FILE 7241@serveurexterne
```

### 8.4. Phylowin

Phylowin est un logiciel à interface graphique couleur pour les applications de phylogénie moléculaire.

### 8.5. Seaview

Seaview est un éditeur graphique pour étudier les alignements de séquences multiples.

## 8.6. Matlab

Matlab est un logiciel généraliste du monde des maths appliquées. Il est installé sur nymphaea, mais une seule licence réservée à une application opérationnelle est présente.

Matlab n'existant pas en version parallélisée, le seul intérêt de le faire tourner sur nymphaea réside dans la capacité à exécuter des scripts en langage propriétaire matlab.