

Programació i computació paral·lela

Ivan Roderó Castro
Francesc Guim Bernat

PID_00191902



Els textos i imatges publicats en aquesta obra estan subjectes –llevat que s'indiqui el contrari– a una llicència de Reconeixement-NoComercial-SenseObraDerivada (BY-NC-ND) v.3.0 Espanya de Creative Commons. Podeu copiar-los, distribuir-los i transmetre'ls públicament sempre que en citeu l'autor i la font (FUOC. Fundació per a la Universitat Oberta de Catalunya), no en feu un ús comercial i no en feu obra derivada. La llicència completa es pot consultar a <http://creativecommons.org/licenses/by-nc-nd/3.0/es/legalcode.ca>

Índex

Introducció	5
Objectius	6
1. Models de programació per a memòria compartida	7
1.1. Programació amb processos/fluxos	7
1.2. OpenMP	12
1.2.1. Directives OpenMP	14
1.2.2. Creació de fluxos	14
1.2.3. Clàusules de compartició de variables	15
1.2.4. Clàusules de divisió del treball	16
1.2.5. Directives de sincronització	19
1.2.6. Funcions i variables	21
1.2.7. Entorn de compilació i execució	21
2. Models de programació gràfica	24
2.1. CUDA	24
2.1.1. Arquitectura compatible amb CUDA	24
2.1.2. Entorn de programació	26
2.1.3. Model de memòria	28
2.1.4. Definició de nuclis	32
2.1.5. Organització de fluxos	33
2.2. OpenCL	36
2.2.1. Model de paral·lisme per a dades	36
2.2.2. Arquitectura conceptual	38
2.2.3. Model de memòria	39
2.2.4. Gestió de nuclis i de dispositius	40
3. Models de programació per memòria distribuïda	44
3.1. MPI	44
3.1.1. Comunicadors	46
3.1.2. Comunicacions punt a punt	47
3.1.3. Comunicacions col·lectives	51
3.1.4. Compilació i execució	56
3.2. Llenguatges PGAS	57
3.2.1. UPC	57
3.2.2. Co-Array Fortran	59
3.2.3. Titanium	60
4. Esquemes algorítmics paral·lels	62
4.1. Paral·lisme de dades	62

4.2. Partició de dades	64
4.3. Esquemes paral·lels en arbre	65
4.4. Computació en <i>pipeline</i>	67
4.5. Esquema mestre-esclau	70
4.6. Computació síncrona	71
5. Entorns d'execució d'aplicacions paral·leles.....	73
5.1. Sistemes de gestió de cues	74
5.1.1. PBS (<i>portable batch system</i>)	74
5.2. Planificació	80
5.2.1. Polítiques de planificació basades en <i>backfilling</i>	83
Bibliografia.....	87

Introducció

En aquest mòdul didàctic estudiarem els models principals de programació d'aplicacions paral·leles utilitzats en la computació d'altres prestacions, els esquemes bàsics de programació d'aplicacions paral·leles i els entorns utilitzats per a l'execució d'aquest tipus d'aplicacions en el context de la computació d'altres prestacions.

En els primers apartats ens centrarem en els models de programació tant per a memòria compartida com distribuïda. En relació amb els models de memòria compartida estudiarem principalment OpenMP i també farem èmfasi en la programació gràfica, la qual té cada cop més presència en la computació d'altres prestacions. A continuació estudiarem els models principals per a memòria distribuïda, tant de pas de missatges (MPI) com els que proporcionen abstraccions d'espais d'adreces compartits utilitzant memòria distribuïda.

Un cop presentats els models de programació més rellevants per a la computació d'altres prestacions estudiarem algunes de les tècniques més populars que permeten el desenvolupament d'aplicacions paral·leles, tant per al seu disseny inicial com per a convertir aplicacions seqüencials en paral·leles. En concret, estudiarem tant tècniques per a l'obtenció de paral·lelisme (per exemple, a partir del paral·lelisme de dades) com diversos patrons d'algoritmes paral·lels.

Finalment estudiarem l'entorn d'execució d'aplicacions paral·leles per a la computació d'altres prestacions. Estudiarem els components de programari utilitzats en els sistemes paral·lels orientats a les altres prestacions i també polítiques de planificació que permeten optimitzar l'execució de les aplicacions paral·leles i la utilització dels sistemes d'altres prestacions.

Objectius

Els materials didàctics d'aquest mòdul contenen les eines necessàries per a assolir els objectius següents:

- 1.** Conèixer els models de programació per a memòria compartida i saber programar aplicacions paral·leles amb OpenMP.
- 2.** Aprendre els conceptes fonamentals per a programar dispositius GPU amb els models de programació per a computació gràfica CUDA i OpenCL.
- 3.** Aprendre els conceptes fonamentals per a programar aplicacions paral·leles amb MPI i conèixer models de programació per a memòria distribuïda basats en espais d'adreces compartits com ara PGAS.
- 4.** Aprendre les tècniques i patrons bàsics per al disseny d'algoritmes paral·lels i ser capaç de desenvolupar aplicacions basades amb aquestes tècniques i els models de programació estudiats en aquest mòdul didàctic.
- 5.** Conèixer el funcionament i els components que formen els sistemes de gestió d'aplicacions en sistemes paral·lels per a la computació d'altres prestacions.
- 6.** Aprendre els conceptes fonamentals de les polítiques de planificació de treballs en entorns d'altres prestacions.

1. Models de programació per a memòria compartida

Tal com vam veure en el mòdul “Introducció a la computació d’altes prestacions” d’aquesta assignatura, els sistemes multiprocessador o multinucli de memòria compartida són aquells en els quals qualsevol dels processadors o nuclis poden accedir a qualsevol espai de la memòria. Per tant, hi ha un únic espai d’adreces i, per tant, a cada ubicació de memòria té una única adreça dins d’un rang d’adreces. Les alternatives principals per a la programació de sistemes de memòria compartida són:

- Utilitzar processos convencionals oferts pel sistema operatiu.
- Utilitzar fluxos, com per exemple *Pthreads*.
- Utilitzar un nou llenguatge de programació. Un exemple és el llenguatge Ada, tot i que no és una solució gaire estesa.
- Utilitzar una biblioteca amb un llenguatge de programació existent.
- Modificar la sintaxi d’un llenguatge de programació seqüencial existent per tal de crear un llenguatge de programació paral·lela. Un exemple és UPC¹.
- Utilitzar un llenguatge de programació seqüencial existent i complementar-lo amb directives de compilació per tal d’especificar paral·lisme. Un exemple és OpenMP.

⁽¹⁾De l’anglès *unified parallel C*.

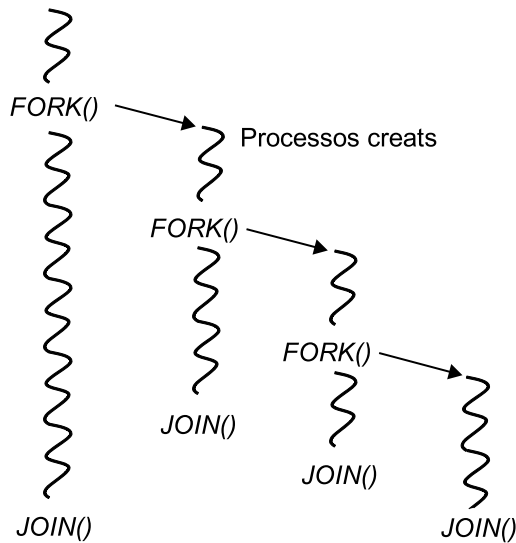
En les seccions següents estudiarem les característiques bàsiques de la programació basada en processadors/fluxos i els possibles problemes associats. A continuació estudiarem OpenMP com a model de programació per a memòria compartida adequat per a la computació d’altes prestacions.

1.1. Programació amb processos/fluxos

Un dels possibles mecanismes per a implementar paral·lisme amb memòria compartida és l’ús de processos convencionals utilitzant el model *fork-join* tal com mostra la figura 1.

Figura 1. Model *fork-join* utilitzat per a implementar paral·lelisme amb processos

Programa principal



Es tracta d'un model SPMD² en què, generalment, el codi corresponent al procés pare és diferent del que executen els processos creats que s'acostumen a anomenar *esclaus*. A més, el procés pare és el que s'encarrega d'esperar la finalització dels processos esclaus (*JOIN* en la figura 1). L'estructura del codi que implementa el model *fork-join* es mostra a continuació.

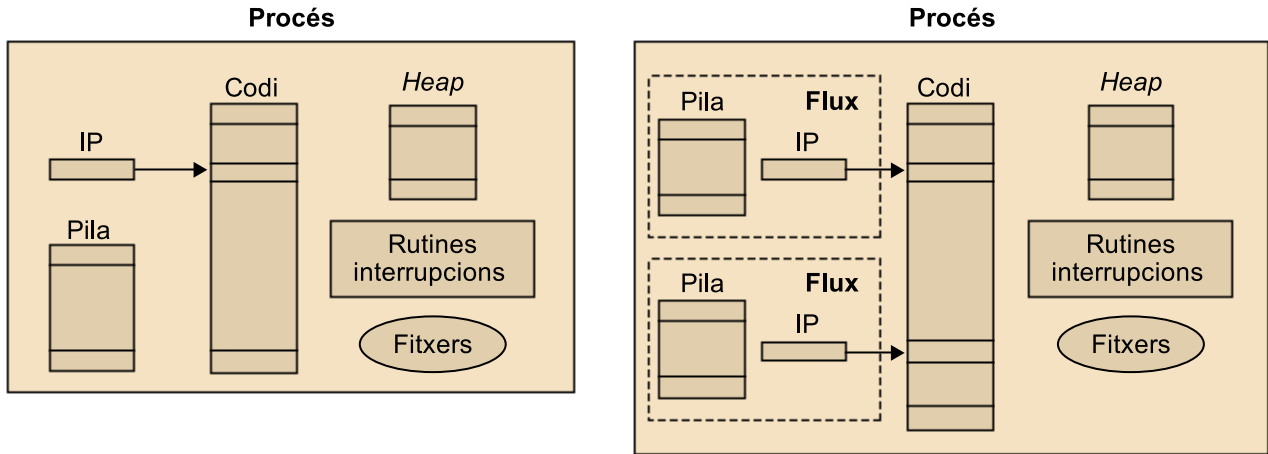
⁽²⁾De l'anglès *single-program multiple-data*

```
pid = fork();
if (pid==0) {
    // codi a executar per un procés esclau
}
else{
    // codi a executar pel procés pare
}
if (pid==0) exit;
else wait(0);
...
```

El problema principal que fa que aquesta solució no sigui efectiva és la gran sobrecàrrega associada a la creació i gestió dels processos. En canvi, el model introduït anteriorment s'aplica majoritàriament en la programació per a memòria compartida.

L'alternativa més evident a utilitzar processos tradicionals és utilitzar fluxos. Mentre que els processos són molt pesats, ja que són codis completament independents amb les seves pròpies variables i estructures de memòria, els fluxos comparteixen un espai global de memòria entre les diferents rutines tal com mostra la figura 2.

Figura 2. Diferències entre processos i fluxos

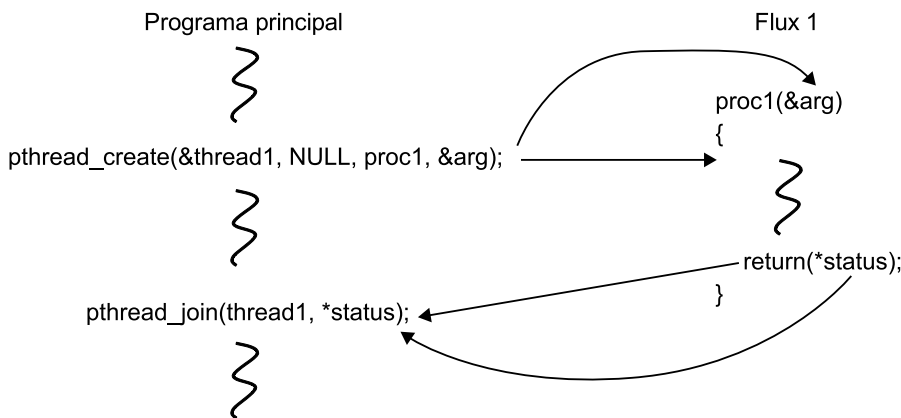


Esquerra: procés convencional. Dreta: procés amb dos fluxos

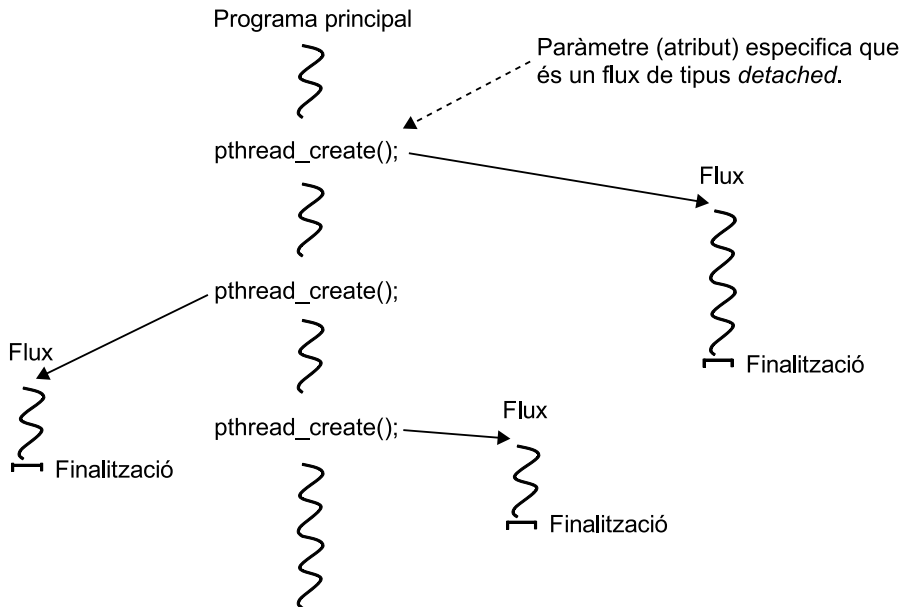
En un altre mòdul de l'assignatura vam presentar la biblioteca de fluxos *Pthreads*, que compleix els estàndards POSIX. També vam veure un codi que exemplificava el model *fork-join*, que es pot representar tal com mostra la figura 3.

Vegeu també
 La biblioteca de fluxos *Pthreads* s'ha vist en el mòdul "Introducció a la computació d'altres prestacions" d'aquesta assignatura.

Figura 3. Execució d'un flux mitjançant *Pthreads*



La biblioteca de *Pthreads* també permet utilitzar fluxos que no necessiten que el programa principal hagi d'esperar la seva finalització (mitjançant *join*). Aquests tipus de fluxos s'anomenen *detached* i quan finalitzen són destruïts i els seus recursos alliberats. La figura 4 mostra un esquema de funcionament de fluxos de tipus *detached*.

Figura 4. Execució de fluxos de tipus *detached*

Quan s'utilitzen fluxos cal tenir en compte que l'espai de memòria és compartit i poden haver-hi problemes tant de consistència com de rendiment.

Podem dir que una rutina és segura si l'execució de múltiples instàncies de la rutina en diferents fluxos produeixen resultats correctes. Així doncs, caldrà assegurar-se que les rutines són segures quan accedeixen a dades compartides. En l'exemple de la figura 5 es poden veure dos fluxos que incrementen el valor de la variable compartida x . Per a dur a terme la instrucció cal llegir primer la variable x , després calcular $x+1$ i finalment escriure el resultat a x . Com que tots dos fluxos llegeixen el valor original de x a la vegada, el valor final de x només reflectirà un únic increment.

Figura 5. Exemple d'accés a dades compartides mitjançant fluxos

Instrucció	Flux 1	Flux 2
$x = x+1;$	llegir x	llegir x
	calcular $x+1$	calcular $x+1$
	escriure a x	escriure a x

Temps ↓

Un mecanisme per a assegurar que només un flux accedeix a un recurs específic al mateix temps és definir seccions del codi que només es puguin executar un sol cop a la vegada. Aquest tipus de seccions de codi s'anomenen *seccions crítiques* i el mecanisme associat a la utilització d'aquestes es coneix com a *exclusió mútua*.

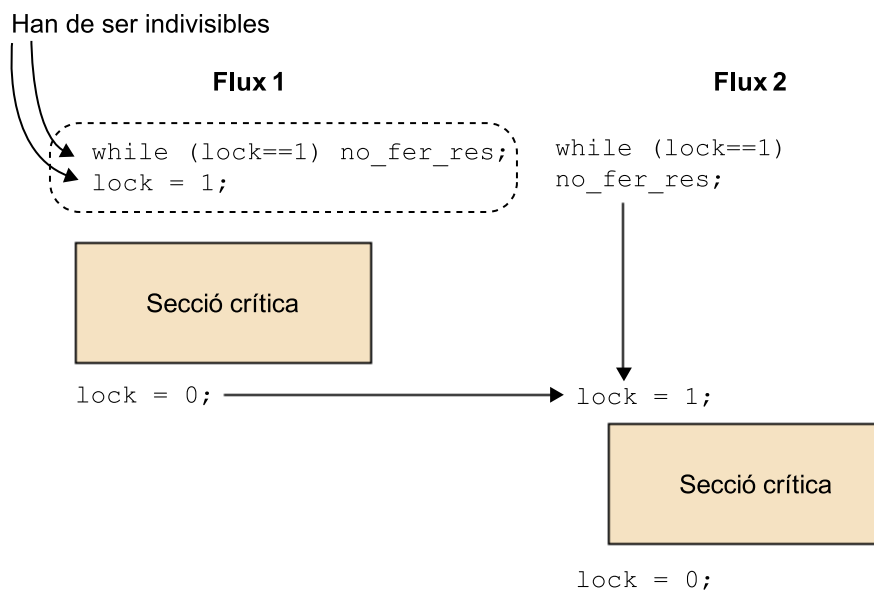
El mecanisme més senzill per a assegurar exclusió mútua en seccions crítiques és la utilització de *locks*. Un *lock* és una variable d'un bit que quan és 1 indica que un flux ha entrat en la secció crítica i quan és 0 indica que cap flux no

hi ha entrat. Una possible implementació d'exclusió mútua és mitjançant una espera activa (*busy waiting*) que consisteix en els passos següents, que també il·lustra la figura 6:

- Esperar fins que el *lock* indiqui que no hi ha cap flux en la secció crítica.
- Bloquejar la secció crítica canviant el valor del *lock* a 1.
- Executar el codi de la secció crítica.
- Desbloquejar la secció crítica i canviar el valor del *lock* un altre cop a 0.

L'inconvenient principal de l'espera activa és que consumeix molts cicles de CPU que no són realment útils.

Figura 6. Esquema d'espera activa per a accedir a una secció crítica des de dos fluxos



La biblioteca de *Pthreads* implementa els *locks* directament amb variables *lock* d'exclusió mútua que s'anomenen variables de tipus *mutex*. Les seccions crítiques es defineixen com es mostra a continuació:

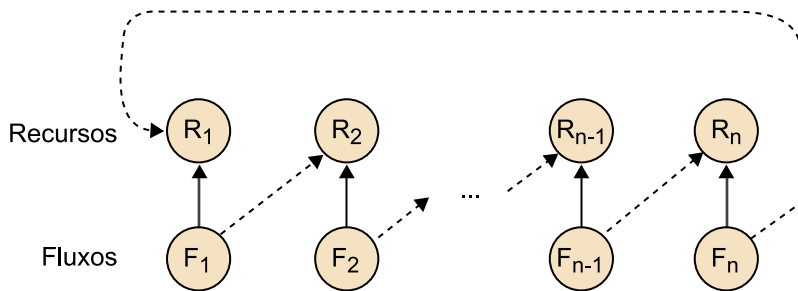
```
pthread_mutex_lock(&mutex1);
    Secció crítica
pthread_mutex_unlock(&mutex1);
```

Si un flux arriba a la crida `pthread_mutex_lock` quedarà bloquejat fins que el *lock* indicat estigui obert. Si hi ha diversos fluxos bloquejats en aquest lloc, quan el *lock* s'obri el sistema determinarà el flux que pot accedir a la secció crítica. Cal tenir en compte que només el flux que tanca un *lock* el pot tornar a obrir.

Un dels problemes principals de l'exclusió mútua és el conegut com a *inter-bloqueig*³. L'interbloqueig de dos fluxos es produeix quan cadascun dels dos necessita el recurs que bloqueja l'altre flux. Aquest fenomen també es pot produir de manera circular involucrant múltiples fluxos, com il·lustra la figura 7. Aquest fenomen es coneix com a *abraçada mortal*.

⁽³⁾En anglès, *deadlock*.

Figura 7. Abraçada mortal involucrant n fluxos



Per tal de prevenir els interbloqueigs, la biblioteca de *Pthreads* ofereix l'operació `pthread_mutex_trylock()`, que permet consultar si un *lock* està obert o tancat abans de bloquejar el flux. Aquesta operació tanca el *lock* i retorna 0 si prèviament estava obert, o bé retorna `EBUSY` si està tancat.

Activitat

Busqueu per la vostra banda altres mecanismes per a prevenir interbloqueigs, com ara semàfors o monitors.

1.2. OpenMP

És possible desenvolupar programes paral·lels per a memòria compartida amb *Pthreads*, però caldrà gestionar els fluxos i tenir en compte problemes de sincronització com els que hem vist anteriorment. Òbviament, això no afavoreix la productivitat del programador i fa que es necessitin abstraccions de més alt nivell que facilitin la programació i l'execució eficient d'aplicacions paral·leles com és la que proporciona OpenMP.

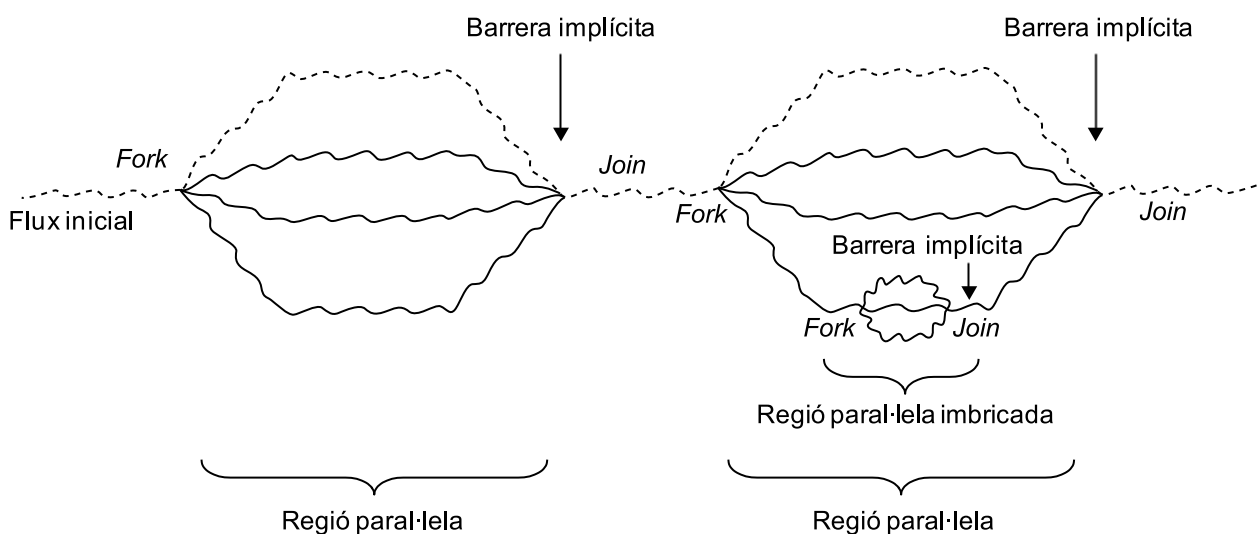
OpenMP consisteix en una interfície de programació que estén llenguatges de programació, com ara C/C++ i Fortran, mitjançant l'ús de directives o "pragmes". OpenMP és l'estàndard actual per a la programació de sistemes de memòria compartida, com ara sistemes multinucli i computadors d'altres prestacions basats en multiprocessadors amb memòria virtual compartida. Típicament OpenMP s'utilitza en sistemes amb un nombre no massa gran de processadors, tot i que hi ha implementacions de sistemes de memòria compartida amb centenars de nuclis. També hi ha versions d'OpenMP per a d'altres tipus de sistemes, com ara clústers i acceleradors, però en aquests casos el rendiment dels programes OpenMP pot ser inferior al dels que utilitzen entorns de programació específics per a aquests sistemes.

OpenMP per a clústers

El sistema blacklight del Pittsburgh Supercomputing Center exposa 4.096 nuclis amb un model de memòria compartida.

El model d'execució d'OpenMP és *fork-join*, com il·lustra la figura 8. Això vol dir que inicialment el programa treballa amb un únic flux fins que arriba a una regió paral·lela, és a dir, quan arriba al constructor `#pragma omp parallel` que posa en marxa un nombre de fluxos esclaus (part *fork* del model). El flux que treballa inicialment és el flux mestre d'aquest conjunt d'esclaus. Els fluxos estan enumerats des de 0 fins al nombre de fluxos -1. El mestre i els esclaus treballen en paral·lel en el bloc que apareix a continuació del constructor. En acabar la regió paral·lela hi ha una sincronització de tots els fluxos, els esclaus moren i únicament queda el flux mestre (part *join* del model). A partir d'aquest moment el mestre continua treballant seqüencialment tret que comenci una altra regió paral·lela. Tal com il·lustra la figura, també cal tenir en compte que hi pot haver més d'un nivell de regions paral·leles (regió paral·lela imbricada).

Figura 8. Model *fork-join* d'OpenMP



El codi següent mostra un exemple de programa en OpenMP en què podem observar alguns elements necessaris en programes OpenMP, com ara la inclusió del fitxer de capçaleres `omp.h`, en què es defineixen les funcions d'OpenMP, la utilització de directives (comencen per `#pragma omp`) per a indicar al compilador la manera en què s'ha de distribuir el treball o la utilització de clàusules per a indicar operacions concretes com per exemple la reducció. En concret, el codi retorna la suma de la multiplicació dels elements de dos vectors.

```
#include <omp.h>
#include <stdio.h>
#include <stdlib.h>

int main (int argc, char *argv[])
{
    int i, n;
    float a[100], b[100], sum;

    /* Algunes inicialitzacions - regió seqüencial */
    n = 100;
```

```
for (i=0; i < n; i++)
    a[i] = b[i] = i * 1.0;
sum = 0.0;

#pragma omp parallel for reduction(+:sum) /* Regió paral·lela */
for (i=0; i < n; i++)
    sum = sum + (a[i] * b[i]);

printf("Suma = %f\n",sum); /* Regió seqüencial */
}
```

1.2.1. Directives OpenMP

Les directives OpenMP segueixen les convencions dels estàndards per a directives de compilació en C/C++, diferencien entre majúscules i minúscules, només es pot especificar un nom de directiva per directiva, i cada directiva s'aplica almenys a la sentència que segueix, que pot ser un bloc estructurat. En directives llargues es pot utilitzar el caràcter \ al final de la línia.

El format general és el següent:

```
#pragma omp nom_directiva [clàusules, ...]
```

en què

- `#pragma omp` es demana a totes les directives OpenMP per a C/C++.
- `nom_directiva` és un nom vàlid de directiva, i ha d'aparèixer després del `pragma` i abans de qualsevol clàusula. En el nostre exemple és `parallel for`.
- `[clàusules, ...]` són opcionals. Les clàusules poden anar en qualsevol ordre i repetir-se quan sigui necessari, a menys que hi hagi alguna restricció. En el nostre cas són `reduction` i `private`.

1.2.2. Creació de fluxos

Tal com s'ha indicat anteriorment, la directiva amb què es creen els fluxos esclaus és `parallel` i s'utilitza de la manera següent:

```
#pragma omp parallel [clàusules]
    bloc
```

Amb aquest `pragma`:

- Es crea un grup de fluxos i el flux que els posa en marxa pren el rol de flux mestre.
- El nombre de fluxos a crear s'obté per la variable d'entorn `OMP_NUM_THREADS` o explícitament amb una crida a la biblioteca tal com veurem més endavant. Si el valor és igual a 0 s'executa de manera seqüencial.
- Hi ha una barrera implícita al final de la regió, de manera que el flux mestre espera que acabin tots els esclaus per a continuar amb l'execució seqüencial.
- Quan dins d'una regió hi ha un altre constructor `parallel`, cada esclau crearia un altre grup de fluxos esclaus dels quals seria el mestre. Això s'anomena *paral·lisme imbricat* i, tot i que es poden programar d'aquesta manera, en algunes implementacions d'OpenMP la creació de grups d'esclaus dins d'una regió paral·lela no està suportada.
- Les clàusules de compartició de variables que suporta la directiva `parallel` són principalment: `private`, `firstprivate`, `lastprivate`, `default`, `shared`, `copyin` i `reduction`. El significat d'aquestes clàusules i d'altres de compartició es veurà més endavant.

1.2.3. Clàusules de compartició de variables

En els casos anteriors hem vist que les directives OpenMP tenen clàusules amb què poden indicar com es du a terme la compartició de variables, les quals estan totes en memòria compartida. Cada directiva té una sèrie de clàusules que s'hi poden aplicar. Aquestes són:

- `private(l·lista)`. Les variables de la llista són privades dels fluxos, que vol dir que cada flux té una variable privada amb aquest nom (les quals poden tenir valors diferents en diferents fluxos). Les variables no s'inicialitzen abans d'entrar en la regió paral·lela (per tant, no podem esperar un valor inicial concret en entrar en la regió paral·lela) i no es guarda el seu valor en sortir de la regió paral·lela (per tant, no podem esperar que l'últim valor d'una variable privada es conservi en finalitzar la regió paral·lela).
- `firstprivate(l·lista)`. Les variables són privades dels fluxos i s'inicialitzen en entrar a la regió paral·lela amb el valor que tingués la variable corresponent del flux mestre.
- `lastprivate(l·lista)`. Són privades dels fluxos i en sortir de la regió paral·lela queden amb el valor de l'última iteració (si estem en un bucle

`for paral·lel`) o secció (veurem més endavant el funcionament de les seccions).

- `shared(llista)`. Indica les variables compartides per tots els fluxos. Per defecte totes són compartides, per la qual cosa en la majoria dels casos no fa realment falta utilitzar aquesta clàusula.
- `default(shared|none)`. Indica com seran les variables per defecte. Si s'especifica `none` caldrà indicar explícitament amb la clàusula `shared` les que es vulgui que siguin compartides.
- `reduction(operador:llista)`. Les variables de la llista s'obtenen per l'aplicació de l'operador, que ha de ser associatiu.
- `copyin(llista)`. S'utilitza per a assignar el valor de la variable en el mestre a variables del tipus `threadprivate`.

1.2.4. Clàusules de divisió del treball

Una vegada s'han generat amb `parallel` un conjunt de fluxos esclaus, aquests i el mestre poden treballar en la resolució d'un problema en paral·lel. Hi ha dues maneres de dividir el treball entre el fluxos: mitjançant les directives `for` i `sections`. Mentre que la primera directiva fa referència al paral·lelisme de dades, la segona fa referència al paral·lelisme funcional.

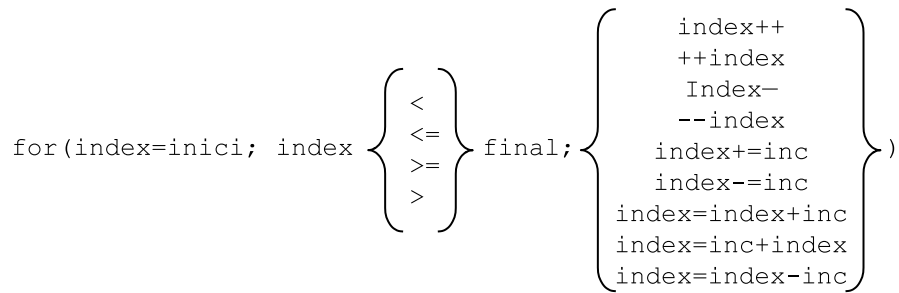
La directiva `for` té la forma següent:

```
#pragma omp for [clàusules]
    bucle for
```

en què:

- Les iteracions les executen en paral·lel els fluxos que ja hi ha, creats prèviament amb `parallel`.
- El bucle `for` ha de tenir una forma especial (forma canònica), tal com mostra la figura 9. La part d'inicialització ha de tenir una assignació; la part de l'increment una suma o resta; la de l'avaluació és la comparació d'una variable entera amb signe amb un valor, utilitzant un comparador major o menor (pot incloure igual); i els valors que apareixen en les tres parts del `for` han de ser enters.

Figura 9. Forma canònica dels bucles en OpenMP



- Hi ha una barrera implícita al final del bucle, tret que s'utilitzi la clàusula `nowait`.
- Les clàusules de compartició de variables que admet són `private`, `first-private`, `lastprivate` i `reduction`.
- Pot aparèixer una clàusula `schedule` per a indicar de quina manera es divideixen les iteracions del `for` entre els fluxos, és a dir la política de planificació.

Les polítiques de planificació disponibles a OpenMP són:

- `schedule(static, mida)`. Les iteracions es divideixen segons la mida de bloc que s'indica. Per exemple, si el nombre d'iteracions del bucle és de 100 i aquestes estan numerades de 0 a 99, i la mida de bloc és 2, es consideren 50 blocs cadascun dels quals amb dues iteracions (blocs amb iteracions 0-1, 2-3, etc.), i si disposem de 4 fluxos, els blocs d'iteracions s'assignen als fluxos cíclicament, amb la qual cosa al flux 0 corresponen les iteracions 0-1, 8-9, etc., al flux 1 les iteracions 2-3, 10-11, etc., al 2 les 4-5, 12-13, etc., i al 3 les 6-7, 14-15, etc. Si no s'indica cap mida de bloc, les iteracions es divideixen per igual entre els fluxos amb blocs de mida màxima: si tenim 12 iteracions i 4 fluxos, s'assignen al flux 0 les iteracions 0-2, a l'1 les 3-5, al 2 les 6-8 i al 4 les 9-11.
- `schedule(dynamic, mida)`. Les iteracions s'agrupen segons la mida de bloc especificada i s'assignen als fluxos dinàmicament quan van acabant el seu treball. En el cas anterior de 100 iteracions, mida de bloc 2 i 4 fluxos, s'assignen inicialment al flux i les iteracions $2i$ i $2i+1$, i a partir d'aquí la resta de blocs de dues iteracions s'assignen als fluxos segons es vagin quedant sense feina. Quan el volum de computació de cada iteració no es coneix a priori, pot ser preferible utilitzar una assignació dinàmica, ja que permet reduir el desbalanceig.
- `schedule(guided, mida)`. Les iteracions s'assignen dinàmicament als fluxos però amb mida de bloc decreixent fins a arribar a la mida de bloc

que s'indica (si no s'indica cap valor per defecte és 1). La mida inicial de bloc i la manera en què decreix depèn de la seva implementació.

- `schedule(runtime)`. Decideix la política de planificació en temps d'execució mitjançant la variable d'entorn `OMP_SCHEDULE`.

La figura 10 mostra les traces obtingudes de l'execució d'una aplicació OpenMP utilitzant les polítiques de planificació `static`, `dynamic` i `guided`. Es pot apreciar que el temps d'execució (eix de les *x*) pot variar molt significativament en funció de la política de planificació.

Figura 10. Traces de l'execució d'una mateixa aplicació OpenMP amb diferents polítiques de planificació per a la distribució del treball



La directiva `sections` té la forma següent:

```
#pragma omp sections [clàusules]
```

```

{
    [#pragma omp section]
    bloc
    [#pragma omp section]
    bloc
    ...
}

```

en què:

- Cada secció s'executa mitjançant un flux. La manera en què es distribueixen les seccions entre els fluxos depèn de la implementació concreta d'OpenMP.
- Hi ha una barrera final tret que s'utilitzi la clàusula `nowait`.
- Les clàusules de compartició de variables que admet són `private`, `firstprivate`, `lastprivate` i `reduction`.

Atès que `for` i `sections` s'utilitzen amb molta freqüència immediatament després d'haver creat els fluxos amb `parallel`, hi ha dues directives combinades que s'utilitzen molt i que s'especifiquen a continuació:

```

#pragma omp parallel for [clàusules]
    bucle for

```

i

```

#pragma omp parallel sections [clàusules]

```

Aquestes són una manera abreviada de posar la directiva `parallel` que conté una única directiva `for` o `sections`. Les dues tenen les mateixes clàusules que `for` i `sections`, excepte `nowait`, que en aquest cas no està permesa ja que acaba el `parallel`, amb el qual el mestre ha d'esperar que acabin tots els esclaus per a seguir l'execució seqüencial.

1.2.5. Directives de sincronització

Dins d'una regió paral·lela pot ser necessari sincronitzar els fluxos en l'accés a algunes variables o zones de codi per a evitar situacions d'interbloqueig. OpenMP ofereix diverses primitives amb aquesta finalitat:

- `single`. El codi afectat per la directiva l'executarà un únic flux. Els fluxos que no estan treballant durant l'execució de la directiva esperen al final. Admet les clàusules `private`, `firstprivate` i `nowait`. No està permès fer bifurcacions cap a un bloc o des d'un bloc `single`. És útil per a seccions

de codi que poden ser no segures per a la seva execució paral·lela (per exemple, per a entrada/sortida).

- `master`. El codi només l'executa el flux mestre. La resta de fluxos no executen aquesta secció de codi.
- `critical`. Protegeix una secció de codi perquè hi pugui accedir un únic flux a la vegada. S'hi pot associar un nom de la forma:

```
#pragma omp critical [nom]
```

de manera que hi pot haver seccions crítiques protegint zones diferents del programa, de manera que a seccions amb noms diferents hi poden accedir diversos fluxos a la vegada. Les regions que tinguin el mateix nom es tracten com la mateixa regió. Totes les que no tenen nom es consideren la mateixa. No està permès fer bifurcacions cap a un bloc o des d'un bloc `critical`.

- `atomic`. Assegura que una posició de memòria es modifiqui sense que múltiples fluxos intentin escriure-la de manera simultània. S'aplica a la sentència que segueix a la directiva. Només s'assegura en mode exclusiu l'actualització de la variable, però no l'avaluació de l'expressió.
- `barrier`. Sincronitza tots els fluxos. Quan un flux arriba a la barrera, espera que arribin els altres i, quan han arribat tots, aquestes segueixen la seva execució.
- `threadprivate`. S'utilitza perquè variables globals es converteixin en locals i persistents en un flux per mitjà de múltiples regions paral·leles.
- `ordered`. Assegura que el codi s'executi en l'ordre en què les iteracions s'executen en la seva forma seqüencial. Pot aparèixer només una vegada en el context d'una directiva `for` o `parallel for`. No està permès fer bifurcacions cap a un bloc o des d'un bloc `ordered`. Només hi pot haver un únic flux executant-se simultàniament en una secció `ordered`. Una iteració d'un bucle no pot executar la mateixa directiva `ordered` més d'una vegada, i no ha d'executar més d'una directiva `ordered`. Un bucle amb una directiva `ordered` ha de contenir una clàusula `ordered`.
- `flush`. Té la forma:

```
#pragma omp flush [llista-de-variables]
```

i assegura que el valor de les variables s'actualitza en tots els fluxos en què són visibles. Si no hi ha una llista de variables s'actualitzaran totes.

1.2.6. Funcions i variables

Les funcions que típicament s'utilitzen més en OpenMP són les relacionades amb el nombre de fluxos, és a dir, `omp_set_num_threads`, `omp_get_num_threads` i `omp_get_thread_num`. Altres funcions importants són:

- `omp_get_max_threads`. Obté la màxima quantitat possible de fluxos.
- `omp_get_num_procs`. Retorna el nombre màxim de processadors que es poden assignar al programa.
- `omp_in_parallel`. Retorna un valor diferent de 0 si s'executa dins d'una regió paral·lela.

Es pot posar o treure el que el nombre de fluxos es pugui assignar dinàmicament en les regions paral·leles (`omp_set_dynamic`), o es pot comprovar si està permès l'ajust dinàmic, amb la funció `omp_get_dynamic`, que retorna un valor diferent de 0 si està permès l'ajust dinàmic del nombre de fluxos.

Es pot permetre desautoritzar el paral·lelisme imbricat amb `omp_set_nested`, o comprovar si està permès, amb `omp_get_nested`, que retornen un valor diferent de 0 si n'està.

També hi ha funcions per a la gestió de *locks* com les que es mostren a continuació:

- `void omp_init_lock(omp_lock_t *lock)` per a inicialitzar un *lock*, que s'inicialitza com a no bloquejat.
- `void omp_init_destroy(omp_lock_t *lock)` per a destruir un *lock*.
- `void omp_set_lock(omp_lock_t *lock)` per a bloquejar un *lock*.
- `void omp_unset_lock(omp_lock_t *lock)` per a desbloquejar un *lock*.
- `void omp_test_lock(omp_lock_t *lock)` per a comprovar si un *lock* està bloquejat o no i així poder evitar bloquejos indesitjats.

1.2.7. Entorn de compilació i execució

Per a compilar un programa OpenMP cal disposar d'un compilador que pugui interpretar les directives que apareixen en el codi. OpenMP està suportat pel compilador GCC des de la versió 4.1 i també per d'altres de propietaris, com ara l'ICC d'Intel. En aquest subapartat utilitzarem el llenguatge C per a exemplificar l'ús d'OpenMP. L'opció de compilació per a GCC és `-fopenmp` o `-openmp` i per a ICC és `-openmp`. Així doncs, la línia següent serviria per a compilar un programa OpenMP amb GCC:

```
gcc -fopenmp -o programa programa.c
```

L'execució del fitxer binari (programa en el nostre exemple) té lloc com qualsevol altre programa, però cal determinar quants fluxos intervinen en l'execució paral·lela. Podem utilitzar una variable d'entorn (`OMP_NUM_THREADS`) que ens indica el nombre de fluxos a utilitzar. Si no s'inicialitza aquesta variable, quan executem el nostre programa OpenMP tindrà un valor per defecte, que acostuma a coincidir amb el nombre de nuclis del node en què estem treballant.

Abans de l'execució es poden establir els valor de la variable, per exemple des de l'interpret d'ordres. Per exemple, si establím el valor de la variable com a continuació:

```
export OMP_NUM_THREADS=6
```

Establím el nombre de fluxos a utilitzar en la regió paral·lela a sis, independentment del nombre de nuclis de què disposem. Això vol dir que es poden utilitzar més fluxos que processadors hi ha al nostre sistema. En termes generals, el rendiment de l'execució d'un programa OpenMP, utilitzant més fluxos que processadors disponibles en el sistema, és inferior al que podem aconseguir utilitzant un flux per processador per la sobrecàrrega que genera la creació i gestió dels fluxos. De fet, podem arribar a observar que el temps d'execució del programa paral·lel sigui més elevat que el del programa seqüencial.

És molt important recordar que la mida del problema és vital per a poder explotar el paral·lisme. Així doncs, encara que disposem de diversos nuclis la mida del problema (el nombre d'interval) pot no ser suficientment gran perquè es noti l'efecte de la paral·lelització, especialment si tenim en compte que hi ha zones que s'executen seqüencialment.

Tot i això, encara que no disposem de suficients processadors per a executar tots els fluxos, podria ser que s'obtingués un rendiment millor gràcies a unes millors particions i localitat de les dades. Per exemple, pot succeir que, en dividir el problema en diversos fluxos, les dades a tractar a cada iteració d'un bucle hi càpiguen en una pàgina de memòria cau i, en conseqüència, es pugui accelerar l'execució de manera molt significativa.

Activitat

Com que actualment la majoria de computadors personals disposen de múltiples nuclis, podeu experimentar i buscar explicació als temps que s'obtenen de l'execució d'algun programa OpenMP variant el nombre de fluxos, essent aquest nombre tant més petit com més gran que el nombre de nuclis disponibles.

Hi ha quatre variables d'entorn. A més d'`OMP_NUM_THREADS`, també hi ha disponible:

- `OMP_SCHEDULE` indica el tipus de planificació per a `for` i `parallel for`.
- `OMP_DYNAMIC` autoritza o desautoritza l'ajust dinàmic del nombre de fluxos.

- OMP_NESTED autoritza o desautoritza el paral·lelisme imbricat. Per defecte no està autoritzat.

2. Models de programació gràfica

En els darrers anys, la capacitat de càlcul dels processadors gràfics⁴ ha esdevingut més elevada que la dels processadors multinucli més avançats per a dur a terme certes tasques. Això ha fet que aquest tipus de dispositius estiguin esdevenint molt populars per al còmput d'algoritmes de propòsit general i no sols per a la generació de gràfics. La computació de propòsit general sobre GPU es coneix popularment com a GPGPU⁵.

⁽⁴⁾GPU, de l'anglès *graphics processing unit*.

⁽⁵⁾De l'anglès *general-purpose computing on graphics processing unit*.

Les GPU proporcionen rendiments molt elevats només per a certes aplicacions per les seves característiques arquitecturals i de funcionament. De manera general, podem dir que les aplicacions que poden aprofitar millor les capacitats de les GPU són les que compleixen les dues condicions següents:

- Treballen sobre vectors de dades grans.
- Tenen un paral·lelisme de gra fi tipus SIMD.

Un dels inconvenients principals a l'hora de treballar amb GPU és la dificultat per al programador per a transformar programes dissenyats per a CPU tradicionals en programes que es puguin executar de manera eficient en una GPU. Per aquest motiu s'han desenvolupat models de programació, tant propietaris (CUDA) com oberts (OpenCL), que proporcionen al programador un nivell d'abstracció, més proper a la programació per a CPU, que li simplifiquen considerablement la tasca.

En aquest apartat estudiarem els models principals de programació per a GPU orientats a aplicacions de propòsit general, específicament CUDA i OpenCL, que són els estàndards actuals per a la programació de processadors gràfics.

2.1. CUDA

CUDA⁶ és una especificació inicialment propietària desenvolupada per Nvidia com a plataforma per als seus productes GPU. CUDA inclou les especificacions de l'arquitectura i un model de programació associat. En aquest subapartat estudiarem el model de programació CUDA; no obstant això, ens referirem a dispositius compatibles amb CUDA en les GPU que implementen l'arquitectura i especificacions definides a CUDA.

⁽⁶⁾De l'anglès *compute unified device architecture*.

2.1.1. Arquitectura compatible amb CUDA

La figura 11 mostra l'arquitectura d'una GPU genèrica compatible amb CUDA. Aquesta arquitectura està organitzada en una sèrie de multiprocessadors⁷ els quals tenen una quantitat elevada de fluxos d'execució. En la figura, dos SM

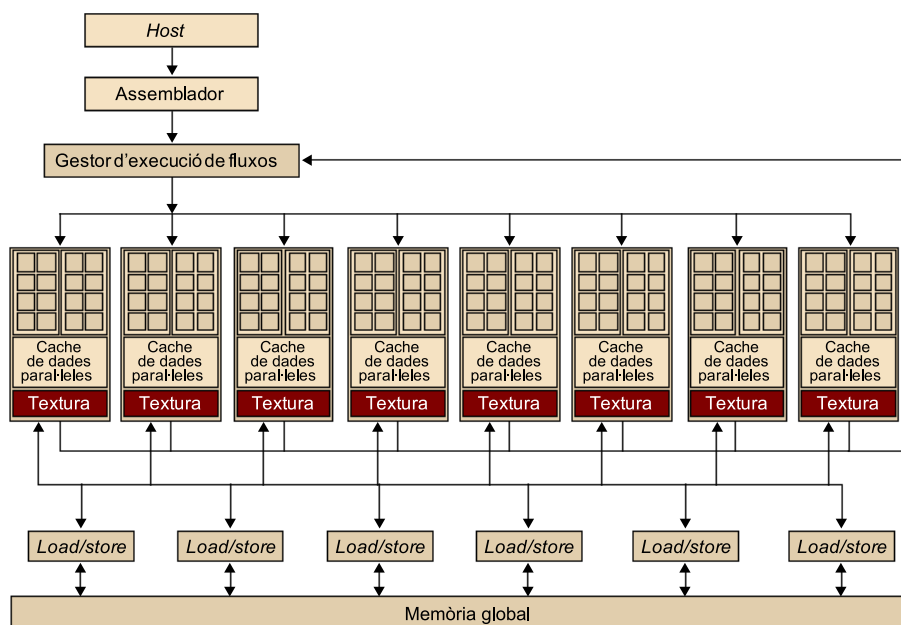
⁽⁷⁾SM, de l'anglès *streaming multiprocessors*.

formen un bloc, tot i que el nombre d'SM per bloc depèn de la implementació concreta del dispositiu. A més, cada SM de la figura té un nombre de *streaming processors* (SP) que comparteixen la lògica de control i memòria cache d'instruccions.

La GPU té una memòria DRAM de tipus GDDR⁸, la qual està indicada en la figura 11 com a memòria global. Aquesta memòria GDDR es diferencia de la memòria DRAM de la placa base del computador en què essencialment s'utilitza per a gràfics (*framebuffer*). Per a aplicacions gràfiques aquesta manté les imatges de vídeo i la informació de les textures. En canvi, per a càlculs de propòsit general aquesta funciona com a memòria externa amb molta amplada de banda però amb una latència una mica més elevada que la memòria típica del sistema. Tot i així, per a aplicacions massivament paral·leles, l'amplada de banda més gran compensa la latència més elevada.

⁽⁸⁾De l'anglès *graphics double data rate*.

Figura 11. Esquema de l'arquitectura d'una GPU genèrica compatible amb CUDA



Noteu que aquesta arquitectura genèrica és molt similar a la G80, la qual suporta fins a 768 fluxos per SM, el qual suma un total de 12.000 fluxos en un únic xip.

L'arquitectura GT200, posterior a la G80, té 240 SP i supera el TFlop de pic teòric de rendiment. Com que els SP són massivament paral·lels, es poden arribar a utilitzar encara més fluxos per aplicació que la G80. La GT200 suporta 1.024 fluxos per SM i en total suma prop de 30.000 fluxos per xip. Per tant, la tendència mostra clarament que el nivell de paral·lisme que suporten les GPU està augmentant ràpidament. Serà molt important, doncs, intentar explotar aquest nivell tan elevat de paral·lisme quan es desenvolupin aplicacions de propòsit general per a GPU.

2.1.2. Entorn de programació

La CUDA es va desenvolupar per tal d'augmentar la productivitat en el desenvolupament d'aplicacions de propòsit general per a GPU. Des del punt de vista del programador, el sistema està compost per un *host*, que és una CPU tradicional, com per exemple un processador d'arquitectura Intel, i d'un o més dispositius⁹ que són GPU.

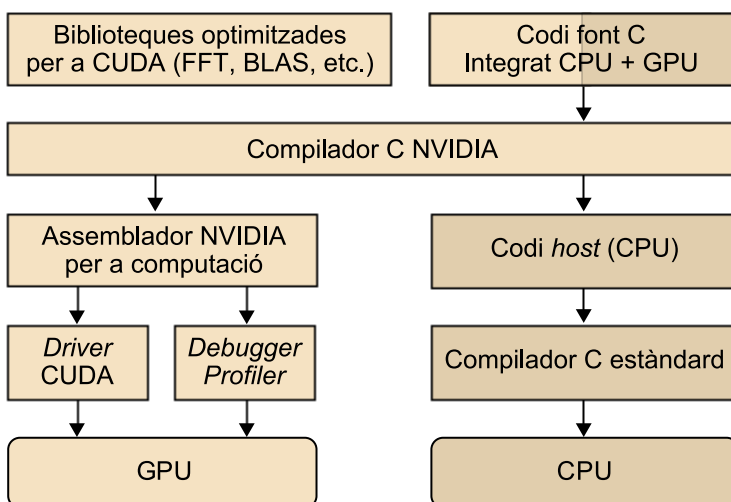
⁹En anglès, *devices*.

La CUDA pertany al model SIMD, per tant, està pensada per a explotar el paral·lisme de dades. Això vol dir que un conjunt d'operacions aritmètiques es poden executar sobre un conjunt de dades de manera simultània. Afortunadament, moltes aplicacions tenen parts amb un nivell molt elevat de paral·lisme de dades.

Un programa en CUDA consisteix en una o més fases que es poden executar o bé en el *host* (CPU) o bé en el dispositiu GPU. Les fases en què hi ha molt poc o gens de paral·lisme entre dades s'implementen en el codi que s'executarà en el *host*, i les fases amb un nivell de paral·lisme entre dades elevat s'implementen en el codi que s'executarà en el dispositiu.

Tal com mostra la figura 12, el compilador de NVIDIA (*nvcc*) s'encarrega de proporcionar la part del codi corresponent al *host* i al dispositiu durant el procés de compilació. El codi corresponent al *host* és simplement codi ANSI C, que es compila mitjançant el compilador de C estàndard del *host*, com si fos un programa per a CPU convencional. El codi corresponent al dispositiu també és ANSI C però amb extensions que inclouen paraules clau per a poder definir funcions que tracten les dades en paral·lel. Aquestes funcions s'anomenen *nuclis* o *kernels*.

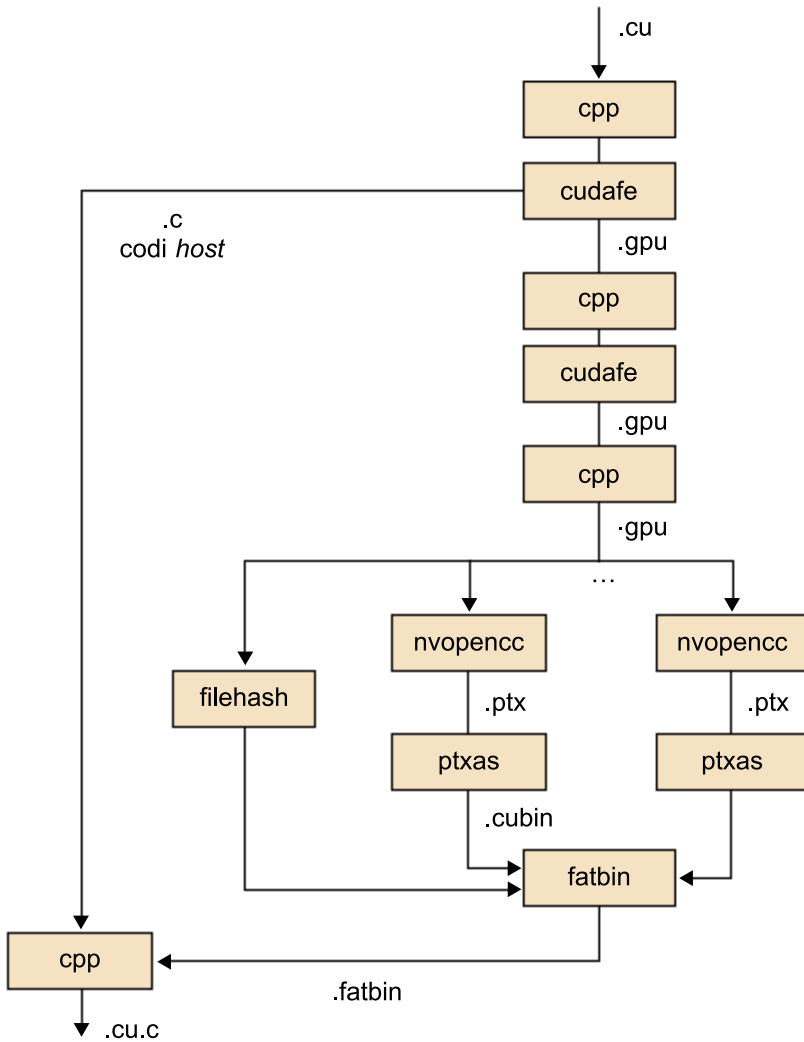
Figura 12. Esquema de blocs de l'entorn de compilació de la CUDA



El codi del dispositiu es torna a compilar amb *nvcc* i llavors ja es pot executar en el dispositiu GPU. En situacions en què no hi ha cap dispositiu disponible o bé el nucli és més apropiat per a una CPU, també es poden executar els nuclis

en una CPU convencional mitjançant eines d'emulació que proporciona la plataforma CUDA. La figura 13 mostra totes les etapes del procés de compilació i la taula 1 descriu com el compilador `nvcc` interpreta els diferents tipus de fitxers d'entrada.

Figura 13. Passos en la compilació de codi CUDA, des de `.cu` fins a `.cu.c`



Taula 1. Interpretació de `nvcc` dels fitxers d'entrada

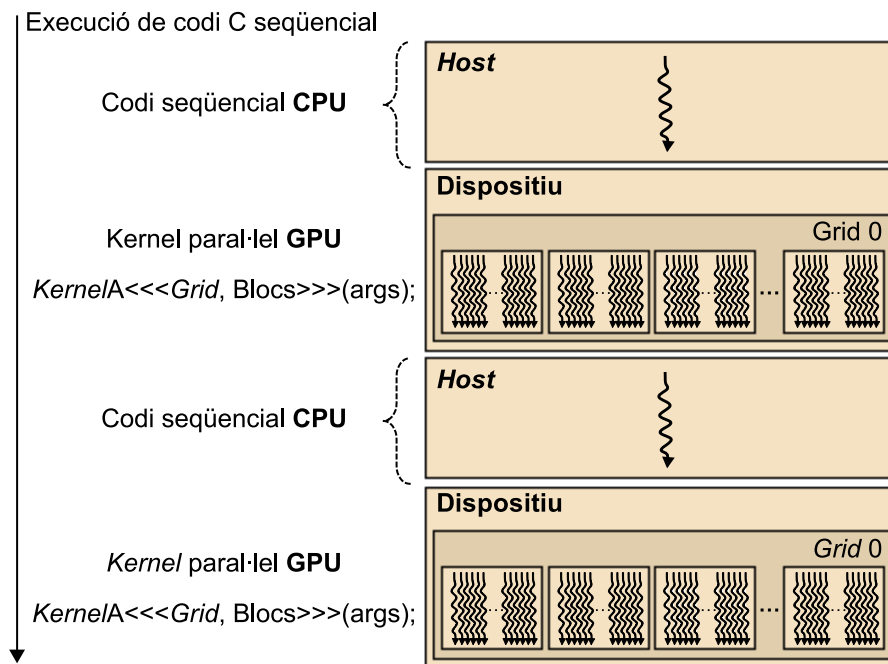
<code>.cu</code>	Codi font CUDA que conté tant el codi del <i>host</i> com les funcions del dispositiu.
<code>.cup</code>	Codi font CUDA preprocessat que conté tant el codi del <i>host</i> com les funcions del dispositiu.
<code>.c</code>	Fitxer de codi font C.
<code>.cc, .cxx, .cpp</code>	Fitxer de codi font C++.
<code>.gpu</code>	Fitxer intermedi gpu.
<code>.ptx</code>	Fitxer assemblador intermedi ptx.
<code>.o, .obj</code>	Fitxer d'objecte.
<code>.a, .lib</code>	Fitxer de biblioteca.

.res	Fitxer de recurs.
.so	Fitxer d'objecte compartit.

Per a explotar el paral·lisme entre dades, els nuclis han de generar una quantitat de fluxos d'execució força elevada (per exemple, unes desenes o centenars de milers de fluxos). Hem de tenir en compte que els fluxos de CUDA són molt més lleugers que els fluxos de CPU. De fet podem assumir que per a generar i planificar aquests fluxos només necessitem uns pocs cicles de rellotge a causa del suport de maquinari, en contrast amb els fluxos convencionals per a CPU que normalment requereixen milers de cicles.

L'execució d'un programa típic CUDA es mostra en la figura 14. L'execució comença amb execució en el *host* (CPU). Quan s'invoca un nucli, l'execució es mou cap al dispositiu (GPU) on es genera un nombre molt elevat de fluxos. El conjunt de tots aquest fluxos que es generen quan s'invoca un nucli s'anomena *grid*. En la figura 14 es mostren dos *grids* de fluxos. Un nucli finalitza quan tots els seus fluxos finalitzen la seva execució en el *grid* corresponent. Un cop finalitzat el nucli, l'execució del programa continua en el *host* fins que s'invoca un altre nucli.

Figura 14. Etapes en l'execució d'un programa típic CUDA



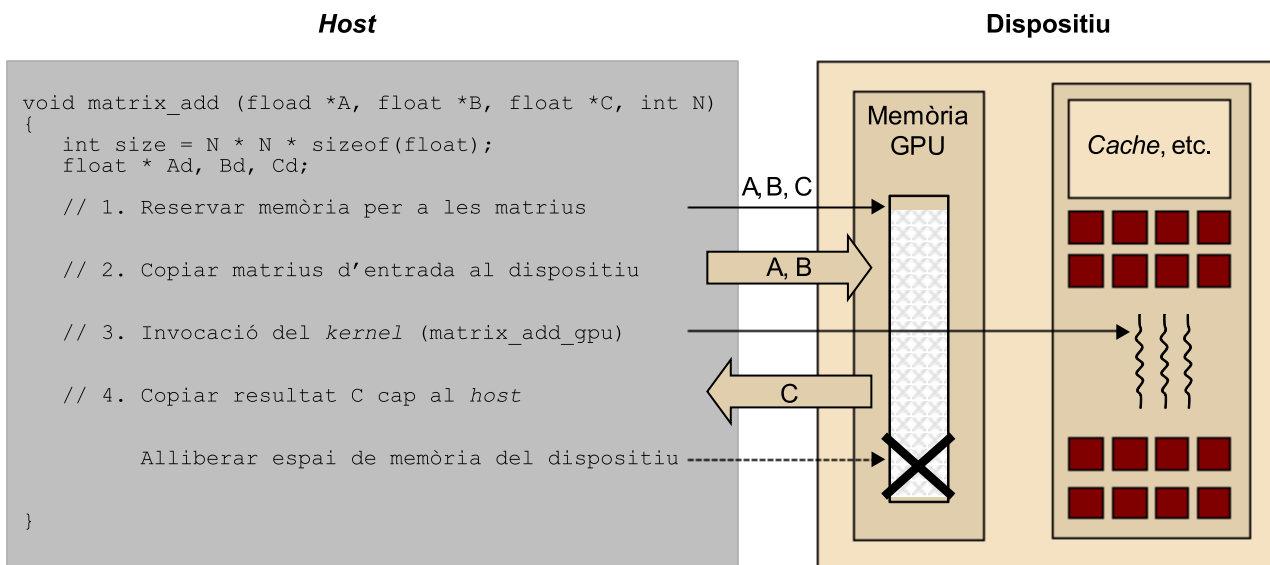
2.1.3. Model de memòria

És important remarcar que la memòria del *host* i del dispositiu són espais de memòria completament separats. Això reflecteix la realitat que els dispositius són típicament targetes que tenen la seva pròpia memòria DRAM. Per a executar un nucli en el dispositiu GPU normalment cal seguir aquests passos:

- Reservar memòria en el dispositiu (pas 1 de la figura 15).
- Transferir les dades necessàries des del *host* a l'espai de memòria assignat al dispositiu (pas 2 de la figura 15).
- Invocar l'execució del nucli en qüestió (pas 3 de la figura 15).
- Transferir les dades amb els resultat des del dispositiu cap al *host* i alliberar la memòria del dispositiu (si ja no és necessària) un cop finalitzada l'execució del nucli (pas 4 de la figura 15).

L'entorn CUDA proporciona una interfície de programació que simplifiquen aquestes tasques al programador. Per exemple, només cal especificar quines dades cal transferir des del *host* al dispositiu i viceversa.

Figura 15. Esquema dels passos per a l'execució d'un nucli en una GPU



Durant tot aquest apartat utilitzarem el mateix programa d'exemple que fa la suma de dues matrius. L'exemple següent mostra el codi corresponent a la implementació seqüencial en C estàndard de la suma de matrius per a arquitectura de tipus CPU.

```

void matrix_add_cpu (float *A, float *B, float *C, int N)
{
  int i, j, index;
  for (i=0; i<N; i++){
    for (j=0; j<N; j++){
      index = i+j*N;
      C[index] = C[index] + B[index];
    }
  }
}

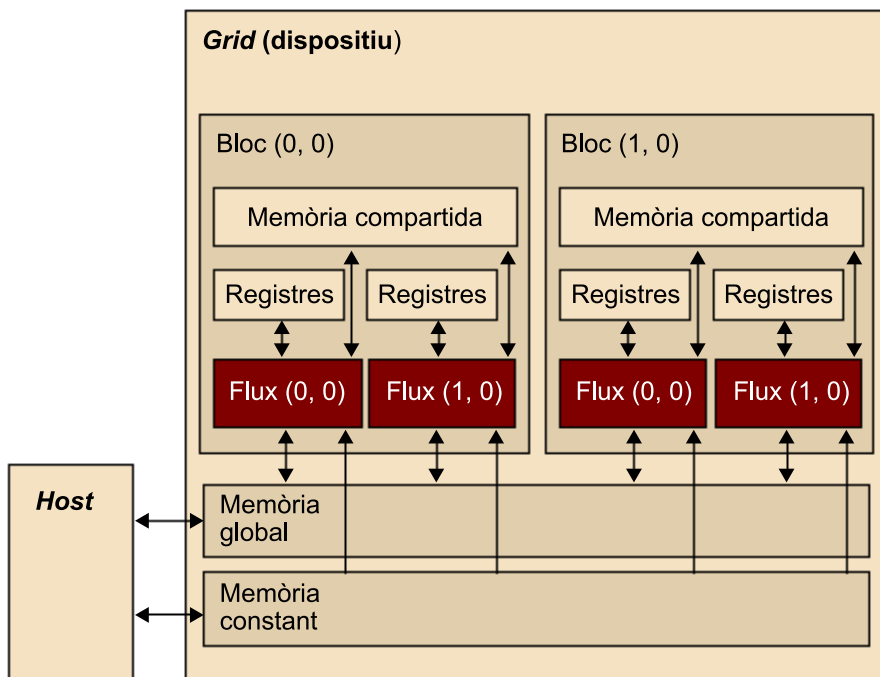
```

```
int main() {
    matrix_add_cpu(a, b, c, N);
}
```

La figura 16 mostra el model de memòria de CUDA exposat al programador en termes d'assignació, transferència i utilització dels diferents tipus de memòria del dispositiu. En la part inferior de la figura es troben les memòries de tipus global i constant. Aquestes memòries són aquelles a les quals el *host* pot transferir dades de manera bidireccional, per *grid*. Des del punt de vista del dispositiu, es pot accedir a diferents tipus de memòria amb els modes següents:

- Accés de lectura/escriptura a la memòria global, per *grid*.
- Accés només de lectura a la memòria constant, per *grid*.
- Accés de lectura/escriptura als registres, per flux.
- Accés de lectura/escriptura a la memòria local, per flux.
- Accés de lectura/escriptura a la memòria compartida, per bloc.
- Accés de lectura/escriptura.

Figura 16. Model de memòria de CUDA



La interfície de programació per a assignar i alliberar memòria global del dispositiu consisteix en dues funcions bàsiques: `cudaMalloc()` i `cudaFree()`. La funció `cudaMalloc()` es pot cridar des del codi del *host* per a assignar un espai de memòria global per a un objecte. Com haureu observat, aquesta funció és molt similar a la funció `malloc()` de la biblioteca de C estàndard, ja que CUDA és una extensió del llenguatge C i cerca mantenir les interfícies tan similars a les originals com sigui possible.

La funció `cudaMalloc()` té dos paràmetres. El primer paràmetre és l'adreça del punter cap a l'objecte un cop s'hagi assignat l'espai de memòria. Aquest punter és genèric i no depèn de cap tipus d'objecte, per tant s'haurà de fer *cast* a tipus (`void **`). El segon paràmetre és la mida de l'objecte que es vol assignar en bytes. La funció `cudaFree()` allibera l'espai de memòria de l'objecte indicat com a paràmetre de la memòria global del dispositiu. El codi següent mostra un exemple de com s'utilitzen aquestes dues funcions. Després de fer el `cudaMalloc()`, `Matriu` apunta a una regió de la memòria global del dispositiu que s'hi ha assignat.

```
float *Matriu;
int mida = AMPLADA * LLARGÀRIA * sizeof(float);
cudaMalloc((void **) &Matriu, mida);
...
cudaFree(Matriu);
```

Un cop que un programa ha assignat memòria global del dispositiu per als objectes o estructures de dades del programa, es poden transferir les dades que faran falta per a la computació des del *host* cap al dispositiu. Això es fa mitjançant la funció `cudaMemcpy()`, que permet transferir dades entre memòries. Cal tenir en compte que la transferència és asíncrona.

La funció `cudaMemcpy()` té quatre paràmetres. El primer és un punter en l'adreça de destinació en què s'han de copiar les dades. El segon paràmetre apunta a les dades que s'han de copiar. El tercer paràmetre especifica el nombre de bytes que s'han de copiar. Finalment, el quart paràmetre indica el tipus de memòria involucrat en la còpia, que pot ser un dels següents:

- `cudaMemcpyHostToHost`: de la memòria del *host* cap a la memòria del mateix *host*.
- `cudaMemcpyHostToDevice`: de la memòria del *host* cap a la memòria del dispositiu.
- `cudaMemcpyDeviceToHost`: de la memòria del dispositiu cap a la memòria del *host*.
- `cudaMemcpyDeviceToDevice`: de la memòria del dispositiu cap a la memòria del dispositiu.

Cal tenir en compte que aquesta funció es pot utilitzar per a copiar dades de la memòria d'un mateix dispositiu però no entre diferents dispositius. A continuació es mostra un exemple de transferència de dades entre *host* i dispositiu basat en l'exemple de la figura 15:

```
void matrix_add (float *A, float *B, float *C, int N)
{
    int size = N * N * sizeof(float);
    float * Ad, Bd, Cd;
```

```

// 1. Reservar memòria per a les matrius
cudaMalloc((void **) &Ad, size);
cudaMalloc((void **) &Bd, size);
cudaMalloc((void **) &Cd, size);
// 2. Copiar matrius d'entrada al dispositiu
cudaMemcpy(Ad, A, size, cudaMemcpyHostToDevice);
cudaMemcpy(Bd, B, size, cudaMemcpyHostToDevice);
...
// 4. Copiar resultat C cap al host
cudaMemcpy(C, Cd, size, cudaMemcpyDeviceToHost);
...

```

A més dels mecanismes d'assignació de memòria i transferència de dades, CUDA també suporta diferents tipus de variables. Els diferents tipus de variables que utilitzen els diversos tipus de memòria són utilitzades en diversos àmbits i tindran cicles de vida diferents, com resumeix la taula 2.

Taula 2. Diferents tipus de variables de CUDA

Declaració variables	Tipus de memòria	Àmbit	Cicle de vida
Per defecte (diferents de vectors)	Registre	Flux	Nucli
Vectors per defecte	Local	Flux	Nucli
<code>__device__, __shared__, int SharedVar;</code>	Compartida	Bloc	Nucli
<code>__device__, int GlobalVar;</code>	Global	<i>Grid</i>	Aplicació
<code>__device__, __constant__, int ConstVar;</code>	Constant	<i>Grid</i>	Aplicació

2.1.4. Definició de nuclis

El codi que s'executa en el dispositiu (nucli) és la funció que executen els diferents fluxos durant la fase paral·lela, cadascun en el rang de dades que li correspon. Cal recordar que la CUDA segueix el model SPMD i, per tant, tots els fluxos executen el mateix codi. A continuació es mostra la funció o nucli de la suma de matrius i la seva crida.

```

__global__ matrix_add_gpu (float *A, float *B, float *C, int N)
{
    int i = blockIdx.x * blockDim.x + threadIdx.x;
    int j = blockIdx.y * blockDim.y + threadIdx.y;
    int index = i + j*N;
    if (i<N && j<N){
        C[index] = A[index] + B[index];
    }
}

int main(){

```



```

dim3 dimBlock(blocksize, blocksize);
dim3 dimGrid(N/dimBlock.x, N/dimBlock.y);
matrix_add_gpu<<<dimGrid, dimBlock>>>(a, b, c, N);
}

```

Podem observar la utilització de la paraula clau específica de la CUDA `__global__` davant de la declaració de `matrix_add_gpu()`. Aquesta paraula clau indica que aquesta funció és un nucli i que cal cridar-lo des del *host* per a generar el *grid* de fluxos que executarà el nucli en el dispositiu. A més de `__global__` hi ha dues altres paraules clau que es poden utilitzar davant de la declaració d'una funció:

- La paraula clau `__device__` indica que la funció declarada és una funció CUDA de dispositiu. Una funció de dispositiu s'executa únicament en un dispositiu CUDA i només es pot cridar des d'un nucli o des d'una altra funció de dispositiu. Aquestes funcions no poden tenir ni crides recursives ni crides indirectes a funcions mitjançant punters.
- La paraula clau `__host__` indica que la funció és una funció de *host*, és a dir, una funció simple de C que s'executa en el *host* i, per tant, que pot ser cridada des de qualsevol funció de *host*. Per defecte totes les funcions en un programa CUDA són funcions de *host* si és que no s'especifica cap paraula clau en la definició de la funció.

Les paraules clau `__host__` i `__device__` es poden utilitzar simultàniament en la declaració d'una funció. Aquesta combinació fa que el compilador generi dues versions de la mateixa funció: una que s'executa en el *host* i que només es pot cridar des d'una funció de *host*, i una altra que s'executa en el dispositiu i que només es pot cridar des del dispositiu o funció de nucli.

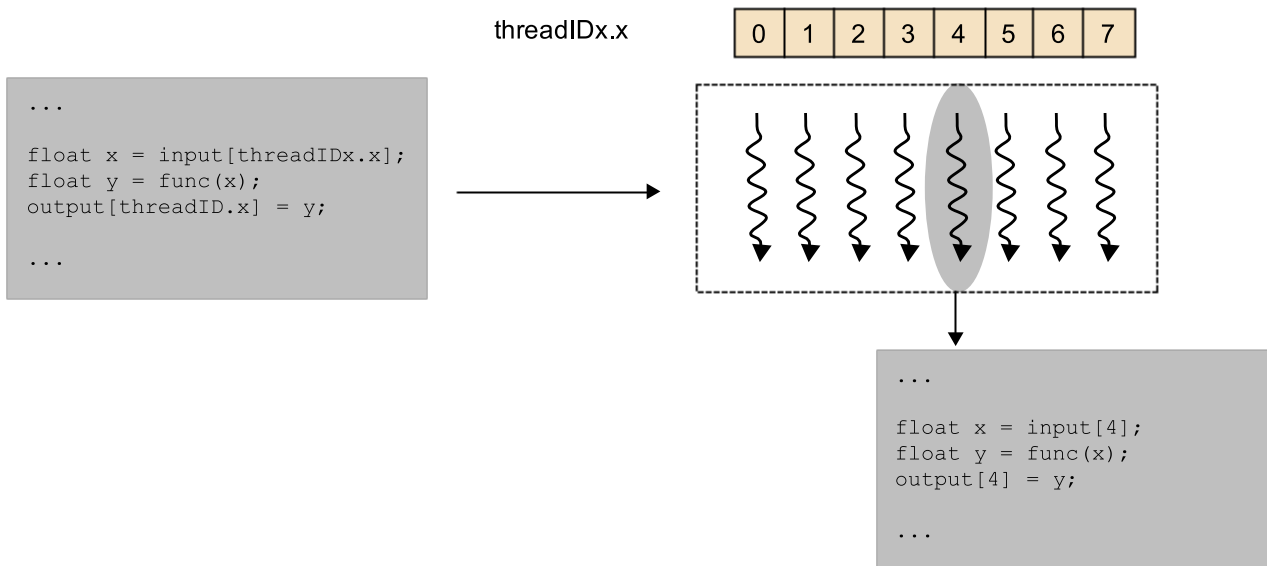
2.1.5. Organització de fluxos

En la CUDA, un nucli s'executa mitjançant un conjunt de fluxos (per exemple, un vector o una matriu de fluxos). Com que tots els fluxos executen el mateix nucli (model SIMT⁽¹⁰⁾) es necessita un mecanisme que permeti diferenciar-los, i així poder assignar la part corresponent de les dades a cada flux d'execució. La CUDA incorpora paraules clau per a fer referència a l'índex d'un flux (per exemple `threadIdx.x` i `threadIdx.y` si tenim en compte dues dimensions).

⁽¹⁰⁾De l'anglès *single instruction multiple threads*.

La figura 17 mostra com el nucli fa referència a l'identificador de flux i que durant la seva execució, a cadascun dels fluxos l'identificador es substitueix pel valor que li correspon. Per tant, les variables `threadIdx.x` i `threadIdx.y` tindran diferents valors per a cadascun dels fluxos d'execució. Noteu que les coordenades reflecteixen l'organització multidimensional dels fluxos d'execució tot i que l'exemple de la figura 17 només fa referència a una dimensió (`threadIdx.x`).

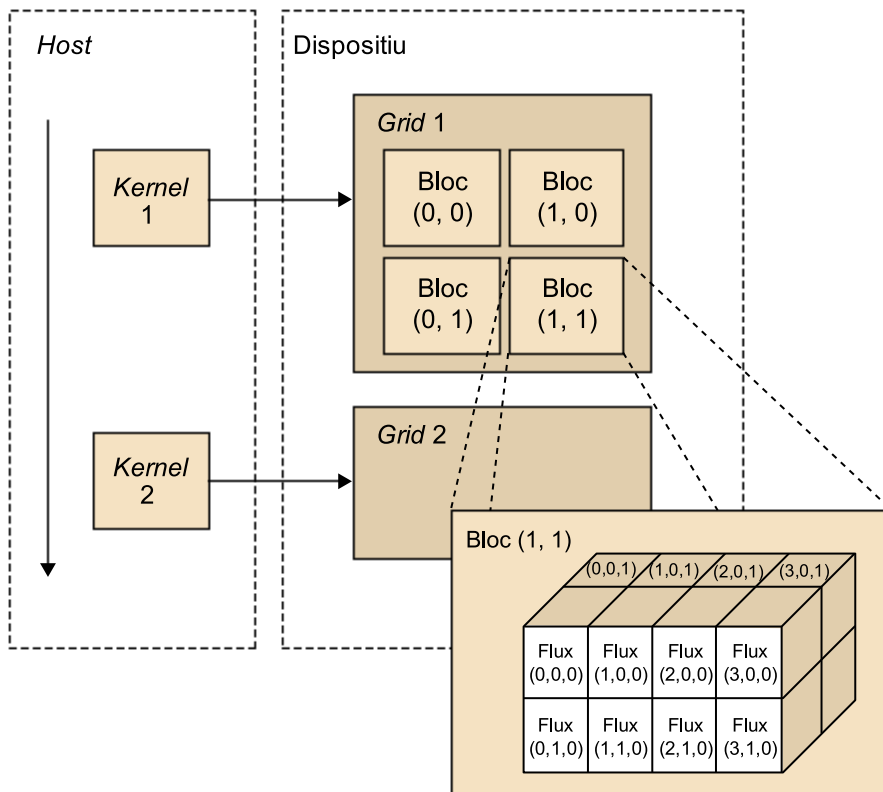
Figura 17. Execució d'un nucli CUDA en un vector de fluxos



Normalment, els *grid* que s'utilitzen en CUDA estan formats per molts fluxos (de prop de milers o fins i tot de milions de fluxos). Els fluxos d'un *grid* estan organitzats en una jerarquia de dos nivells, com es pot veure en la figura 18. També es pot observar com el kernel 1 crea el Grid 1 per a la seva execució. El nivell superior d'un *grid* consisteix en un o més blocs de fluxos. Tots els blocs d'un *grid* tenen el mateix nombre de fluxos i han d'estar organitzats de la mateixa manera. En la figura 18, el Grid 1 consta de quatre blocs i està organitzat com una matriu de 2x2 blocs.

Cada bloc d'un *grid* té una coordenada única en un espai de dues dimensions mitjançant les paraules clau `blockIdx.x` i `blockIdx.y`. Els blocs s'organitzen en un espai de tres dimensions amb un màxim de 512 fluxos d'execució. Les coordenades dels fluxos d'execució en un bloc s'identifiquen amb tres índexs (`threadIdx.x`, `threadIdx.y` i `threadIdx.z`) tot i que no totes les aplicacions necessiten utilitzar les tres dimensions de cada bloc de fluxos. En la figura 18, cada bloc està organitzat en un espai de 4x2x2 fluxos d'execució.

Figura 18. Exemple d'organització dels fluxos d'un grid en CUDA



Quan el *host* fa una crida a un nucli s'han d'especificar les dimensions del *grid* i dels blocs de fluxos mitjançant paràmetres de configuració. El primer paràmetre especifica les dimensions del *grid* en termes de nombre de blocs, i el segon especifica les dimensions de cada bloc en termes de nombre de fluxos. Tots dos paràmetres són de tipus `dim3`, que essencialment és una estructura de C amb tres camps (x , y , z) de tipus enter sense signe. Com que els *grids* són grups de blocs en dues dimensions, el tercer camp de paràmetres de configuració del *grid* s'ignora (qualsevol valor serà vàlid). A continuació teniu un exemple en què dues variables d'estructures de tipus `dim3` defineixen el *grid* i els blocs de l'exemple de la figura 18:

```
// Configuració de les dimensions de grid i blocs
dim3 dimGrid(2, 2, 1);
dim3 dimBlock(4, 2, 2);

// Invocació del kernel (suma de matrius)
matrix_add_gpu<<<dimGrid, dimBlock>>>(a, b, c, N);
```

La CUDA també ofereix un mecanisme per a sincronitzar els fluxos d'un mateix bloc mitjançant la funció de tipus barrera `__syncthreads()`. Quan es crida la funció `__syncthreads()`, el flux que l'executa quedarà bloquejat fins que tots els fluxos del seu bloc arribin a aquest mateix punt. Això serveix per a assegurar que tots els fluxos d'un bloc han completat una fase abans de passar a la següent.

2.2. OpenCL

L'OpenCL és una interfície estàndard, oberta, lliure i multiplataforma per a la programació paral·lela. La motivació principal per al desenvolupament de l'OpenCL va ser la necessitat de simplificar la tasca de programació portable i eficient de la creixent quantitat de plataformes heterogènies, com ara la CPU multinucli, la GPU o fins i tot sistemes encastats. L'OpenCL va ser concebuda per Apple tot i que la va acabar desenvolupant el grup Khronos, que és el mateix que va impulsar l'OpenGL i n'és responsable.

L'OpenCL consisteix en tres parts: l'especificació d'un llenguatge multiplataforma, una interfície per a l'entorn de computació i una interfície per a coordinar la computació paral·lela entre processadors heterogenis. L'OpenCL utilitza un subconjunt de C99 amb extensions per al paral·lelisme i utilitza l'estàndard de representació numèrica IEEE 754 per a garantir la interoperabilitat entre plataformes.

Hi ha moltes similituds entre l'OpenCL i la CUDA tot i que l'OpenCL té un model de gestió de recursos més complex, ja que suporta múltiples plataformes i portabilitat entre diferents fabricants. L'OpenCL suporta models de paral·lelisme tant per a dades com per a tasques. En aquest subapartat ens centrarem en el model de paral·lelisme per a dades, el qual és equivalent al de la CUDA.

2.2.1. Model de paral·lelisme per a dades

De la mateixa manera que en la CUDA, un programa en l'OpenCL està format per dues parts: els nuclis que s'executen en un o més dispositius, i un programa *host* que invoca i controla l'execució dels nuclis. Quan s'invoca un nucli, el codi s'executa en feines elementals⁽¹⁾ que corresponen als fluxos de la CUDA. Les feines elementals i les dades associades a cada feina elemental es defineixen a partir del rang d'un espai d'índexs N-dimensional⁽²⁾. Les feines elementals formen grups de feines⁽³⁾ que corresponen al blocs de la CUDA. Les feines elementals tenen un identificador global que és únic. A més, els grups de feines elementals s'identifiquen dins del rang N-dimensional i, per a cada grup, cadascuna de les feines elementals té un identificador local que anirà des de 0 fins a la mida del grup-1. Per tant, la combinació de l'identificador del grup i de l'identificador local dintre del grup també identifica de manera única una feina elemental. La taula 3 resumeix algunes de les equivalències entre l'OpenCL i la CUDA.

⁽¹⁾En anglès, *work items*.

⁽²⁾En anglès, *NDRanges*.

⁽³⁾En anglès, *work groups*.

Taula 3. Algunes correspondències entre OpenCL i CUDA

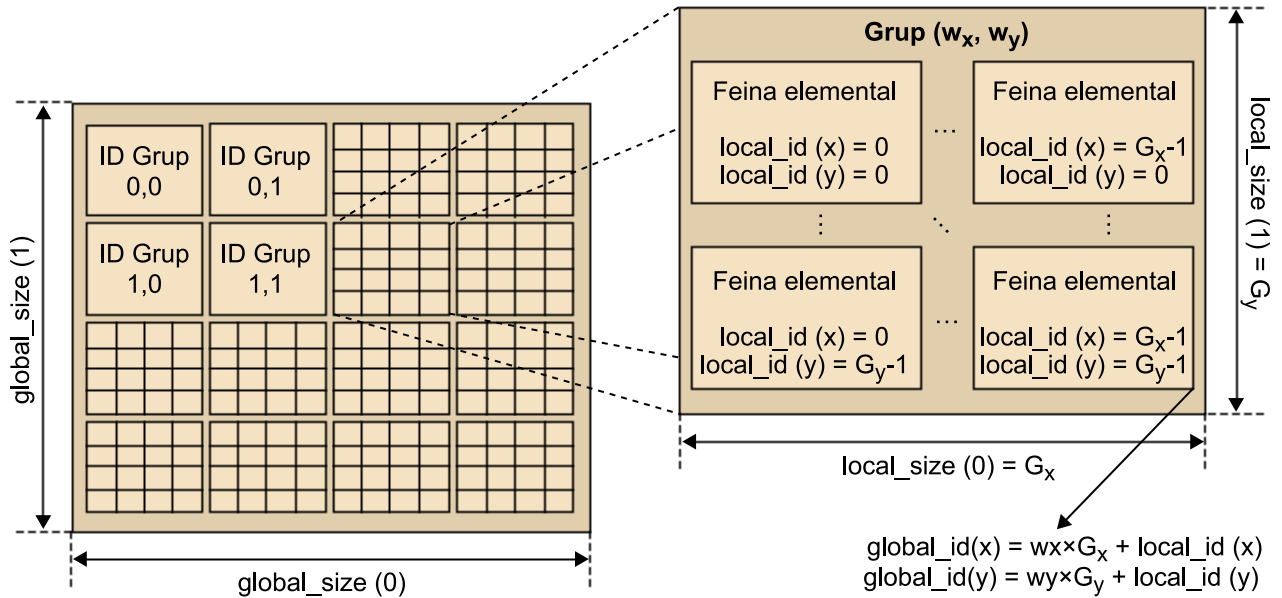
OpenCL	CUDA
Nucli	Nucli
Programa <i>host</i>	Programa <i>host</i>

OpenCL	CUDA
<i>NDRange</i> (rang N-dimensional)	<i>Grid</i>
Feina elemental (<i>work item</i>)	Flux
Grup de feines (<i>work group</i>)	Bloc
<code>get_global_id(0);</code>	<code>blockIdx.x * blockDim.x + threadIdx.x</code>
<code>get_local_id(0);</code>	<code>threadIdx.x</code>
<code>get_global_size(0);</code>	<code>gridDim.x*blockDim.x</code>
<code>get_local_size(0);</code>	<code>blockDim.x</code>

La figura 19 mostra el model de paral·lelisme per a dades de l'OpenCL. El rang N-dimensional (equivalent al *grid* en la CUDA) conté les feines elementals. En l'exemple de la figura, el nucli utilitza un rang de dues dimensions. Mentre que en la CUDA cada flux té uns valors de `blockIdx` i `threadIdx` que es combinen per a identificar el flux, en l'OpenCL disposem d'interfícies per a identificar les feines elementals de les dues maneres que hem vist. Per una banda, la funció `get_global_id()`, donada una dimensió, retorna l'identificador únic de feina elemental en la dimensió especificada. En l'exemple de la figura, les crides `get_global_id(0)` i `get_global_id(1)` retornen l'índex de les feines elementals en les dimensions X i Y, respectivament. Per l'altra banda, la funció `get_local_id()`, donada una dimensió, retorna l'identificador de la feina elemental dins del seu grup en la dimensió especificada. Per exemple, `get_local_id(0)` és equivalent a `threadIdx.x` en la CUDA.

L'OpenCL també disposa de les funcions `get_global_size()` i `get_local_size()` que, donada una dimensió, retornen la quantitat total de feines elementals i la quantitat de feines elementals dintre d'un grup en la dimensió especificada, respectivament. Per exemple, `get_global_size(0)` retorna la quantitat de feines elementals que és equivalent a `gridDim.x*blockDim.x` en la CUDA.

Figura 19. Exemple de rang N-dimensional



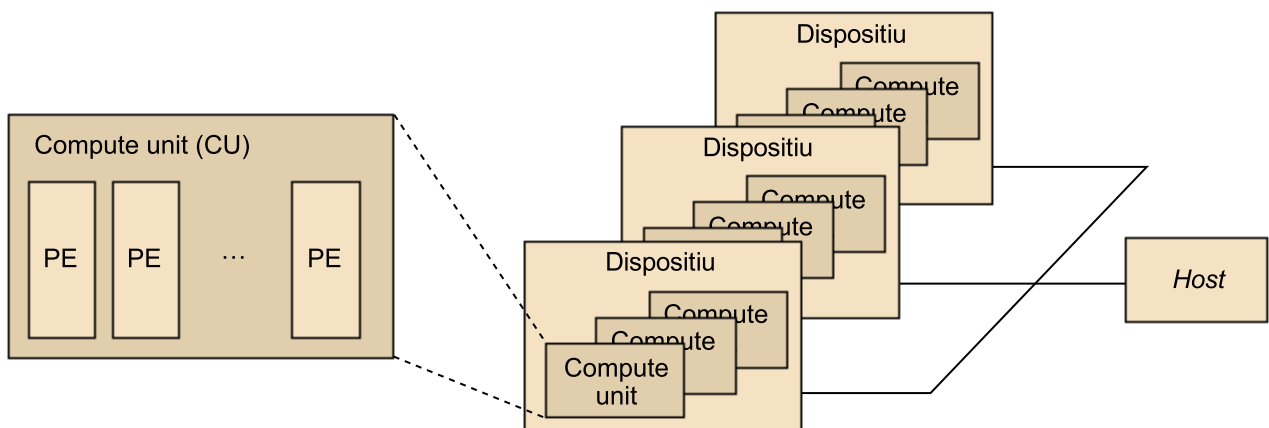
Es poden observar les feines elementals, els grups que formen i identificadors associats.

Les feines elementals dintre del mateix grup es poden sincronitzar utilitzant barreres que són equivalents a `__syncthreads()` de la CUDA. En canvi, les feines elementals de diferents grups no es poden sincronitzar, excepte si és per la terminació del nucli o la invocació d'un de nou.

2.2.2. Arquitectura conceptual

La figura 20 mostra l'arquitectura conceptual de l'OpenCL, la qual està formada per un *host* (típicament una CPU que executa el programa *host*) connectat a un o més dispositius OpenCL. Un dispositiu OpenCL està compost per una o més *compute units* (CU) que corresponen als SM de la CUDA. Finalment, un CU està format per un o més *processing elements* (PE) que corresponen als SP de la CUDA. El programa s'acabarà executant en els PE.

Figura 20. Arquitectura conceptual de l'OpenCL



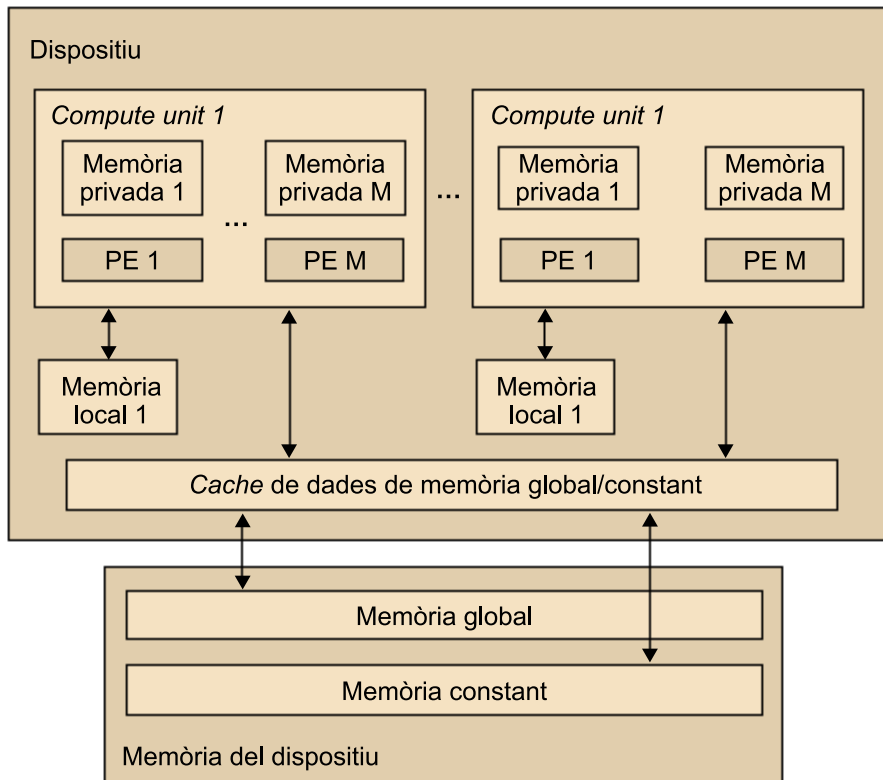
2.2.3. Model de memòria

En un dispositiu OpenCL hi ha disponible una jerarquia de memòria que inclou tipus de memòria diferents, tal com mostra la figura 21: global, constant, local i privada.

- La memòria global és la que es pot utilitzar per a totes les unitats de càlcul d'un dispositiu.
- La memòria constant és la memòria que es pot utilitzar per a emmagatzemar dades constants per a accés només de lectura per a totes les unitats de càlcul d'un dispositiu durant l'execució d'un nucli. El processador *host* és responsable d'assignar i iniciar els objectes de memòria que resideixen en l'espai de memòria.
- La memòria local és la memòria que es pot utilitzar per a les feines elementals d'un grup.
- La memòria privada és la que es pot utilitzar únicament per a una unitat de càlcul única. Això és similar als registres en una única unitat de càlcul o un únic nucli d'una CPU.

Així doncs, tant la memòria global com la constant corresponen als tipus de memòria amb el mateix nom de la CUDA, la memòria local correspon a la memòria compartida de la CUDA i la memòria privada correspon a la memòria local de la CUDA.

Figura 21. Arquitectura i jerarquia de memòria d'un dispositiu OpenCL



Per a crear i manejar objectes en la memòria d'un dispositiu, l'aplicació que s'executa en el *host* utilitza la interfície d'OpenCL, ja que els espais de memòria del *host* i dels dispositius són principalment independents l'un de l'altre. La interacció entre l'espai de memòria del *host* i dels dispositius pot ser de dos tipus: copiant dades explícitament o bé mapant/desmapant regions d'un objecte OpenCL en la memòria. Per copiar dades explícitament, el *host* posa ordres a la cua per a transferir dades entre la memòria de l'objecte i la memòria del *host*. Aquestes ordres de transferència poden ser o bé bloquejants o bé no bloquejants. Quan s'utilitza una crida bloquejant es pot accedir a les dades des del *host* amb seguretat, un cop la crida ha finalitzat. En canvi, per a fer una transferència no bloquejant, la crida a la funció d'OpenCL finalitza immediatament i es desenvia tot i que la memòria des del *host* no és segura per a utilitzar. La interacció mapant/desmapant objectes en la memòria permet que el *host* pugui tenir en el seu espai de memòria una regió corresponent a objectes OpenCL. Les ordres de mapeig/desmapeig també poden ser bloquejants o no bloquejants.

2.2.4. Gestió de nuclis i de dispositius

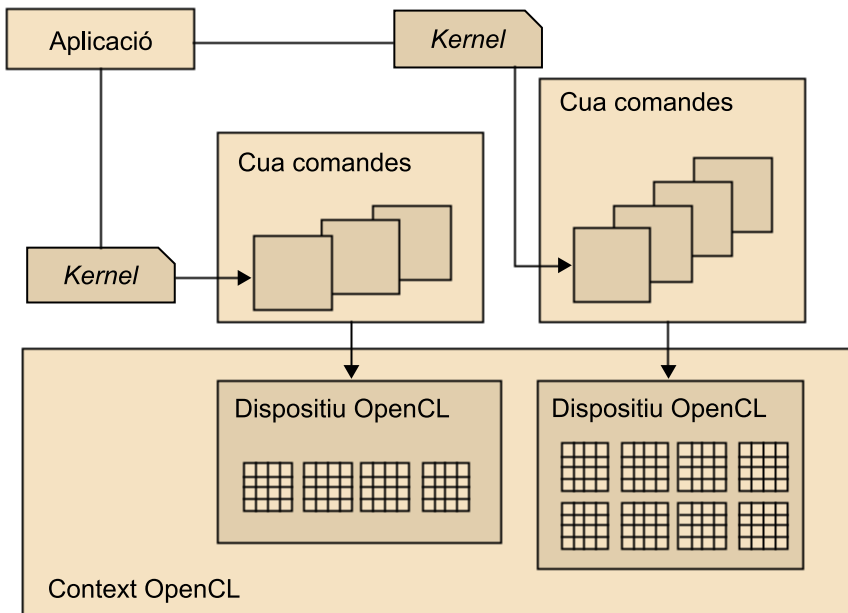
Els nuclis de l'OpenCL tenen la mateixa estructura que els de la CUDA. Per tant, són funcions que es declaren començant amb la paraula clau `__kernel`, que és equivalent a `__global` de la CUDA. A continuació es mostra la implementació del nostre exemple de suma de matrius en OpenCL. Fixeu-vos que els arguments del nucli corresponents a les matrius estan declarats com a `__global`, ja que es troben en la memòria global, i les dues matrius d'entrada també

estan declarades com a `const`, ja que només caldrà fer accessos en la modalitat de lectura. En la implementació del cos del nucli s'utilitza `get_global_id()` en les dues dimensions de la matriu per tal de definir l'índex associat. Aquest índex s'utilitza per a seleccionar les dades que corresponen a cadascuna de les feines elementals que instanciï el nucli.

```
__kernel void matrix_add_opencl ( __global const float *A,
    __global const float *B,
    __global float *C,
    int N) {
    int i = get_global_id(0);
    int j = get_global_id(1);
    int index = i + j*N;
    if (i<N && j<N){
        C[index] = A[index] + B[index];
    }
}
```

El model de gestió de dispositius de l'OpenCL és molt més sofisticat que el de la CUDA, ja que l'OpenCL permet abstraure diferents plataformes maquinari. Els dispositius es gestionen mitjançant contextos. Com mostra la figura 22, per a gestionar un o més dispositius (dos en l'exemple de la figura), primer cal crear un context que contingui els dispositius mitjançant o bé la funció `clCreateContext()` o bé la funció `clCreateContextFromType()`. Normalment cal indicar a les funcions `CreateContext` la quantitat i el tipus de dispositius del sistema mitjançant la funció `clGetDeviceIDs()`. Per a executar qualsevol tasca en un dispositiu, primer cal crear una cua d'ordres per al dispositiu mitjançant la funció `clCreateCommandQueue()`. Un cop s'ha creat una cua per al dispositiu, el codi que s'executa en el *host* pot inserir-hi un nucli i els paràmetres de configuració associats. Un cop el dispositiu estigui disponible per a executar el nucli següent de la cua, aquest nucli s'eliminarà de la cua i passarà a ser executat.

Figura 22. Gestió de dispositius en OpenCL mitjançant contextos



A continuació es mostra el codi corresponent al *host* per a executar la suma de matrius mitjançant el nucli de l'exemple de codi anterior. Suposem que la definició i inicialització de variables s'ha dut a terme correctament i que la funció del nucli `matrix_add_opencl()` està disponible. En el primer pas es crea un context i, un cop s'han obtingut els dispositius disponibles, es crea una cua d'ordres que utilitzarà el primer dispositiu disponible dels que hi hagi en el sistema. En el segon pas es defineixen els objectes que ens caldran en la memòria (les tres matrius A, B i C). Les matrius A i B es defineixen de lectura, la matriu C es defineix d'escriptura. Noteu que en la definició de les matrius A i B s'utilitza l'opció `CL_MEM_COPY_HOST_PTR`, que indica que les dades de les matrius A i B es copiaran dels punters especificats (`srcA` i `srcB`). En el tercer pas es defineix el nucli que s'executarà posteriorment mitjançant `clCreateKernel` i també s'especifiquen els arguments de la funció `matrix_add_opencl`. A continuació s'envia el nucli a la cua d'ordres definida prèviament i, finalment, es llegeixen els resultats de l'espai de memòria corresponent a la matriu C per mitjà del punter `dstC`. En aquest exemple no hem entrat en detall en molts dels paràmetres de les diferents funcions de la interfície de l'OpenCL, per tant es recomana repassar l'especificació d'aquests que està disponible en la pàgina web <http://www.khronos.org/opencl/>.

```
main() {
    // Inicialització de variables, etc.
    (...)

    // 1. Creació del context i cua al dispositiu
    cl_context context = clCreateContextFromType(0, CL_DEVICE_TYPE_GPU, NULL, NULL, NULL);
    // Per obtenir la llista de dispositius GPU associats al context
    size_t cb;
    clGetContextInfo(context, CL_CONTEXT_DEVICES, 0, NULL, &cb);
    cl_device_id *devices = malloc(cb);
```

```
clGetContextInfo( context, CL_CONTEXT_DEVICES, cb, devices, NULL);
cl_cmd_queue cmd_queue = clCreateCommandQueue(context, devices[0], 0 , NULL);

// 2. Definició dels objectes a memòria (matrius A, B i C)
cl_mem memobjs[3];
memobjs[0] = clCreateBuffer(context, CL_MEM_READ_ONLY | CL_MEM_COPY_HOST_PTR,
    sizeof(cl_float)*n, srcA, NULL);
memobjs[1] = clCreateBuffer(context, CL_MEM_READ_ONLY | CL_MEM_COPY_HOST_PTR,
    sizeof(cl_float)*n, srcB, NULL);
memobjs[2] = clCreateBuffer(context, CL_MEM_WRITE_ONLY, sizeof(cl_float)*n, NULL, NULL);

// 3. Definició del kernel i arguments
cl_program program = clCreateProgramWithSource(context, 1, &program_source, NULL, NULL);
cl_int err = clBuildProgram(program, 0, NULL, NULL, NULL, NULL);
cl_kernel kernel = clCreateKernel(program, "matrix_add_opencl", NULL);
err = clSetKernelArg(kernel, 0, sizeof(cl_mem), (void *)&memobjs[0]);
err |= clSetKernelArg(kernel, 1, sizeof(cl_mem), (void *)&memobjs[1]);
err |= clSetKernelArg(kernel, 2, sizeof(cl_mem), (void *)&memobjs[2]);
err |= clSetKernelArg(kernel, 3, sizeof(int), (void *)&N);

// 4. Invocació del kernel
size_t global_work_size[1] = n;
err = clEnqueueNDRangeKernel(cmd_queue, kernel, 1, NULL, global_work_size,
    NULL, 0, NULL, NULL);

// 5. Lectura dels resultats (matriu C)
err = clEnqueueReadBuffer(context, memobjs[2], CL_TRUE, 0, n*sizeof(cl_float),
    dstC, 0, NULL, NULL);

(...)
```

3. Models de programació per memòria distribuïda

En aquest apartat ens centrarem en els models principals de programació per a memòria distribuïda, MPI, que és l'estàndard *de facto* des de fa anys, i els models PGAS, que intenten donar solució a reptes importants que presenten les característiques dels futurs sistemes d'altres prestacions.

3.1. MPI

L'MPI¹⁴ és una biblioteca de pas de missatges que es pot utilitzar tant en programes C com Fortran, des dels quals es fan crides a funcions de l'MPI per a la gestió i comunicació de processos.

La figura 23 mostra l'estructura general d'un programa MPI en què es pot apreciar com inicialment hi ha un únic procés que executa el codi de manera seqüencial fins que s'inicialitza l'entorn MPI, que és quan el codi paral·lel comença. En la secció de codi paral·lel es farà la tasca concreta i s'utilitzarà pas de missatges per a comunicar els diferents processos MPI. Finalment, un cop acabat el treball que es vol executar en paral·lel, s'indica la finalització del codi paral·lel i es torna a executar la resta del codi de manera seqüencial fins que acaba l'execució del programa.

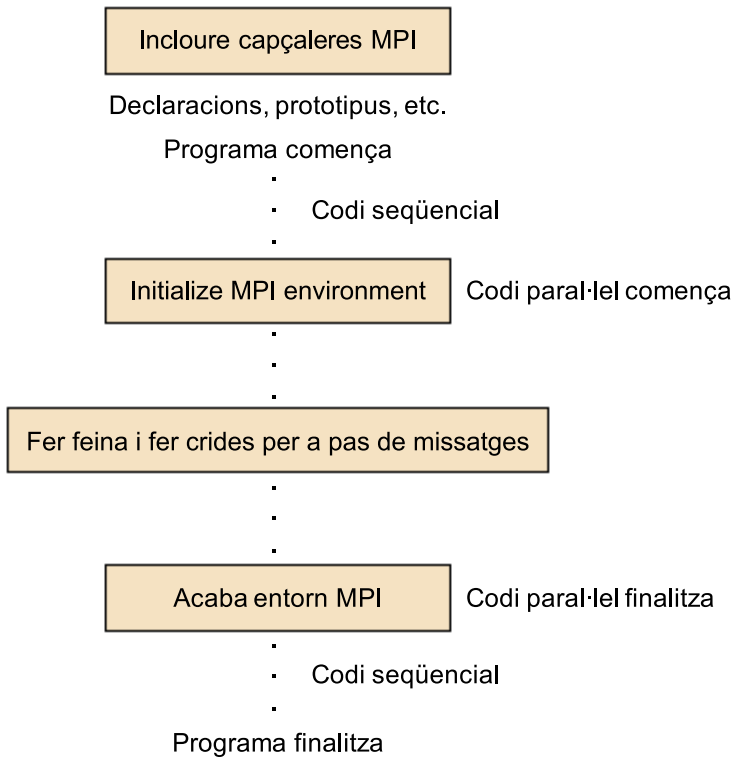
Cal tenir present que, a diferència de l'OpenMP, el nombre de processos no es pot establir en temps d'execució, sinó que s'especifica quan es vol executar el programa tal com veurem més endavant. També cal recordar que tots els processos MPI executaran exactament el mateix codi.

⁽¹⁴⁾De l'anglès *message passing interface*.

Vegeu també

La biblioteca MPI s'ha vist en el mòdul "Introducció a la computació d'altres prestacions" d'aquesta assignatura.

Figura 23. Estructura general d'un programa MPI



A continuació es mostra l'exemple bàsic de programa MPI (variant del típic "hello world") en què es pot observar l'estructura general d'un programa MPI. Aquesta és una versió similar a la que s'ha vist en el primer mòdul de l'assignatura però exemplificant-ne la utilització en Fortran.

```

program hello

use fmpi
! include "mpif.h"
call MPI_INIT( ierr )
call MPI_COMM_RANK( MPI_COMM_WORLD, myid, ierr )
call MPI_COMM_SIZE( MPI_COMM_WORLD, numprocs, ierr )

write (*,*) "Hello from ",myid
write (*,*) "Numprocs is ",numprocs

call MPI_FINALIZE(ierr)
stop
end
  
```

Una de les maneres més habituals d'utilitzar l'MPI és mitjançant el model SPMD. En aquest cas caldrà diferenciar el codi que executa el procés mestre (típicament el procés amb rang igual a zero) del que executen els processos esclaus, tal com es mostra a continuació.

```
main (int argc, char *argv[])
```

```
{
    MPI_Init(&argc, &argv);
    MPI_Comm_rank(MPI_COMM_WORLD, &myrank); // Consulta el rang del procés

    if (myrank == 0)
        master(); // codi que executa el procés mestre
    else
        slave(); // codi que executen els processos esclaus

    MPI_Finalize();
}
```

3.1.1. Comunicadors

Un comunicador permet definir una sèrie de processos entre els quals es poden realitzar comunicacions. Cada procés pot estar en diferents comunicadors i tindrà un identificador en cadascun (també anomenat *rang*), estant els identificadors entre 0 i el nombre de processos en el comunicador menys 1. En posar-se en marxa l'MPI, tots els processos són a `MPI_COMM_WORLD`, la qual cosa és una constant que identifica un comunicador que inclou a tots els processos.

A part de `MPI_COMM_WORLD`, també es poden definir comunicadors amb un nombre inferior de processos, per exemple per a comunicar les dades en cadascuna de les files de processos en una graella, amb la qual cosa es pot fer *broadcast* en les diferents files de processos sense que les comunicacions dels uns afectin els altres. Tot i això, també es poden fer servir en les comunicacions punt a punt.

Hi ha dos tipus de comunicadors: els intracomunicadors i els intercomunicadors. Els intracomunicadors s'utilitzen per a enviar missatges entre els processos en aquest comunicador, i els intercomunicadors s'utilitzen per a enviar missatges entre processos en diferents comunicadors. Nosaltres ens centrem en comunicacions entre processos d'un mateix comunicador. Les comunicacions entre processos en comunicadors diferents poden tenir sentit si en dissenyar biblioteques es creen comunicadors i un usuari de la biblioteca vol comunicar un procés del seu programa amb un altre de la biblioteca MPI.

Un comunicador consta d'un grup, que és una col·lecció ordenada de processos als quals s'associen identificadors, i un context, que és un identificador que associa el sistema al grup. Addicionalment, a un comunicador s'hi pot associar una topologia virtual.

Activitat

Exploreu com es poden definir topologies virtuals de comunicadors i quina en pot ser la utilitat.

En l'MPI hi ha dos tipus bàsics de comunicacions: les comunicacions punt a punt, que involucren dos processos que intercanvien missatges només entre ells, i les comunicacions col·lectives, que involucren un conjunt de processos que realitzen tots ells una tasca concreta col·lectivament.

3.1.2. Comunicacions punt a punt

En l'exemple anterior els processos no interaccionen. El normal és que els processos no treballin de manera independent, sinó que intercanviïn informació per mitjà de pas de missatges. Per a enviar missatges entre dos processos (un d'origen i un de destinació) s'utilitzen les comunicacions punt a punt. Més endavant es mostra un fragment de codi MPI en què el procés amb rang 0 li envia l'enter x al procés amb rang 1. El programa mostra la manera normal de treballar amb pas de missatges. S'executa el mateix programa en tots els processos, però processos diferents executen parts diferents del codi. En l'exemple, s'utilitzen funcions `MPI_Send` i `MPI_Recv` per a enviar i rebre missatges respectivament. Aquestes dues funcions són bloquejants i les seves definicions es mostren a continuació.

```
int MPI_Send(void *buf, int count, MPI_Datatype datatype,
             int dest, int tag, MPI_Comm comm)

int MPI_Recv(void *buf, int count, MPI_Datatype datatype,
             int source, int tag, MPI_Comm comm,
             MPI_Status *status)
```

Els paràmetres es descriuen a continuació:

- `buf` conté l'inici de la zona de memòria de la qual s'agafaran les dades a enviar o bé on s'emmagatzemaran les dades que es reben.
- `count` indica el nombre de dades a enviar o l'espai disponible per a rebre'n (no el nombre de dades del missatge que es rep, ja que la mida la determina qui envia).
- `datatype` és el tipus de dades a transferir, i ha de ser un tipus MPI (`MPI_Datatype`), que en el nostre exemple és `MPI_CHAR`.
- `dest` i `source` són els identificadors dels processos als quals s'envia i dels quals es rep el missatge respectivament. Es pot utilitzar la constant `MPI_ANY_SOURCE` per a indicar que es pot rebre des de qualsevol procés.
- `tag` s'utilitza per a diferenciar entre missatges, i ha de coincidir en el procés que envia i el que rep. Es pot utilitzar `MPI_ANY_TAG` per a indicar que el missatge és compatible amb missatges amb qualsevol identificador.

- `comm` és el comunicador dins del qual es fa la comunicació. És del tipus `MPI_Comm`, i en l'exemple s'utilitza l'identificador de comunicador format per a tots els processos (`MPI_COMM_WORLD`).
- `status` referencia una variable de tipus `MPI_Status`. En el programa no s'utilitza, però conté informació del missatge que s'ha rebut, i es pot consultar per a identificar alguna característica del missatge, per exemple, la seva longitud, el procés origen, etc.

A continuació es mostra un fragment de codi MPI en què el procés amb rang 0 li envia l'enter x al procés amb rang 1:

```
MPI_Comm_rank(MPI_COMM_WORLD,&myrank); /* find rank */
if (myrank == 0) {
    int x;
    MPI_Send(&x, 1, MPI_INT, 1, msgtag, MPI_COMM_WORLD);
}
else if (myrank == 1) {
    int x;
    MPI_Recv(&x, 1, MPI_INT, 0, msgtag, MPI_COMM_WORLD, status);
}
```

En els missatges s'han d'utilitzar tipus de dades MPI. Per raons de portabilitat, l'MPI redefineix els seus tipus de dades elementals. La taula següent mostra els tipus estàndard per a C i Fortran:

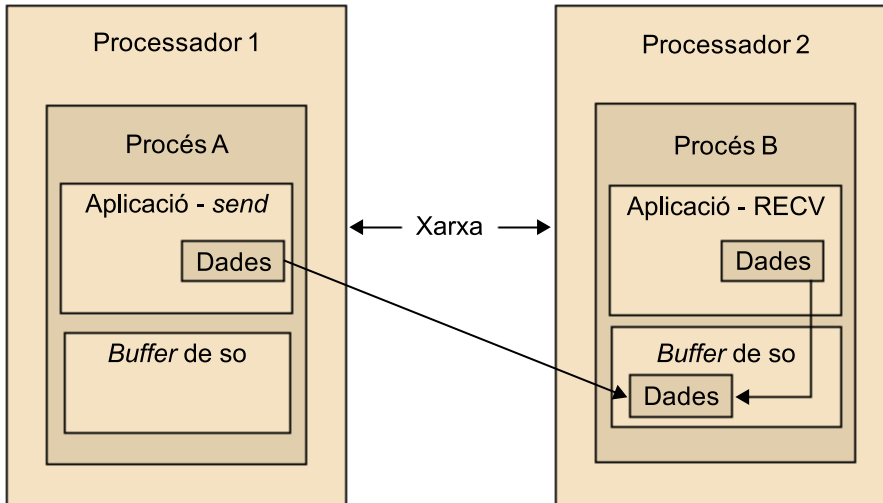
Taula 4. Tipus de dades MPI per a C i Fortran

Tipus de dades C		Tipus de dades Fortran	
MPI_CHAR	signed char	MPI_CHARACTER	character(1)
MPI_WCHAR	wchar_t wide character		
MPI_SHORT	signed short int		
MPI_INT	signed int	MPI_INTEGER MPI_INTEGER1 MPI_INTEGER2 MPI_INTEGER4	integer integer*1 integer*2 integer*4
MPI_LONG	signed		
MPI_LONG_LONG_INT MPI_LONG_LONG	signed long long int		
MPI_SIGNED_CHAR	Signed char		
MPI_UNSIGNED_CHAR	Unsigned char		
MPI_UNSIGNED_SHORT	Unsigned short int		
MPI_UNSIGNED	Unsigned int		
MPI_SIGNED_CHAR	Signed char		

Tipus de dades C		Tipus de dades Fortran	
MPI_UNSIGNED_CHAR	Unsigned char		
MPI_FLOAT	Float	MPI_REAL MPI_REAL2 MPI_REAL4 MPI_REAL8	Real Real*2 Real*4 Real*8
MPI_DOUBLE	Double	MPI_DOUBLE_PRECISION	Double precision
MPI_LONG_DOUBLE	Long double		
MPI_C_COMPLEX MPI_C_FLOAT_COMPLEX	Float_Complex	MPI_COMPLEX	Complex
MPI_C_DOUBLE_COMPLEX	Double_Complex	MPI_DOUBLE_COMPLEX	Double complex
MPI_C_LONG_DOUBLE_COMPLEX	Long double_Complex		
MPI_C_BOOL	_Bool	MPI_LOGICAL	Logical
MPI_C_LONG_DOUBLE_COMPLEX	Long double_Complex		
MPI_INT8_T MPI_INT16_T MPI_INT32_T MPI_INT64_T	Int_8_t Int_16_t Int_32_t Int_64_t		
MPI_UINT8_T MPI_UINT16_T MPI_UINT32_T MPI_UINT64_T	Uint_8_t uint_16_t uint_32_t uint_64_t		
MPI_BYTE	8 binary digits	MPI_BYTE	8 binary digits
MPI_PACKED	Data packed or un- packed with MPI_PACK()/ MPI_UNPACK	MPI_PACKED	Data packed or un- packed with MPI_PACK()/ MPI_UNPACK

La comunicació que hem vist anteriorment mitjançant `MPI_Send` i `MPI_Recv` és de tipus bloquejant, ja que es fa una sincronització entre el procés que envia i el que rep el missatge. Amb `MPI_Send` i `MPI_Recv` el procés que envia el missatge es continua executant un cop acabat l'enviament, i el que rep, si quan sol·licita la recepció del missatge, aquest encara no ha arribat, es queda bloquejat fins que el missatge arribi.

L'MPI proporciona altres possibilitats per a enviaments bloquejants, per exemple `MPI_Ssend`, `MPI_Bsend` i `MPI_Rsend`. Aquestes funcions es diferencien en la manera de gestionar l'enviament, podent gestionar el *buffer* de comunicació o determinar que el procés prengui les dades directament de la memòria. En la figura 24 s'il·lustra el funcionament de pas de missatges mitjançant *buffer*.

Figura 24. Funcionament de pas de missatge a MPI mitjançant *buffer*.

La funció `MPI_Sendrecv` combina en una crida l'enviament i la recepció entre dos processos però continua essent una funció bloquejant.

L'MPI també proporciona comunicació no bloquejant, en què el procés receptor sol·licita el missatge i si aquest no ha arribat continua la seva execució. Les funcions són `MPI_Isend` i `MPI_Irecv`. Si el procés receptor no ha rebut el missatge pot arribar un moment en què per força hagi d'esperar-lo per a continuar l'execució. Les definicions d'aquestes dues funcions es mostren a continuació:

```
int MPI_Isend(void *buf, int count, MPI_Datatype datatype,
             int dest, int tag, MPI_Comm comm)

int MPI_Irecv(void *buf, int count, MPI_Datatype datatype,
              int source, int tag, MPI_Comm comm,
              MPI_Status *status)
```

També es disposa de la funció `MPI_Wait()` per a esperar l'arribada d'un missatge, i de `MPI_Test()` per a comprovar si l'operació s'ha completat.

A continuació es mostra el mateix exemple del desenvolupament del fragment de codi MPI en què el procés amb rang 0 li envia l'enter x al procés amb rang 1, però ara utilitzant crides asíncrones en canvi de crides síncrones o bloquejants.

```
MPI_Comm_rank(MPI_COMM_WORLD, &myrank); /* find rank */
if (myrank == 0) {
    int x;
    MPI_Isend(&x, 1, MPI_INT, 1, msgtag, MPI_COMM_WORLD, req1);
    compute(); // Realitza algun càlcul mentre es realitza l'enviament
    MPI_Wait(req1, status);
}
else if (myrank == 1) {
```

```
int x;
MPI_Recv(&x, 1, MPI_INT, 0, msgtag, MPI_COMM_WORLD, status);
}
```

3.1.3. Comunicacions col·lectives

A més de les comunicacions punt a punt, l'MPI ofereix una sèrie de funcions per a dur a terme comunicacions en què intervenen tots els processos d'un comunicador. Sempre que sigui possible fer les comunicacions per mitjà d'aquestes comunicacions col·lectives, és convenient fer-ho així, ja que facilita la programació evitant possibles errors, i si les funcions estan optimitzades per al sistema on estem treballant és possible que es desenvolupin programes més eficients amb comunicacions col·lectives i no pas amb comunicacions punt a punt.

Algunes de les comunicacions col·lectives més útils es presenten a continuació:

a) Barrera:

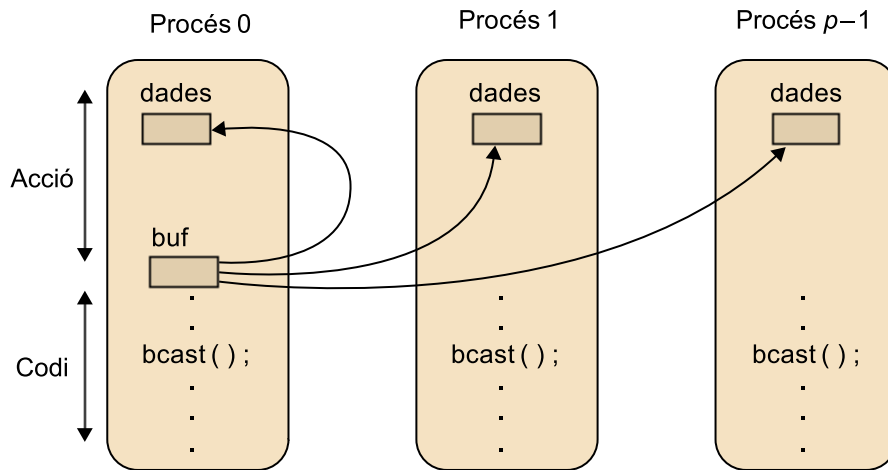
```
int MPI_Barrier(MPI_Comm comm)
```

Estableix una barrera. Tots els processos esperen que tots arribin a la barrera, per a continuar l'execució un cop hi hagin arribat. S'utilitza un comunicador que estableix el grup de processos que s'estan sincronitzant. Totes les funcions de comunicacions col·lectives retornen un codi d'error.

b) *Broadcast*:

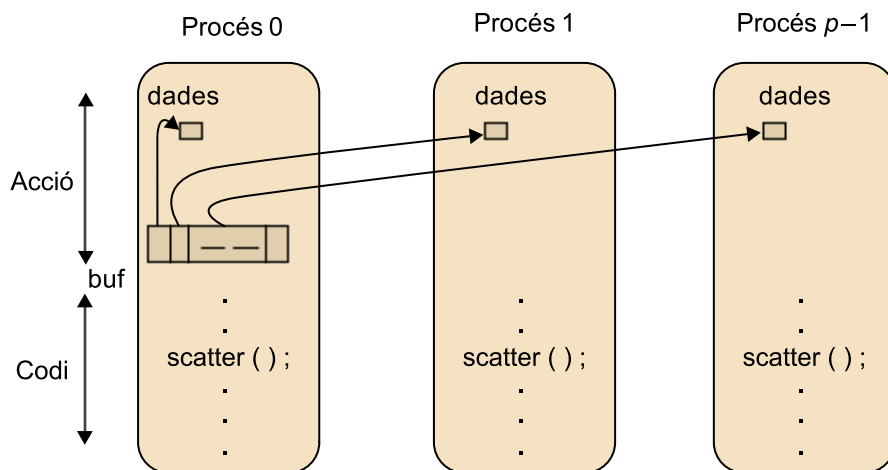
```
int MPI_Bcast(void *buffer, int count,
              MPI_Datatype datatype, int root, MPI_Comm comm)
```

Du a terme una operació de *broadcast* (comunicació un a tots), en què s'envien `count` dades del tipus `datatype` des del procés arrel (`root`) a la resta de processos en el comunicador. Tots els processos que intervenen han de cridar la funció indicant el procés que actua com a arrel. En l'arrel, les dades que s'envien es prenen de la zona apuntada pel *buffer*, i els que reben emmagatzemen en la memòria reservada en el *buffer*. La figura 25 il·lustra el funcionament d'aquesta operació.

Figura 25. Funcionament de l'operació *broadcast* de l'MPIc) *Scatter*:

```
int MPI_Scatter (void *sendbuf, int sendcount,
                MPI_Datatype sendtype, void *recvbuf,
                int recvcount, MPI_Datatype recvtype,
                int root, MPI_Comm comm)
```

Es pot utilitzar `MPI_Scatter` per a enviar des d'un procés missatges diferents a la resta de processos. En el procés arrel el missatge es divideix en segments de mida `sendcount`, i el segment *i*-èsim s'envia al procés *i*. Si es volen enviar blocs de mides diferents als diferents processos o si els blocs a enviar no estan seguits en la memòria, es pot utilitzar la funció `MPI_Scatterv`. La figura 26 il·lustra el funcionament del `MPI_Scatter`.

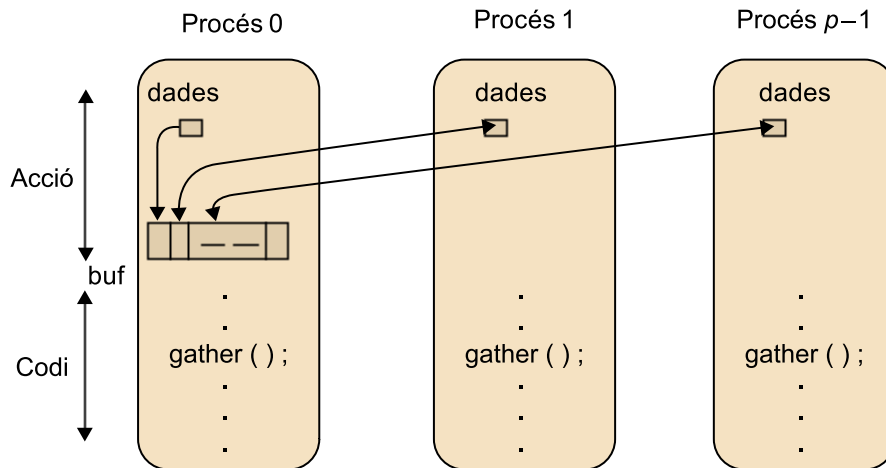
Figura 26. Funcionament de l'operació *scatter* de l'MPId) *Gather*:

```
int MPI_Gather (void *sendbuf, int sendcount,
                MPI_Datatype sendtype, void *recvbuf,
                int recvcount, MPI_Datatype recvtype,
```

```
int root, MPI_Comm comm)
```

Aquesta és la funció inversa del `MPI_Scatter`. Tots els processos (inclosa l'arrel) envien al procés arrel `sendcount` dades des de `sendbuf`, i l'arrel les emmagatzema a `recvbuf` per l'ordre dels processos. Si es vol que cada procés envii blocs de mides diferents s'ha d'utilitzar `MPI_Gatherv`. La figura 27 il·lustra el funcionament de `MPI_Gather`.

Figura 27. Funcionament de l'operació *gather* de l'MPI



Per a enviar blocs de dades de tots a tots els processos s'utilitza:

```
int MPI_Allgather (void *sendbuf, int sendcount,
                  MPI_Datatype sendtype, void *recvbuf,
                  int recvcount, MPI_Datatype recvtype,
                  MPI_Comm comm)
```

en què el bloc enviat per l'i-èsim procés s'emmagatzema com a bloc i-èsim en `recvbuf` en tots els processos. Per a enviar blocs de mides diferents s'utilitza `MPI_Allgatherv`.

Per a enviar blocs de dades diferents als diferents processos s'utilitza:

```
int MPI_Alltoall(void *sendbuf, int sendcount,
                MPI_Datatype sendtype, void *recvbuf,
                int recvcount, MPI_Datatype recvtype,
                MPI_Comm comm)
```

en què, de cada procés i, el bloc j s'envia al procés j, que l'emmagatzema com a bloc i en `recvbuf`. Hi ha el corresponent `MPI_Alltoallv`.

e) Reducció:

```
int MPI_Reduce(void *sendbuf, void *recvbuf, int count,
              MPI_Datatype datatype, MPI_Op op,
```

```
int root, MPI_Comm comm)
```

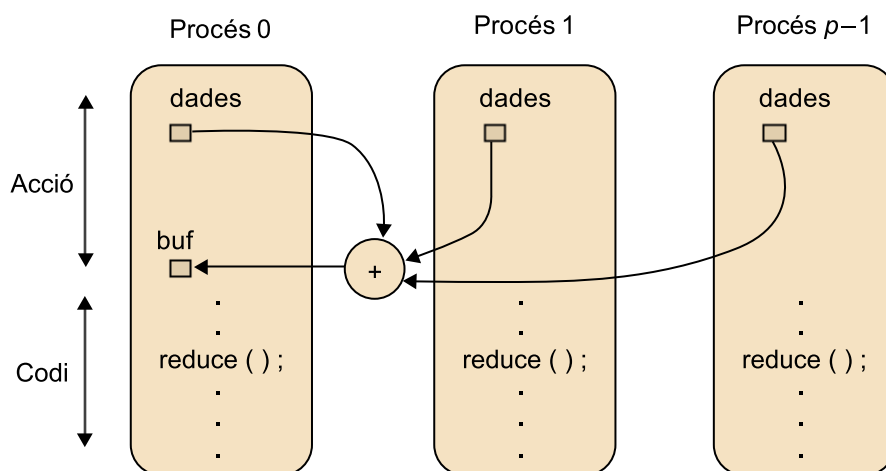
Du a terme una reducció de tots a un, és a dir, es recullen els valors distribuïts en els diferents processos i s'aplica una operació a tots ells. Això vol dir que les dades sobre les quals es farà la reducció es prenen de la zona apuntada per `sendbuf` i el resultat es deixa en l'apuntada per `recvbuf`. Si el nombre de dades (`count`) és superior a 1, l'operació es fa element a element en els del `buffer` per separat. El resultat es deixa en el procés `root`. L'operació que s'aplica a les dades és determinada per `op`. Les operacions que s'admeten per a fer reduccions es mostren en la taula següent:

Taula 5. Tipus de dades MPI per a C i Fortran

Operació	Significat	Tipus C permesos
MPI_MAX MPI_MIN MPI_SUM MPI_PROD	màxim mínim suma producte	Enters i punt flotant Enters i punt flotant Enters i punt flotant Enters i punt flotant
MPI_LAND MPI_LOR MPI_LXOR	AND lògic OR lògic XOR lògic	Enters Enters Enters
MPI_BAND MPI_BOR MPI_BXOR	AND bit a bit OR bit a bit XOR bit a bit	Enters i bytes Enters i bytes Enters i bytes
MPI_MAXLOC MPI_MINLOC	Màxim i localització Mínim i localització	Parelles de tipus Parelles de tipus

La figura 28 il·lustra el funcionament de `MPI_Reduce` i a continuació es mostra un exemple d'utilització d'aquesta operació en el qual s'utilitza per a sumar els valors d'un interval a partir de les sumes parcials de cadascun dels processos MPI i reducció.

Figura 28. Funcionament de l'operació de reducció de l'MPI



```
#include "mpi.h"
#include <stdio.h>
int main(int argc, char **argv)
```

```

{
    int MyProc, tag=1, size;
    char msg='A', msg_recpt ;
    MPI_Status *status ;
    int root ;
    int left, right, interval ;
    int number, start, end, sum, GrandTotal;
    int mystart, myend;
    MPI_Init(&argc, &argv);
    MPI_Comm_rank(MPI_COMM_WORLD, &MyProc);
    MPI_Comm_size(MPI_COMM_WORLD, &size);
    root = 0;
    if (MyProc == root)      /* El procés arrel llegeix els límits de l'interval */
    {
        printf("Give the left and right limits of the interval\n");
        scanf("%d %d", &left, &right);
        printf("Proc root reporting : the limits are : %d %d\n", left, right);
    }
    MPI_Bcast(&left, 1, MPI_INT, root, MPI_COMM_WORLD); /*Bcast limits a tots*/
    MPI_Bcast(&right, 1, MPI_INT, root, MPI_COMM_WORLD);
    if (((right - left + 1) % size) != 0)
        interval = (right - left + 1) / size + 1 ; /*Fixa límits locals de suma*/
    else
        interval = (right - left + 1) / size;
    mystart = left + MyProc*interval ;
    myend = mystart + interval ;
    /* estableix els límits dels intervals correctament */
    if (myend > right) myend = right + 1 ;
    sum = root;          /* Suma localment a cada procés MPI */
    if (mystart <= right)
        for (number = mystart; number < myend; number++) sum = sum + number ;
    /* Fa la reducció al procés arrel */
    MPI_Reduce(&sum, &GrandTotal, 1, MPI_INT, MPI_SUM, root, MPI_COMM_WORLD) ;
    MPI_Barrier(MPI_COMM_WORLD);
    /* El procés arrel retorna els resultat */
    if(MyProc == root)
        printf("Proc root reporting : Grand total = %d \n", GrandTotal);
    MPI_Finalize();
}

```

Quan tots els processos han de rebre el resultat de l'operació s'utilitza la funció:

```

int MPI_Allreduce (void *sendbuf, void *recvbuf, int count,
                  MPI_Datatype datatype, MPI_Op op, MPI_Comm comm)

```

3.1.4. Compilació i execució

En aquest primer exemple veurem com es compila i s'executa un codi MPI. La manera de fer-ho pot variar d'una implementació a una altra. Una manera normal de compilar és la següent (per MPICH, per exemple).

```
mpicc programa.c -o programa
```

en què `mpicc` crida el compilador de C que estigui establert i el munta amb la biblioteca de l'MPI. S'indica el codi a compilar i les opcions que tradicionalment es passen al compilador C.

Un cop generat el codi, la manera d'executar-lo també depèn de la compilació. El més habitual és utilitzar `mpirun` passant-hi el codi a executar i una sèrie d'arguments que indiquen els processos a posar en marxa i la distribució dels processos en els processadors. Per exemple, si es crida de la manera següent:

```
mpirun -np 4 programa
```

es posen en marxa quatre processos (`-np 4`) en què tots executen el mateix codi. No es diu a quins processadors s'assignen els processos, per la qual cosa s'utilitzarà l'assignació per defecte que estigui establerta en el sistema: pot ser que tots els processos s'assignin al mateix node (que pot ser monoprocesador o tenir varis nuclis) o que, si hi ha diferents nodes en el sistema, s'assigni un procés a cada node.

Hi ha diferents maneres d'indicar com s'assignen els processos en executar el programa. Per exemple, mitjançant:

```
mpirun -np 4 -machinefile nodes.txt programa
```

s'indica que els processos s'assignin segons s'indica en el fitxer `nodes.txt`, el qual assenyala els servidors que es poden fer servir per a executar els processos MPI. Es tracta d'una llista amb els noms dels servidors, com per exemple:

```
node1  
node2
```

En aquest exemple, i si els estem executant en el `node0`, els quatre processos s'assignaran en l'ordre següent: el procés 0 al node en què estem executant el programa (`node0`), el procés 1 al `node1`, el 2 al `node2`, i una vegada s'hagi acabat la llista de màquines del fitxer es comença una altra vegada pel principi, de manera que el procés 3 s'assigna al `node1`. La forma d'assignació depèn de la implementació, i pot ser que no s'assigni cap procés al node on som o que aquest entri en l'assignació cíclica un cop acabada la llista de màquines del fitxer. Per exemple, amb l'opció `-nolocal` podem evitar que s'utilitzi el

servidor des del qual estem executant el programa (cosa que pot ser útil, per exemple, en el cas d'estar utilitzant un node *login* que no s'hauria de fer servir per a la computació).

3.2. Llenguatges PGAS

La programació de sistemes distribuïts moltes vegades s'ha equiparat amb el paradigma de pas de missatges. Tot i així, és possible implementar abstraccions de nivell més alt per sobre del sistema de memòria distribuïda. Durant l'última dècada, diferents grups de recerca han desenvolupat aquestes abstraccions, de la qual cosa ha resultat una família de llenguatges coneguts com a *llenguatges PGAS*¹⁵.

⁽¹⁵⁾De l'anglès *partitioned global address space*.

L'espai d'adreces global fa referència al fet que el llenguatge proporciona un únic espai de memòria per sobre de les memòries virtuals de les màquines distribuïdes. Aquests llenguatges òbviament no proporcionen memòria compartida, ja que no es pot esperar que el maquinari pugui encarregar-se de la seva coherència. Tot i així, l'espai d'adreces global permet definir estructures de dades globals, cosa que millora el model de memòria distribuïda en el qual el programador ha de recordar un conjunt d'estructures de dades independents que actuen com una estructura de dades distribuïda difícil de mantenir. Amb els llenguatges PGAS, els programadors es poden centrar en el comportament d'un únic procés i poden distingir les dades locals de les dades no locals (remotes), però aquests llenguatges permeten simplificar la programació, ja que eliminen els detalls del pas de missatges, i millora així la productivitat del programador. Són els compiladors els que generen totes les crides de comunicació quan hi ha referències a adreces de memòria que no són locals.

Les implementacions de llenguatges PGAS més importants i que a continuació estudiarem són UPC, Co-Array Fortran i Titanium, les quals estenen C, Fortran i Java, respectivament.

3.2.1. UPC

L'UPC¹⁶ va ser desenvolupat prop de l'any 2000. L'UPC proporciona al programador una visió global de l'espai de memòria en què quan els elements (per exemple matrius) es defineixen com a compartits (*shared*) es distribueixen en les memòries dels diferents processos, de tal manera que els elements compartits es distribueixen cíclicament o amb blocs cíclics. Aquest tipus de distribució fa que es pugui balancejar millor la càrrega, però en general també fa que es redueixi la localitat de les dades. Tot i així, podem millorar la localitat de les dades assignant els fragments de les dades compartides directament al processador de tal manera que moltes dades contigües quedin en el mateix procés.

⁽¹⁶⁾De l'anglès *unified parallel C*.

Com que l'UPC estén el llenguatge de programació C es poden utilitzar apuntadors. Els apuntadors d'UPC poden ser o bé privats (locals respecte al procés) o bé compartits entre tots els processos. Com que els apuntadors poden ser privats o compartits i, a més, poden apuntar a dades privades o compartides, hi ha quatre tipus d'apuntadors tal com mostra la figura 29.

Figura 29. Tipus d'apuntador en l'UPC

		Propietat de l'apuntador	
		<i>Private</i>	<i>Shared</i>
Propietat de la referència	<i>Private</i>	<i>Private-private</i> , p1	<i>Private-shared</i> , p2
	<i>Shared</i>	<i>Shared-private</i> , p3	<i>Shared-shared</i> , p4

Aquestes propietats estan associades als tipus del llenguatge, com es mostra a continuació, en què es poden veure diversos exemples de declaració d'apuntadors en UPC:

```
int *p1; /* apuntador privat apuntant a dades locals */
shared int *p2; /* apuntador privat apuntant a l'espai compartit */
int *shared p3; /* apuntador compartit apuntant a dades locals */
shared int *shared p4; /* apuntador compartit apuntant a l'espai compartit */
```

El codi de la suma de vectors que es mostra a continuació il·lustra l'ús del segon cas (apuntador privat apuntant a l'espai compartit).

```
shared int v1[N], v2[N], v1v2sum[N];

void main()
{
    int i;
    shared int *p1, *p2;
    p1=v1;
    p2=v2;
    upc_forall(i=0; i<N; i++; p1++; p2++; i)
    {
        v1v2sum[i]=*p1+*p2;
    }
}
```

El codi anterior també il·lustra l'ús d'un altre element bàsic de l'UPC, el `upc_forall`. Aquesta abstracció distribueix les iteracions de bucles de C normal (`for`) en els processos utilitzant una clàusula d'afinitat (que és la quarta clàusula del `upc_forall`). La clàusula d'afinitat del bucle indica on s'han executat les iteracions del bucle. En l'exemple anterior, el procés que executa

la iteració i -èsima és la que està associada a i (el propietari). El `upc_forall` és una operació global en què la major part del codi s'executa de manera local en un procés.

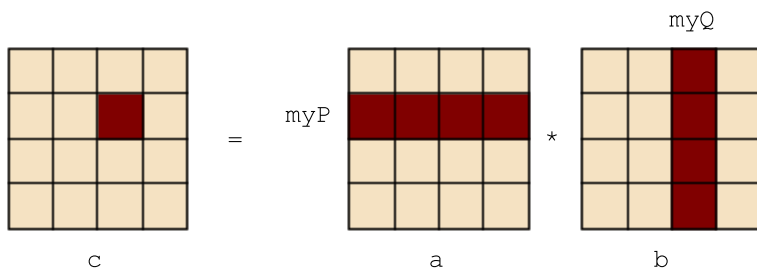
3.2.2. Co-Array Fortran

Co-Array Fortran (CAF) és una extensió de Fortran desenvolupada al final dels anys noranta. El CAF es caracteritza per la seva elegància i simplicitat. La idea principal és afegir una extensió al llenguatge anomenada `co-array`. El `co-array` és el mecanisme per a la comunicació entre processadors. Aprofitant que el Fortran utilitza parèntesis per a referir-se a les matrius, el CAF afegeix claudàtors als noms de variables per a referir-se a co-arrays. Per exemple, el codi que es mostra a continuació defineix tres matrius amb co-arrays.

```
real, dimension(n,n) [p,*]::a,b,c
...
do k=1,n
  do q=1,p
    c(i,j) [myP,myQ]=c(i,j) [myP,myQ]+a(i,k) [myP,q]*b(k,j) [q,myQ]
  enddo
enddo
```

L'especificació dels co-arrays indica que la memòria d'una variable que es distribueix per tots els processos s'anomena *imatge* en CAF. Per tant, en l'exemple anterior, a cada procés hi ha assignats fragments de les matrius a , b i c . A més, el fet que el co-array es defineixi per tenir dues dimensions fa que els processos es puguin organitzar en l'estructura lògica de dues dimensions de p files i q columnes. Cal tenir en compte que el símbol `*` indica que cada fila de processos s'ha d'omplir al màxim, per tant, si `num_images()=pq`, llavors hi ha p files i q columnes de processos. Cada procés trobarà la seva ubicació en l'organització de dues dimensions cridant la rutina `this_image()` i definint els dos paràmetres, `myP` i `myQ`, per a indicar un fragment concret de la matriu. La gestió de la partició de dades, inclosa la inicialització dels seus valors, és responsabilitat del programador. En la figura 30 es mostra un exemple en què `myP=2` i `myQ=3`.

Figura 30. Exemple d'organització de co-array amb CAF



El codi anterior mostra que cada procés, definit per $[myP, myQ]$, actualitza el seu fragment de la matriu quan calcula el producte escalar. Per a indicar les files s'accedeix a les dades en els processos $[myP, q]$, i per a indicar les columnes s'accedeix a les dades en els processos $[q, myQ]$. Les crides de comunicació per a aquests accessos remots les genera el compilador, cosa que fa que la programació se simplifiqui substancialment.

3.2.3. Titanium

El llenguatge Titanium (Ti) és una extensió del Java que es pot executar en computadors de memòria distribuïda. Té un model de memòria similar a l'UPC i, de la mateixa manera que altres llenguatges PGAS, genera codi de comunicació entre nodes basant-se en una biblioteca de comunicacions unidireccional. Una de les qualitats que el diferencia dels altres llenguatges PGAS que hem vist és el fet que és orientat a objectes. També es diferencia del Java pel fet d'afegir les anomenades *regions*, les quals suporten una gestió de memòria d'altres prestacions com a alternativa al mecanisme de *garbage collection*. Hi ha altres funcionalitats del Java que s'han restringit o limitat, però les funcionalitats més típiques són suportades. Per exemple, s'han afegit matrius de dues dimensions per tal de fer el Ti més apropiat per a la computació científica.

Una prestació d'eficiència important del Ti és la iteració desordenada, *foreach*. Aquesta simplifica la tasca tant del programador com del compilador i funciona permetent iterar sobre el domini d'una variable. A continuació es mostra com s'utilitza per a multiplicar matrius.

```
public static void matMul( double [2d] a,
                          double [2d] b,
                          double [2d] c)
{
    foreach (ij in c.domain())
    {
        double [1d] aRowi=a.slice(1, ij[1]);
        double [1d] bColj=b.slice(2, ij[2]);
        foreach (k in aRowi.domain())
        {
            c[ij]+=aRowi[k]*bColj[k];
        }
    }
}
```

D'aquesta manera *foreach*, al contrari de *forall*, permet la concurrència sobre múltiples índexs dins d'un mateix bloc.

El Titanium força la sincronització global mitjançant barreres i el concepte de variable compartida entre tots els processos que s'anomenen *single*. Per exemple, en una simulació és força habitual que cada procés faci els seus càl-

culs sobre dades locals i periòdicament coordini les seves accions. La barreira assegura que tots els processos paren els seus càlculs al mateix temps per a actualitzar o llegir de la memòria. La utilització de variables tipus `single` garanteix que tots els processos estiguin en la mateixa fase del còmput. Per exemple, una simulació de partícules podria utilitzar aquest concepte tal com es mostra a continuació:

```
int single stepCount=0;
int single endCount=100;
for (; stepCount<endCount; stepCount++)
{
    Llegeix partícules remotes
    Calcula les forces de la meva
    Ti.barrier();
    Escriu les meves partícules utilitzant noves forces
    Ti.barrier();
}
```

4. Esquemes algorítmics paral·lels

En aquest apartat analitzarem esquemes algorítmics dels que més s'utilitzen en la resolució de problemes en entorns paral·lels. Estudiarem un conjunt d'esquemes de referència que constitueix la base per a la resolució d'una gamma ampla de problemes reals.

4.1. Paral·lisme de dades

Amb aquesta tècnica es treballa amb moltes dades que es tractaran de manera igual o similar. Típicament es tracta d'algoritmes numèrics en què les dades es troben en vectors o matrius i es treballa amb aquests fent un mateix processament. La tècnica és apropiada per a memòria compartida i normalment s'obté el paral·lisme dividint el treball dels bucles entre els diferents fluxos (per exemple, en OpenMP `#pragma omp for`). Els algoritmes en què apareixen bucles on no hi ha dependències de dades, o aquestes dependències es poden evitar, són adequats per a paral·lismes de dades. Si es distribueixen les dades entre els processos treballant en memòria distribuïda parlem de partició de dades.

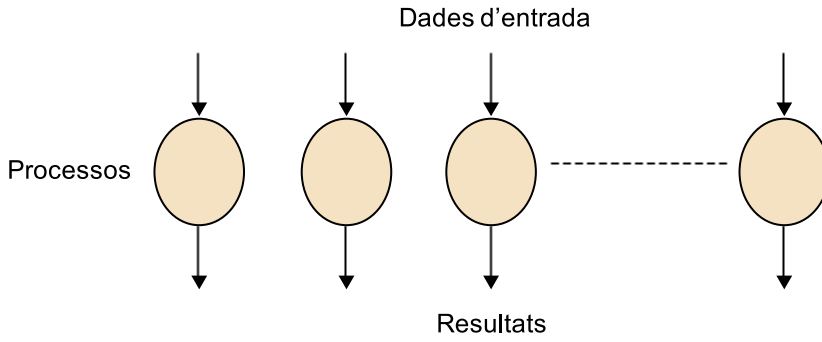
Aquest tipus de paral·lisme és el que hi ha en la multiplicació de matrius en l'OpenMP, ja que en aquest cas es paral·lilitza la multiplicació simplement indicant amb `#pragma omp parallel for` en el bucle més extern que cada flux calculi una sèrie de files de la matriu resultat. L'aproximació d'una integral també segueix aquest paradigma.

La suma seqüencial de n nombres es fa amb un únic bucle `for`, el qual es pot paral·lilitzar amb `#pragma omp parallel for`, amb la clàusula `reduction` per a indicar la manera en què els fluxos accedeixen a la variable en què s'emmagatzema la suma. La paral·lilització es fa d'una manera molt simple, i es parla de paral·lisme implícit, ja que el programador no és el responsable de la distribució del treball ni de l'accés a les variables compartides.

En alguns casos ens pot interessar (per exemple, per a disminuir els costos de gestió dels fluxos) posar en marxa un treball per a cadascun dels fluxos, i que el programador s'encarregui de la distribució i accés a les dades, amb la qual cosa es complica la programació. En aquest cas parlem de paral·lisme explícit.

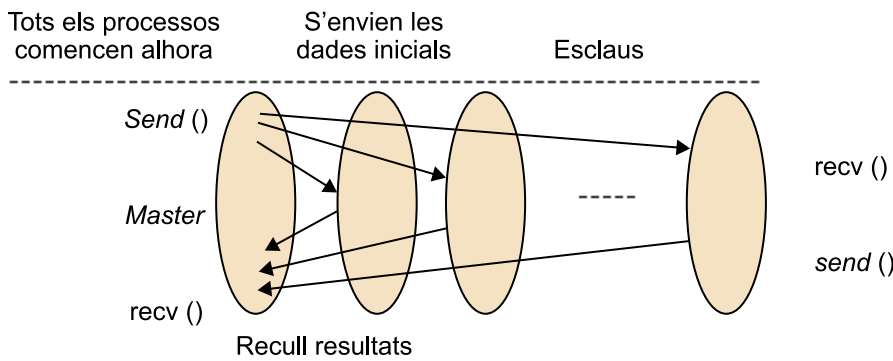
Aquest tipus de problemes que estan compost d'un conjunt de tasques independents, com mostra la figura 31, també es coneix com a *embarrassingly parallel*. En aquest tipus d'aplicacions no hi ha o hi ha molt poca comunicació entre processos i cada procés pot fer la seva tasca sense necessitar cap interacció amb els altres processos.

Figura 31. Model d'execució de múltiples processos amb tasques independents



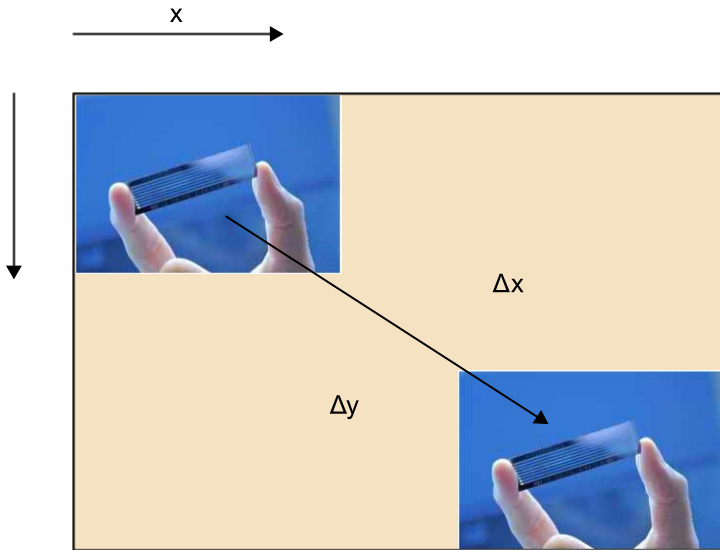
Si pensem en una implementació estàtica amb l'MPI d'aquest model basat en el patró mestre-esclau, veurem que el procés mestre crea un conjunt de processos esclaus que se sincronitzen mitjançant el pas de missatges al principi, fan la seva tasca independent i finalment es tornen a sincronitzar al final quan el procés mestre recull els resultats dels altres processos, com il·lustra la figura 32.

Figura 32. Implementació estàtica de tasques independents amb l'MPI basada en mestre-esclau



Tal com hem comentat anteriorment, aplicacions com ara el processament d'imatges són típiques d'aquest model. Un exemple clar és el desplaçament d'una imatge en el qual cada píxel s'ha de desplaçar un increment concret, com mostra la figura 33.

Figura 33. Exemple d'aplicació de desplaçament d'una imatge



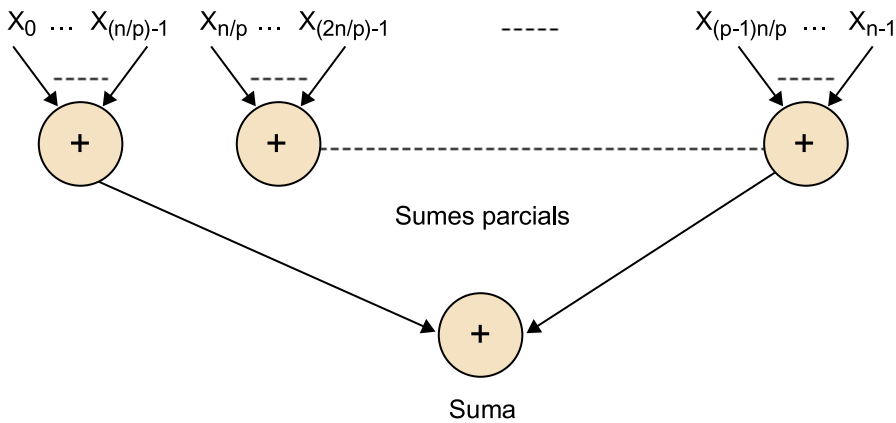
4.2. Partició de dades

En aquest tipus d'aplicacions es treballa amb volums de dades grans que normalment estan en vectors o matrius i s'obté paral·lisme dividint el treball en zones diferents de les dades entre diferents processos. A diferència del paral·lisme de dades, en aquest cas s'han de distribuir les dades entre els processos, i aquesta distribució determina el repartiment del treball. El més normal és dividir l'espai de dades en regions i que en l'algoritme hi hagi un intercanvi de dades entre regions adjacents, per la qual cosa s'han d'assignar regions adjacents als processadors intentant que les comunicacions no siguin gaire costoses. Hi ha comunicacions que comparteixen dades i que estan distribuïdes en la memòria del sistema, per la qual cosa hi ha un cost de comunicacions que sol ser força elevat, de manera que per a obtenir bones prestacions és necessari que el volum de computació sigui molt més gran que el de comunicació.

Una manera concreta de partir dades és la de dividir i conquerir. Aquest tipus de problemes es caracteritzen pel fet de dividir el problema en subproblemes com el problema principal. Normalment es divideix en encara més subproblemes de manera recursiva.

La figura 34 mostra un exemple de suma d'una seqüència de nombres mitjançant un patró model de dividir i conquerir en el qual cada procés fa una suma parcial que finalment és agregada per tal d'obtenir el resultat final.

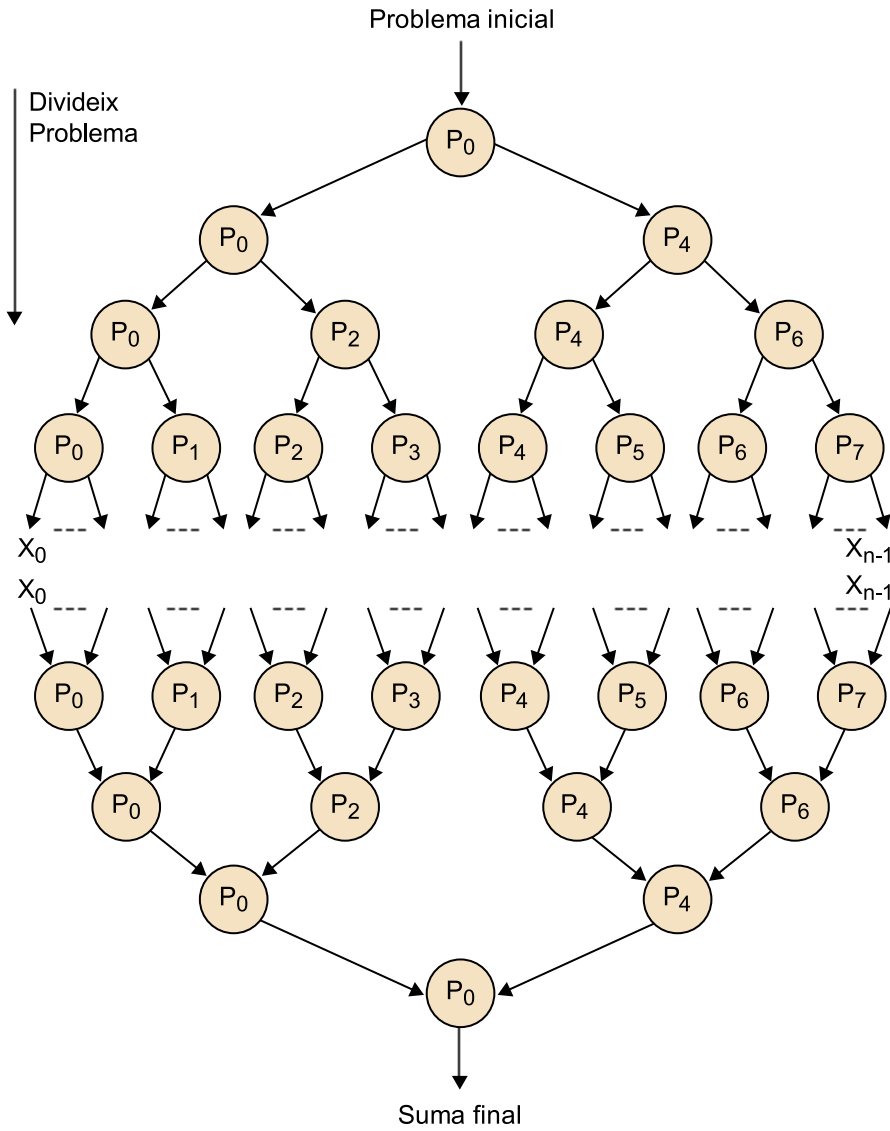
Figura 34. Exemple de suma d'una seqüència de nombres mitjançant un patró model de dividir i conquerir



4.3. Esquemes paral·lels en arbre

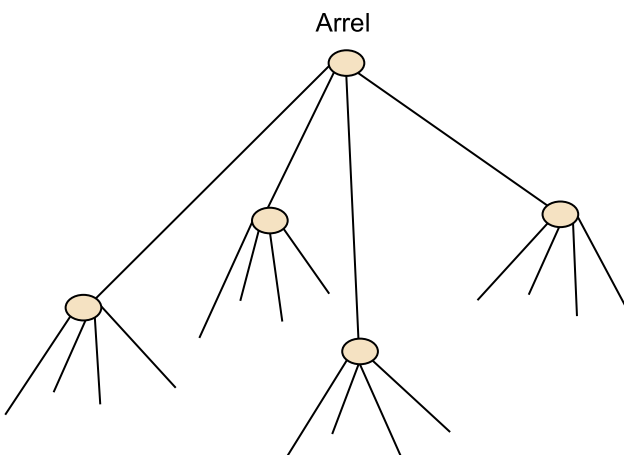
Molts problemes tenen una representació en forma d'arbre o graf i la seva solució s'obté recorrent l'arbre o graf fent computacions. Una tècnica per a desenvolupar programes paral·lels és, donat un problema, obtenir el graf de dependències que representa les computacions a fer (nodes del graf) i les dependències de dades (arestes), i a partir del graf decidir l'assignació de tasques als processos i les comunicacions entre aquests, que estaran determinades per les arestes que comuniquen nodes assignats a processos diferents. Un exemple és la suma prefixa, com mostra la figura 35. Veiem com en la primera fase es divideix el problema a mesura que es generen fills i, en la segona, es recullen els resultats a mesura que es recorre l'arbre.

Figura 35. Exemple de suma prefixa mitjançant l'esquema paral·lel en arbre



Els arbres també poden tenir més de dos fills, com il·lustra la figura 36, en què els nodes tenen quatre fills.

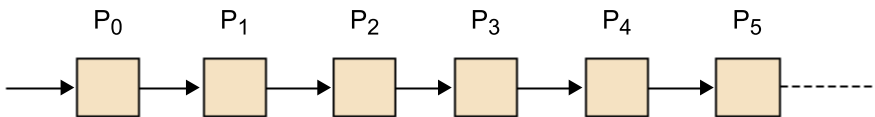
Figura 36. Esquema paral·lel en arbre amb quatre fills



4.4. Computació en *pipeline*

Amb el patró *pipeline* es resol un problema descomponent-lo en una sèrie de tasques successives, de manera que les dades flueixen en una direcció per l'estructura de processos. De fet, es pot veure com una forma de descomposició funcional en què el problema es divideix en diferents funcions que s'han de dur a terme successivament, com mostra la figura 37. Té una estructura lògica de pas de missatges, ja que entre les tasques consecutives s'han de comunicar dades. Pot tenir interès quan no hi ha un únic conjunt de dades a tractar, sinó una sèrie de conjunts de dades, que passen a ser computades una darrera l'altra, o bé problemes en què hi ha un paral·lelisme funcional, amb operacions de diferent tipus sobre el conjunt de dades i que s'han d'executar una darrera l'altra.

Figura 37. Esquema de *pipeline* senzill



Per exemple, suposem que hem de sumar els elements d'un vector mitjançant un algoritme com el següent:

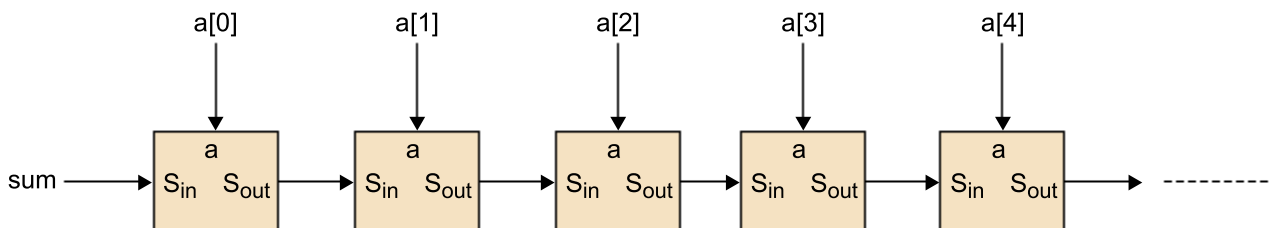
```
for (i = 0; i < n; i++)
    sum = sum + a[i];
```

i podem desplegar el bucle de la manera següent:

```
sum = sum + a[0];
sum = sum + a[1];
sum = sum + a[2];
sum = sum + a[3];
sum = sum + a[4];
...
```

D'aquesta manera ens trobem amb el *pipeline* que mostra la figura 38:

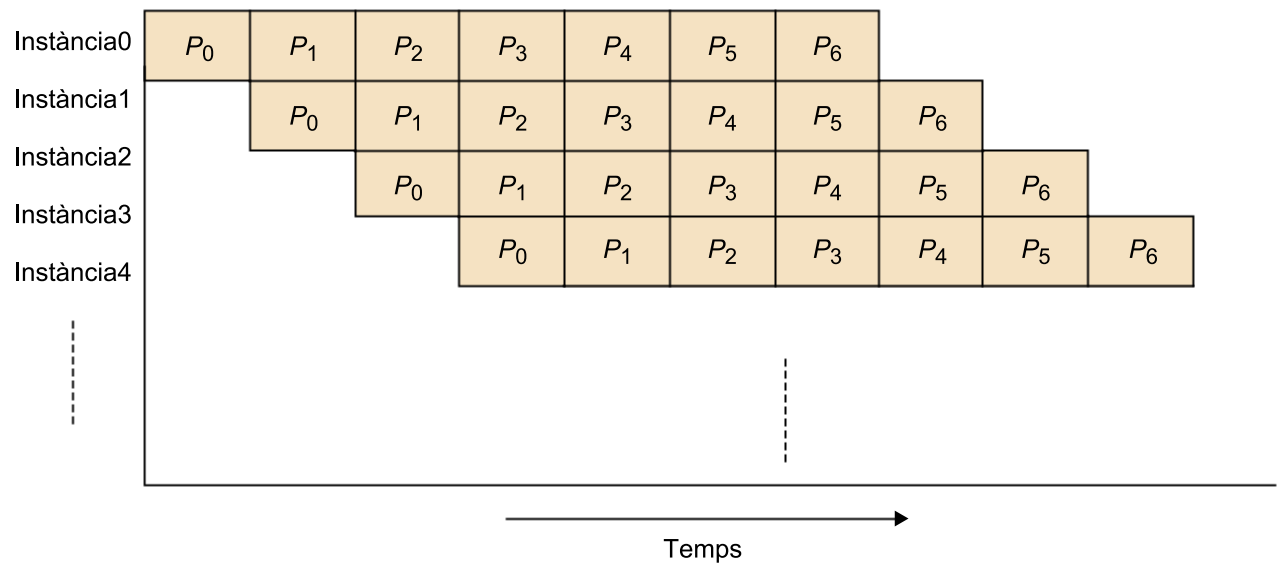
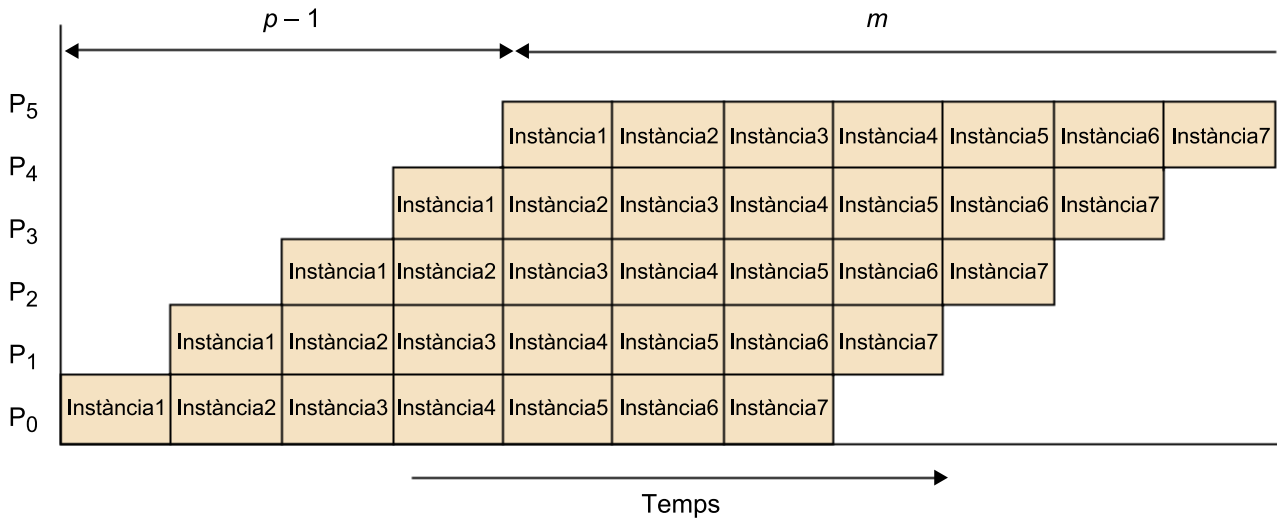
Figura 38. *Pipeline* de l'exemple de la suma dels elements d'un vector



Si suposem que, en general, el problema es pot dividir en una sèrie de tasques seqüencials, l'esquema de *pipeline* pot proporcionar un temps d'execució més curt en els tipus de computacions següents:

1) Si més d'una instància del problema complet es pot executar alhora, tal com mostra el diagrama de la figura 39.

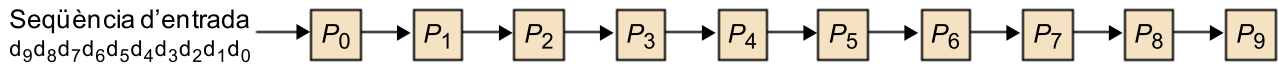
Figura 39. Diagrama espaciotemporal del primer tipus de *pipeline*



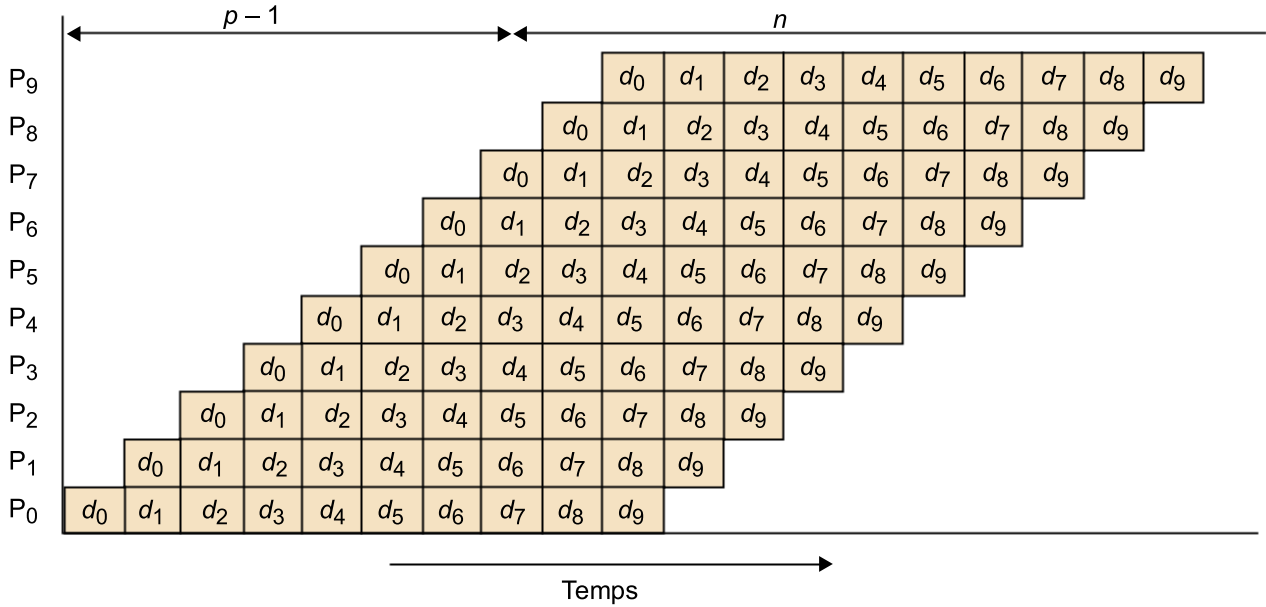
2) Si una sèrie d'elements de dades s'ha de processar i cada una requereix múltiples operacions, com mostra la figura 40.

Figura 40. Diagrama espaciotemporal del segon tipus de pipeline

(a) Estructura pipeline



(b) Diagrama temporal

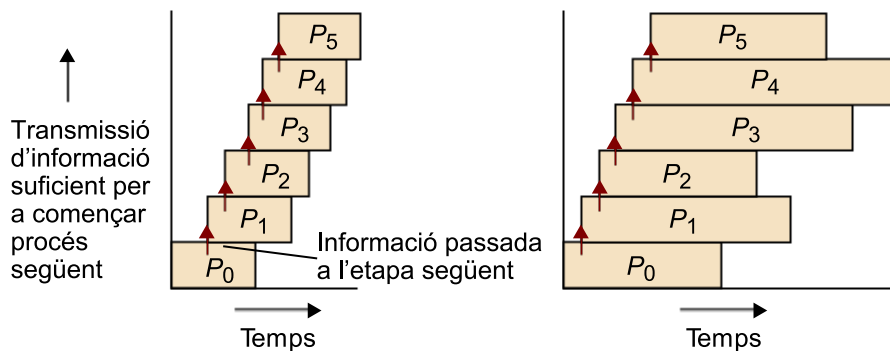


3) Si la informació necessària per a començar l'execució del procés següent es pot passar al procés següent abans que el procés actual hagi acabat les seves operacions tal com il·lustra la figura 41. Un exemple d'aquest tipus és la cerca d'una certa informació en un conjunt de dades.

Figura 41. Diagrama espaciotemporal del tercer tipus de pipeline

(a) Processos amb mateix temps d'execució

(b) Processos amb temps d'execució diferents



4.5. Esquema mestre-esclau

En el paradigma mestre-esclau tenim un flux/procés diferenciats, anomenat *mestre*, que posa en marxa els altres processos, anomenats *esclaus*, els assigna treball i recull les solucions parcials que generen.

En alguns casos es parla de *granja de processos* quan un conjunt de processos treballa de manera conjunta però independentment en la resolució d'un problema. Aquest paradigma pot tenir similituds amb el mestre-esclau si considerem que els processos que constitueixen la granja són els esclaus (o els esclaus i el mestre si és que aquest també intervé en la computació).

Un altre paradigma relacionat amb aquests dos és el dels treballadors replicats. En aquest cas, els treballadors (fluxos o processos) actuen de manera autònoma, executen tasques que possiblement donen lloc a noves tasques que executaria el mateix treballador o un altre. Es gestiona una bossa de tasques, cada treballador pren una tasca, que s'inclou en la bossa, i una vegada acabada la tasca que té assignada en pren una altra de la bossa per a treballar-hi, i així mentre queden tasques per a executar. Es planteja el problema de la terminació, ja que és possible que un treballador sol·liciti una tasca i a la bossa no n'hi hagi més, però això no vol dir que no quedin tasques per executar, ja que és possible que un altre procés estigui computant i generi noves tasques. La condició per a finalitzar serà que no quedin tasques en la bossa i que no hi hagi treballadors treballant. Hi ha diferents maneres d'abordar el problema de la terminació.

Aquesta és la forma típica de treball en OpenMP. Inicialment treballa un únic flux, i al arribar a un constructor paral·lel posa en marxa una sèrie de fluxos esclaus dels que es converteix en el mestre.

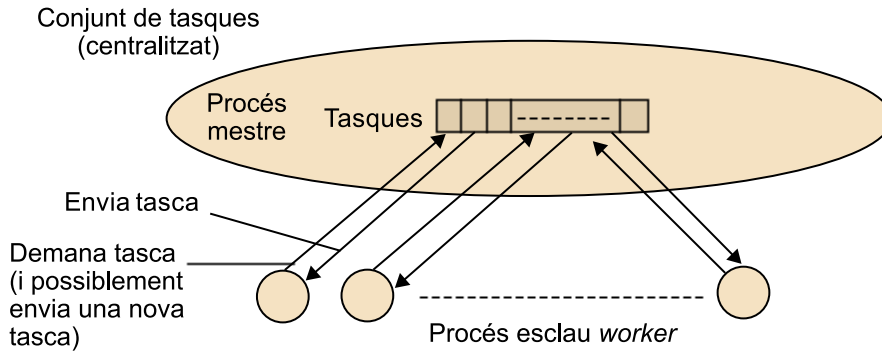
En programació per pas de missatges, també s'acostuma a utilitzar aquest patró. En MPI, en posar-se en marxa els processos amb una ordre com `mpirun -np 8 programa`, s'inicialitzen vuit processos tots iguals, però l'entrada i sortida se sol fer mitjançant un procés, que actua com a mestre, que genera el problema, distribueix les dades, darrera del qual comença la computació, en la qual el mestre pot intervenir o no, i al final el mestre rep els resultats dels esclaus.

L'assignació del treball als esclaus es pot fer de manera estàtica o dinàmica:

- En l'assignació estàtica el mestre decideix els treballs que assigna als esclaus i fa l'enviament.
- En l'assignació dinàmica el mestre genera els treballs i els emmagatzema en una bossa de tasques que s'encarrega de gestionar. Els esclaus demanen treball de la bossa de tasques a mesura que queden lliures per a executar noves tasques. D'aquesta manera s'equilibra la càrrega dinàmicament, però

hi pot haver una sobrecàrrega de la gestió de la bossa i aquesta pot convertir-se en un coll d'ampolla. A més, en alguns problemes, els esclaus en executar una tasca generen noves tasques que s'han incloure en la bossa, de manera que es generen més comunicacions amb el mestre. Un exemple d'això es mostra en la figura 42.

Figura 42. Exemple d'assignació dinàmica mitjançant un conjunt de tasques centralitzat



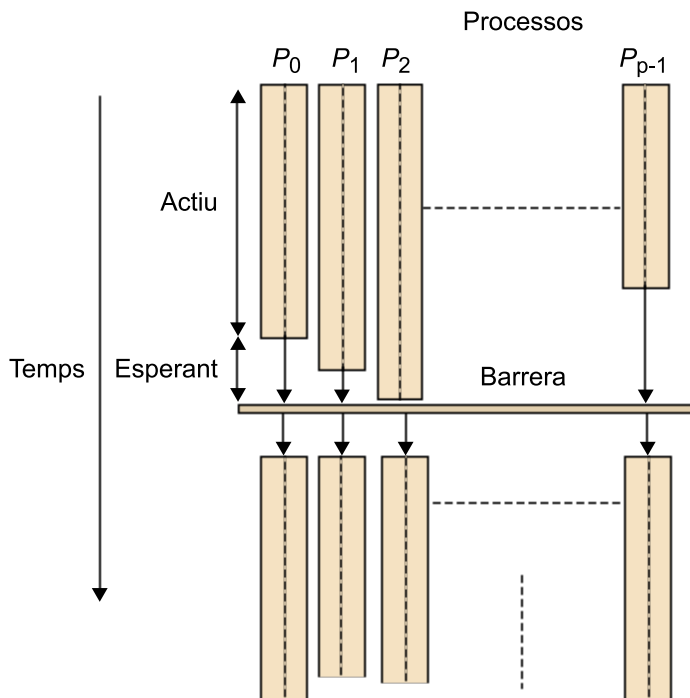
4.6. Computació síncrona

Parlem de *paral·lisme síncron* quan es resol un problema per iteracions successives en què tots els processos se sincronitzen. Les característiques generals són:

- Cada procés du a terme el mateix treball sobre una porció diferent de les dades.
- Part de les dades d'una iteració s'utilitza en la següent, de manera que al final de cada iteració hi ha una sincronització, que pot ser local o global.
- El treball acaba quan es compleix algun criteri de convergència, per al qual normalment es necessita sincronització global.

Hi ha molts problemes que es resolen amb aquest tipus de patró. El mecanisme principal és la utilització de barreres allà on els processos s'han de sincronitzar. Els processos poden continuar l'execució quan tots arriben a aquest punt, com mostra la figura 43.

Figura 43. Processos arribant a una barrera en diferents moments



Activitat

Busqueu possibles implementacions de barreres tant per a memòria compartida com distribuïda. Alguns exemples típics de possibles implementacions són:

- centralitzada (lineal)
- en forma d'arbre
- butterfly*

5. Entorns d'execució d'aplicacions paral·leles

Els sistemes d'altres prestacions actuals són principalment clústers de computadors interconnectats amb xarxes d'altres prestacions. En aquests sistemes normalment hi ha múltiples usuaris que volen accedir als recursos per tal d'executar programes (paral·lels) desenvolupats mitjançant models de programació com els que hem vist anteriorment (per exemple, MPI). Per a organitzar l'accés dels diferents usuaris als recursos i gestionar aquests eficientment (per exemple, per a maximitzar-ne la utilització), els sistemes d'altres prestacions disposen de sistemes de gestió de recursos basats en cues.

De fet, en sistemes d'altres prestacions no té sentit proporcionar als usuaris accés als recursos de manera interactiva, ja que seria un caos: hi hauria conflictes en la utilització dels recursos, poca eficiència en la utilització i compartició de recursos. A més, els usuaris no s'han de preocupar de problemes relacionats amb la gestió dels computadors, com ara el balanceig de la càrrega. Així doncs, els sistemes de computació d'altres prestacions es basen en aplicacions que es poden emmagatzemar en una cua i executar posteriorment en segona pla. Aquests sistemes s'anomenen tradicionalment *sistemes per lot (batch)*.

La figura 44 mostra un esquema d'un clúster d'altres prestacions amb els components essencials per a la seva gestió. També cal tenir en compte que en els sistemes d'altres prestacions normalment hi ha el gran gruix de nodes que s'encarreguen de la computació, però normalment no s'hi pot accedir directament per mesures de seguretat i gestió, i hi ha un conjunt de nodes (que poden ser de les mateixes característiques que la resta de nodes o diferents) que s'utilitzen com a porta d'entrada (*login nodes*). Normalment, aquests nodes d'accés permeten desenvolupar les aplicacions i disposen de sistemes complets amb eines de compilació, depuració, etc. i permeten als usuaris enviar els seus treballs als sistemes de cues per a la seva execució.

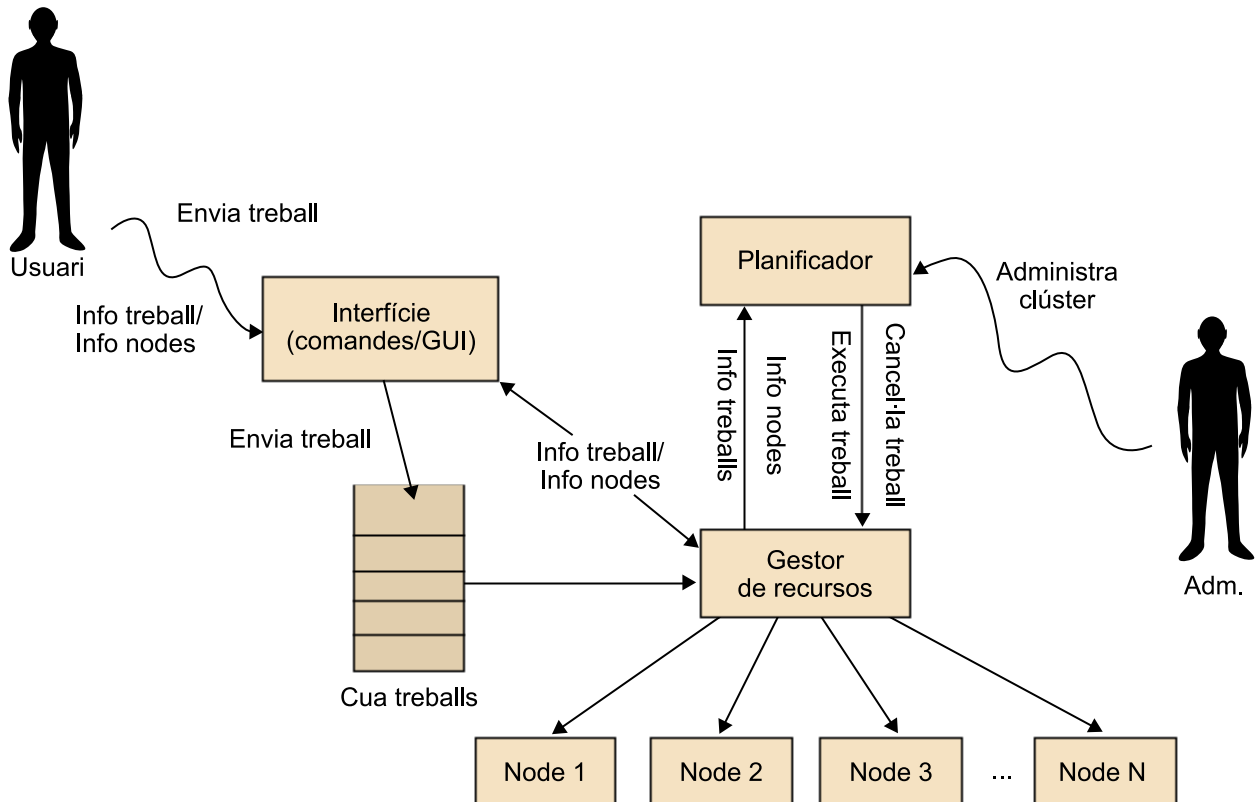
Sistemes d'altres prestacions actuals

Les característiques d'aquests sistemes són la latència reduïda i la gran amplada de banda de la xarxa d'interconnexió.

Exemples de conflictes

Alguns usuaris podrien fer un ús desmesurat dels recursos o fins i tot processos erronis podrien malgastar els recursos sense que es poguessin controlar.

Figura 44. Esquema d'un sistema clúster amb els components de gestió de recursos



5.1. Sistemes de gestió de cues

Els sistemes de gestió de cues són sistemes d'administració, planificació i execució de treballs que s'envien a un clúster. A continuació s'exposen breument els sistemes de gestió de cues PBS, que actualment és el més habitual en la computació d'altres prestacions.

5.1.1. PBS (*portable batch system*)

El PBS¹⁷ és un sistema de cues dissenyat originalment per la NASA que ha arribat a ser uns dels sistemes de cues més populars en la computació d'altres prestacions. De fet, un gran percentatge de sistemes d'altres prestacions actuals utilitzen l PBS com a sistema de cues.

⁽¹⁷⁾De l'anglès *portable batch system*.

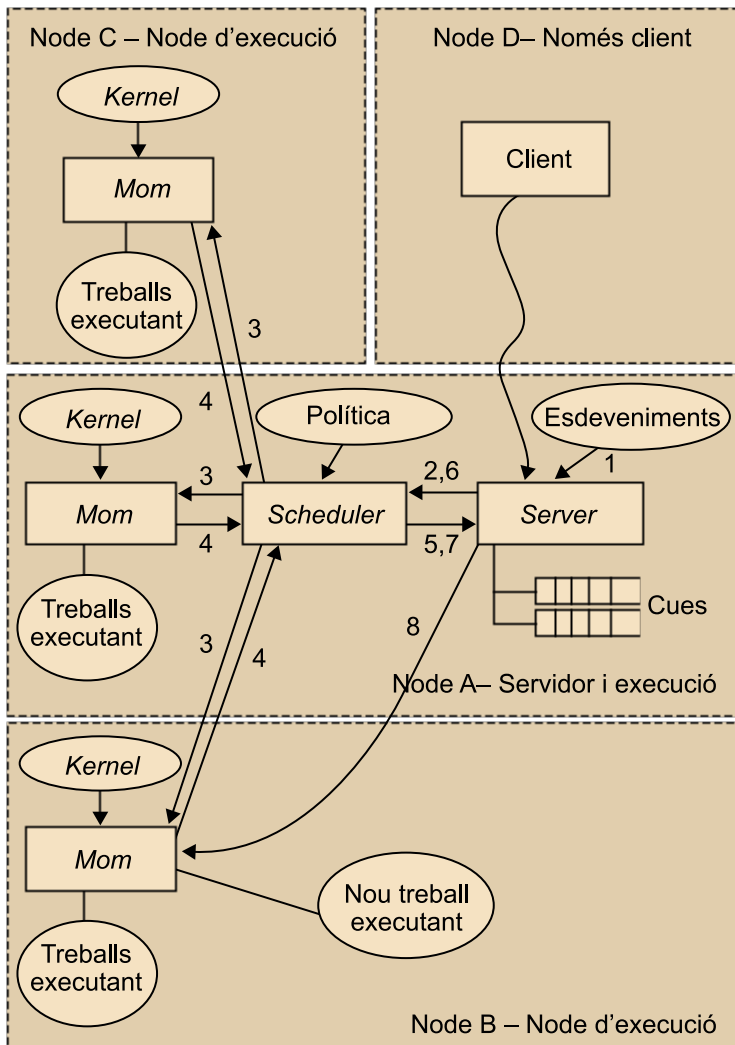
El PBS proporciona una sèrie d'eines per a la gestió de treballs per lots (*batch*), i utilitza una unitat de programació de tasques. A més, permet la planificació de treballs a diferents computadors i també definir i implementar polítiques sobre la utilització d'aquests recursos.

La figura 45 mostra l'esquema del sistema de cues, en què hi podem diferenciar els components que el formen i les interaccions entre aquests. Podem diferenciar clarament tres components bàsics:

- *server*. És un dimoni encarregat de rebre els treballs a executar i d'esperar les ordres de l'usuari.

- *mom* (o *job executor*). És un dimoni que s'executa en cadascun del nodes de còmput i que envia i rep els treballs a executar.
- *scheduler*. És un planificador que decideix quins treballs s'han d'executar en funció de la política que estigui establerta.

Figura 45. Esquema del sistema de cues PBS



Aquests components interactuen mitjançant els passos que es descriuen a continuació d'una manera cíclica:

- Un *event* diu al *server* que comenci un cicle de planificació.
- El *server* envia una petició de planificació a l'*scheduler*.
- L'*scheduler* sol·licita informació sobre els recursos als *mom*.
- Els *mom* retornen la informació sol·licitada a l'*scheduler*.
- L'*scheduler* sol·licita informació sobre el treball al *server*.
- El *server* retorna la informació sobre l'estat del treball a l'*scheduler* i pren la decisió d'executar el treball o no en funció de la política establerta.
- L'*scheduler* envia una petició d'execució al *server*.
- El *server* envia el treball al *mom* per a ser executat.

La manera en què els usuaris interactuen amb el sistema de cues és amb d'un seguit d'instruccions que es poden cridar des del sistema operatiu.

Per a executar un programa l'usuari ha de definir un treball mitjançant un fitxer *script*. A més d'indicar quin és el programa a executar i els arguments, en aquest fitxer l'usuari permet:

- Especificar si és un treball per lots o interactiu. En el cas de ser interactiu no vol dir que s'executi immediatament, sinó que quan s'executa permet a l'usuari d'interactuar amb el programa a diferència dels treballs per lots, que s'executen en segon pla sense que l'usuari hi pugui interactuar directament.
- Definir una llista amb els recursos necessaris.
- Definir la prioritat del treball per a la seva execució.
- Definir el temps d'execució, que és una informació important per a poder planificar els treballs eficientment.
- Especificar si es vol enviar un correu electrònic a l'usuari quan l'execució comença, acaba o és avortada.
- Definir dependències.
- Sincronitzar el treball amb altres treballs.
- Especificar si es vol fer *checkpointing* (en cas que el sistema operatiu ofereixi aquesta possibilitat). *Checkpointing* és especialment important en execucions molt llargues en què es vol poder monitoritzar l'aplicació durant el temps d'execució o bé per qüestions de seguretat contra possibles fallades.

La llista de recursos que es poden especificar en un fitxer de definició d'un treball permeten, a més de garantir que el treball es podrà executar, millorar la utilització del sistema, ja que el planificador pot prendre decisions més acurades. Una llista amb els recursos més comuns (no completa) es descriu a continuació:

- *cput*. Temps màxim de CPU usat per tots els processos del treball.
- *pcput*. Temps màxim de CPU usat per cadascun dels processos del treball.
- *mem*. Quantitat màxima de memòria física utilitzada per a un treball.
- *pmem*. Quantitat màxima de memòria física utilitzada per cadascun del processos d'un treball.
- *vmem*. Quantitat màxima de memòria virtual utilitzada per un treball.
- *pvmem*. Quantitat màxima de memòria virtual utilitzada per cadascun del processos d'un treball.
- *walltime*. Temps d'execució (de rellotge, no de CPU).

Decisions acurades

Un exemple seria no deixar nuclis sense feina perquè no hi ha prou memòria en un node per a assignar-hi un altre treball.

- *file*. Mida màxima de qualsevol fitxer que pot crear el treball.
- *host*. Nom dels nodes computacionals en què s'executa el treball.
- *nodes*. Nombre i/o tipus de nodes a reservar per a l'ús exclusiu del treball.

El PBS proporciona una interfície per a l'interpret d'ordres i una interfície gràfica. Tot i així el més habitual en la computació d'altres prestacions és utilitzar la interfície per mitjà de l'interpret d'ordres per a enviar, monitoritzar, modificar i eliminar treballs. Hi ha diferents tipus d'ordres en funció del tipus d'usuari que les fa servir:

- Ordres d'usuari: `qsub`, `qstat`, `qdel`, `qselect`, `qrerun`, `qorder`, `qmove`, `qhold`, `qalter`, `qmsg`, `qrsls`.
- Ordres d'operació: `qenable`, `qdisable`, `qrun`, `qstart`, `qstop`, `qterm`.
- Ordres d'administració: `qmgr`, `pbsnodes`.

A continuació veurem la sintaxi d'algunes de les ordres més utilitzades des del punt de vista de l'usuari en PBS:

```
qsub [-a date_time] [-A account_string] [-c interval]
     [-C directive_prefix] [-e path] [-h] [-I] [-j join]
     [-k keep] [-l resource_list] [-m mail_options]
     [-M user_list] [-N name] [-o path] [-p priority]
     [-q destination] [-r c] [-S path_list] [-u user_list]
     [-v variable_list] [-V] [-W additional_attributes]
     [-z] [script]
```

L'ordre `qsub` envia un nou treball al servidor del PBS per a la seva execució. Per defecte es considera que el treball es posarà en la cua i s'executarà en segon pla, però també es permet executar el treball (igualment quan planificar ho considerem oportú) de manera interactiva, si s'especifica l'opció `-q`. Típicament, el `script` que s'envia per mitjà de l'ordre és un *shell script* de l'interpret d'ordres com ara `sh` o `csh`.

Les opcions de l'ordre `qsub` que s'indiquen anteriorment permeten especificar atributs que afectaran el comportament del treball. A més, l'ordre `qsub` passarà al treball una llista de variables d'entorn que estaran disponibles durant l'execució del treball. El valor de les variables `HOME`, `LANG`, `LOGNAME`, `PATH`, `MAIL`, `SHELL` i `TZ` es pren de les variables d'entorn de l'ordre `qsub`. Aquests valors s'assignaran a variables que comencen per "PBS_O_". Per exemple, el treball tindrà accés a una variable d'entorn anomenada `PBS_O_HOME`, que tindrà el valor de la variable `HOME` en l'entorn de `qsub`. A part de les variables d'entorn anteriors, també hi ha les següents disponibles per al treball:

- `PBS_O_HOST`. És el nom del servidor en què l'ordre `qsub` s'està executant.
- `PBS_O_QUEUE`. És el nom de la cua inicial a la qual es va enviar el treball.
- `PBS_O_WORKDIR`. És la ruta (absoluta) del directori de treball que s'ha indicat en l'ordre `qsub`.

- `PBS_ENVIRONMENT`. Indica si el treball és en segon pla (`PBS_BATCH`) o interactiu (`PBS_INTERACTIVE`).
- `PBS_JOBID`. És l'identificador de treball que li assigna el sistema de cues.
- `PBS_JOBNAME`. És el nom del treball proporcionat per l'usuari.
- `PBS_NODEFILE`. És el nom del fitxer que conté la llista de nodes que s'assignaran al treball.
- `PBS_QUEUE`. És el nom de la cua des de la qual s'ha executat el treball.

```
qstat [-f][-W site_specific]
      [job_identifier... | destination...]

qstat [-a|-i|-r] [-n] [-s] [-G|-M] [-R] [-u user_list] ñññ[job_identifier... | destination...]

qstat -Q [-f][-W site_specific] [destination...]

qstat -q [-G|-M] [destination...]

qstat -B [-f][-W site_specific] [server_name...]
```

L'ordre `qstat` indica l'estat d'un treball, de les cues i del servidor de cues. L'estat es proporciona mitjançant la sortida estàndard. Quan se sol·licita l'estat d'un treball mitjançant una de les dues opcions que s'han mostrat, `qstat` retorna la informació de cadascun dels treballs indicats mitjançant el seu identificador o tots els treball d'una destinació donada. Cal tenir present que els treballs en què l'usuari no té privilegis no es mostraran. Quan s'utilitza una de les tres últimes opcions que s'han mostrat, `qstat` retornarà la informació de cada destinació.

```
qdel [-W delay] job_identifier ...
```

L'ordre `qdel` elimina els treballs especificats en l'ordre en què es proporcionen. Un treball només el pot eliminar o bé el seu propietari o bé l'operador o l'administrador del sistema de cues. Per a eliminar un treball, el servidor enviarà un senyal `SIGTERM` seguit per un senyal `SIGKILL`.

```
qalter [-a date_time] [-A account_string] [-c interval]
      [-e path] [-h hold_list] [-j join] [-k keep]
      [-l resource_list] [-m mail_options] [-M user_list]
      [-N name] [-o path] [-p priority] [-r c] [-S path]
      [-u user_list] [-W additional_attributes]
      job_identifier...
```

L'ordre `qalter` modifica els atributs d'un o més treballs a partir del seu identificador. Només es poden modificar les opcions que es mostren en la descripció de l'ordre.

Els treballs s'envien mitjançant *scripts* que permeten comunicar-se amb el programa d'enviament de treballs MPI per a especificar-li la quantitat i quins nodes s'han d'utilitzar dels disponibles per al PBS, segons els requeriments de l'usuari. El *server* del PBS no executarà més treballs als nodes ocupats fins que s'hagi acabat el treball actual. A continuació es mostra un exemple de *script* de PBS per a l'execució d'una aplicació MPI mitjançant `mpiexec` (equivalent a `mpirun`, que hem vist anteriorment).

```
#!/bin/sh
#! Exemple de fitxer de definició de treball per enviar mitjançant qsub
#! Les línies que comencen per #PBS són options de l'ordre qsub

#! Nombre de processos (8 en aquest cas, 4 per node)
#PBS -l nodes=4:ppn=2

#! Nombre dels fitxers per a la sortida estàndard i error
#! Si no s'especifiquen, per defecte són <job-name>.o<job_number> i <job-name>.e<job_number>
#PBS -e test.err
#PBS -o test.log

#! Adreça electrònica de l'usuari per a quan el treball acabi o s'avorti
#PBS -m ae

#! Directori de treball
echo Working directory is $PBS_O_WORKDIR
#!cd <working directory>
echo Running on host 'hostname'
echo Time is 'date'
echo Directory is 'pwd'
echo This jobs runs on the following processors:
echo 'cat $PBS_NODEFILE'

MPI executable - it's possible to redirect stdin/stdout of all processes
#! using "<" and ">" - including the double quotes
/usr/local/bin/mpiexec -bg a.out
```

També hi ha altres gestors de cues com ara Loadleveler o SLURM. El Loadleveler és un gestor de cues dissenyat per IBM que es caracteritza per la seva facilitat per a processar treballs en entorns clúster. Permet executar treballs paral·lels (per exemple, MPI) i és fàcilment escalable a milers de processadors. El Loadleveler va ser un dels primers sistemes a incorporar l'algoritme de planificació *backfilling*. SLURM és un gestor de cues de codi obert dissenyat per a clústers Linux de diferents mides.

5.2. Planificació

El problema de la planificació de tasques o treballs en un sistema de computació d'altres prestacions es pot dir que consisteix a, a partir d'una llista de treballs que s'esperen en la cua, les seves característiques –com, per exemple, el nombre de processadors que necessiten, els fitxers executables, els fitxers d'entrada, etc.– i l'estat del sistema, decidir quins treballs s'han d'executar i en quins processadors. Aquest problema es pot dividir en dues etapes que segueixen cadascuna una política concreta: la política de planificació de treballs i la política de selecció de recursos.

Com mostra la figura 46, en la primera etapa de la planificació, el planificador (*scheduler*) ha de decidir quins dels treballs que s'esperen en la cua (*job1*, *job2*, ..., *jobN*) han de començar, tenint en compte el recursos que hi ha disponibles (CPU1, CPU2, CPU3, CPU4). L'algoritme que fa servir el planificador per a fer aquesta selecció s'anomena *política de planificació de treballs*¹⁸. Un cop s'ha seleccionat el treball més apropiat per a ser executat, el planificador demana al gestor de recursos (*resource manager*) que assigni el treball seleccionat (*jobi*) als processadors més adequats, en funció dels requeriments del treball (en la figura 46, el treball *jobi* requereix dos processadors). Llavors, el gestor de recursos seleccionarà els processadors més apropiats (CPU1 i CPU2 en l'exemple) i assignarà processos del treball a cadascun. El gestor de recursos implementa una política de selecció de recursos¹⁹ concreta.

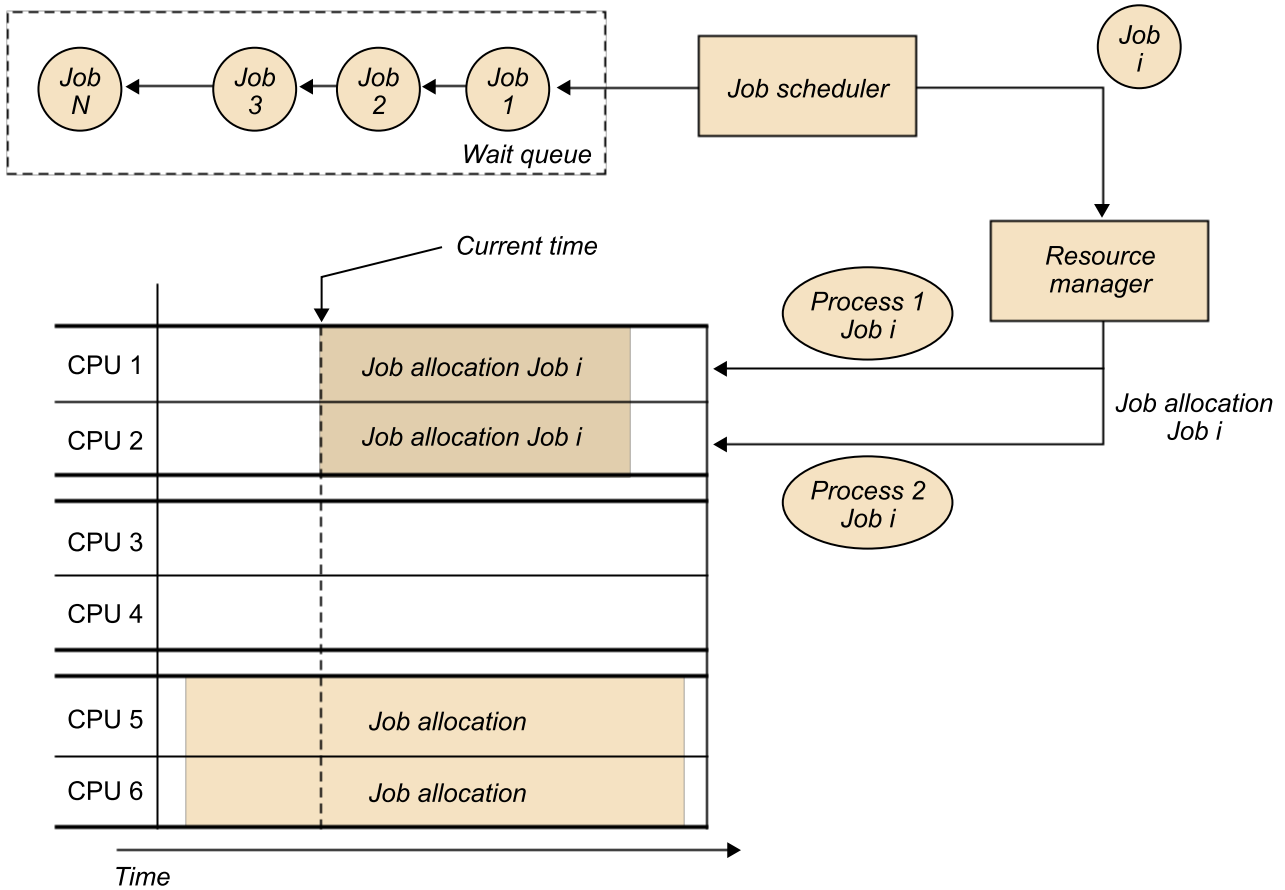
Treball

Utilitzarem el terme *treball* per a referir-nos a l'execució d'una aplicació amb uns paràmetres concrets durant tot el subapartat.

⁽¹⁸⁾En anglès, *job scheduling policy*.

⁽¹⁹⁾En anglès, *resource selection policy*.

Figura 46. El problema de la planificació mitjançant un exemple



S’han proposat moltes tècniques i polítiques de planificació i de selecció de recursos durant les últimes dècades per a optimitzar la utilització dels recursos i reduir els temps de resposta dels treballs. Tal com veurem a continuació, les polítiques de planificació de treballs que més s’han utilitzat en entorns d’altes prestacions són les que estan basades en *backfilling*. S’ha demostrat que aquestes polítiques proporcionen un bon balanç entre el rendiment del sistema (nombre de processadors utilitzats) i el temps de resposta dels treballs. Tot i així, el problema principal de l’algoritme que utilitzen aquestes polítiques és que l’usuari ha de proporcionar una estimació del temps d’execució dels treballs que envia i, normalment, aquestes estimacions no són gens acurades. A més, en moltes situacions l’usuari no té suficients coneixements o informació per a proporcionar aquests paràmetres.

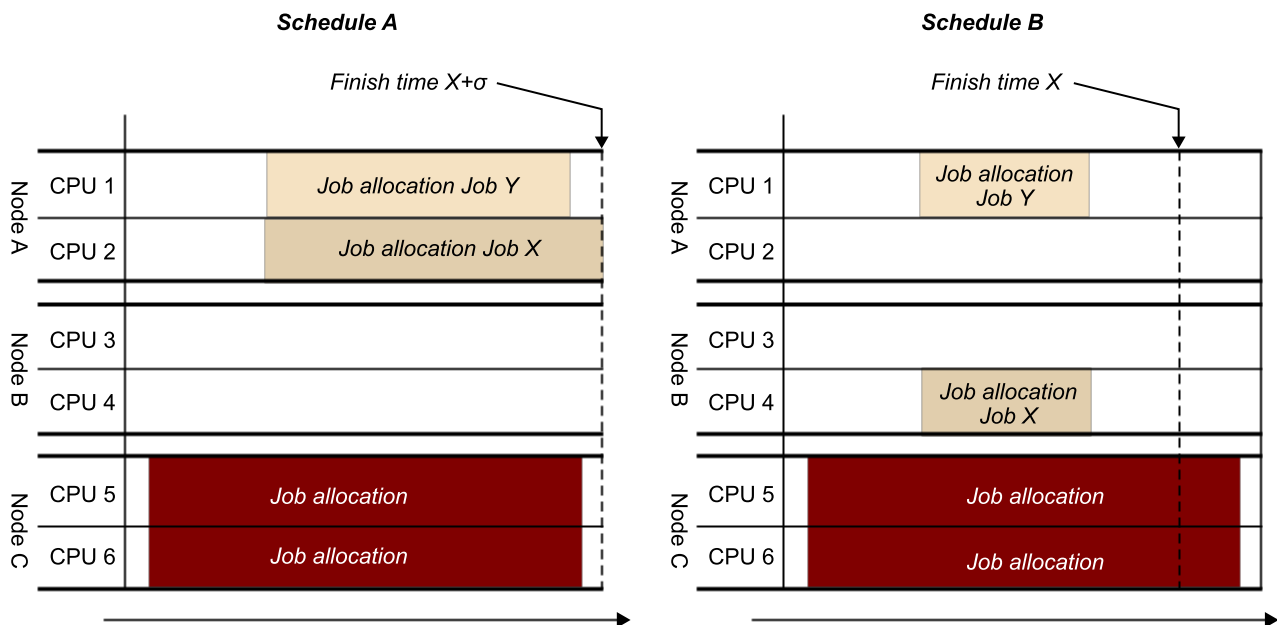
Temps de resposta del treball

És el temps des que s’envia un treball fins que n’acaba l’execució, incloent-hi el temps a la cua.

Per exemple, un usuari no expert en arquitectura de computadors no ha de saber necessàriament que una aplicació pot tenir una certa duració en un tipus d’arquitectura o un 25% més en un altre tipus concret que inicialment pot semblar similar (per exemple, per la quantitat de memòria cau). Tot i així, el temps d’execució d’un treball pot dependre de molts factors com ara el nombre de processadors emprats, l’estat del sistema, les dades d’entrada del programa, etc. També pot succeir que usuaris amb experiència no disposin de suficient informació per a estimar el temps d’execució. Les polítiques de planificació de treballs se centren en la selecció dels treballs que s’han d’executar, però moltes vegades no tenen en compte la ubicació que han de tenir. Aquesta decisió moltes vegades es prenen de manera unilateral mitjançant polítiques de selecció de recursos.

Normalment el que fa el sistema de recursos és ubicar els treballs de dues maneres diferents, o bé un conjunt de filtres predefinits s'apliquen a les característiques dels treballs per a trobar els processadors adequats, o bé els processos dels treballs s'ubiquen en nodes complerts quan l'usuari especifica que els treballs han d'executar-se en nodes no compartits. La primera manera no té en compte l'estat del sistema i, en conseqüència, un treball intensiu de memòria es podria ubicar en un node en què l'amplada de banda a memòria està pràcticament saturada, la qual cosa causaria una degradació important de tots els treballs ubicats en el node. La segona manera normalment implica que la utilització dels recursos és substancialment reduïda perquè normalment els treballs no utilitzen tots els processadors i recursos del node en què estan ubicats. Per tant, en certes arquitectures, quan els treballs comparteixen recursos (per exemple, la memòria) poden resultar amb sobrecàrrega. Això té un impacte negatiu en el rendiment dels treballs que utilitzen els recursos i hi ha un impacte negatiu colateral en el rendiment del sistema. Aquest problema s'il·lustra en la figura 47, en la qual es mostren diferents planificacions dels treballs quan els treballs *job Y* i *job X* s'envien al sistema. Tenint en compte les polítiques de planificació actuals, la planificació més probable seria el cas A. En aquesta, els treballs *job Y* i *job X* s'ubiquen consecutivament en el mateix node (node A). Si els dos treballs són intensius en memòria, hi hauria una reducció en el rendiment de la seva execució deguda a la contenció de recursos. En canvi, si el planificador té en compte l'estat del sistema i els requeriments dels treballs, la planificació més probable seria la del cas B. En aquest darrer cas, tant el treball X com el Y no tindrien penalització en el rendiment de la seva execució per contenció de recursos.

Figura 47. El problema de la utilització de recursos

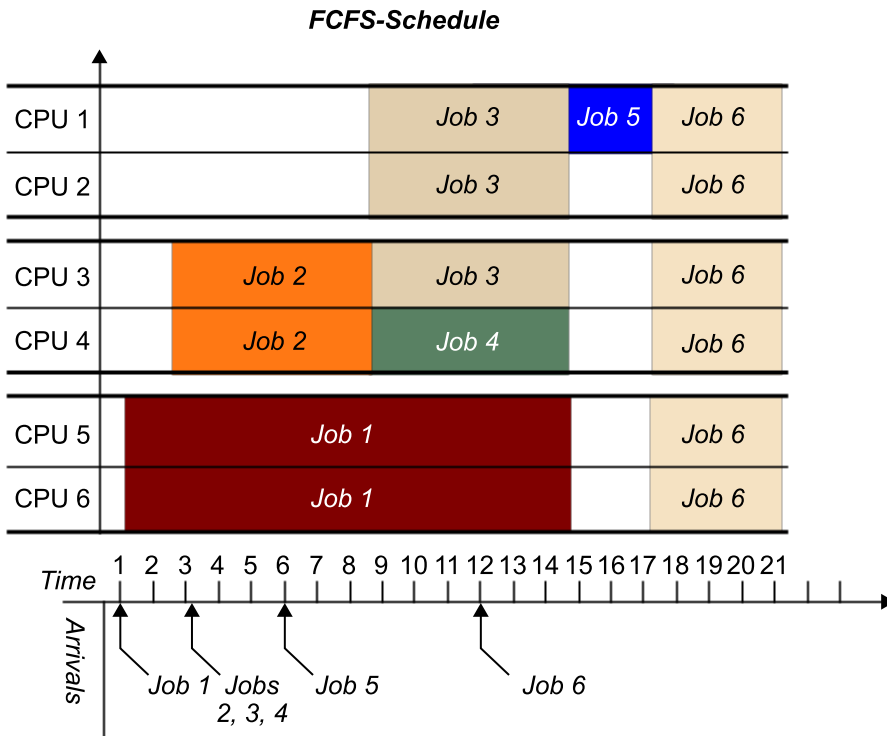


5.2.1. Polítiques de planificació basades en *backfilling*

Des del principi dels anys noranta s'ha investigat molt activament en l'àmbit de la planificació de treballs en computadors d'altres prestacions, i les polítiques basades en *backfilling* són les que s'han utilitzat més extensament en centres d'altres prestacions. El *backfilling* és una optimització de la política de planificació FCFS⁽²⁰⁾. En FCFS, el planificador posa en cua tots els treballs enviats en ordre d'arribada i cadascun d'aquests treballs té associat un nombre de processadors requerits. Quan un treball finalitza la seva execució, si hi ha suficients recursos per a començar l'execució del treball següent de la cua, el planificador pren el treball de la cua i el comença a executar. En cas contrari, el planificador ha d'esperar fins que hi hagi prou recursos disponibles i no es pot executar cap altre treball de la cua. La figura 48 mostra una possible planificació de la política FCFS. El problema principal d'utilitzar aquesta política és que el sistema es fragmenta i la utilització dels recursos pot ser molt baixa.

(20) De l'anglès *first-come-first-serve*.

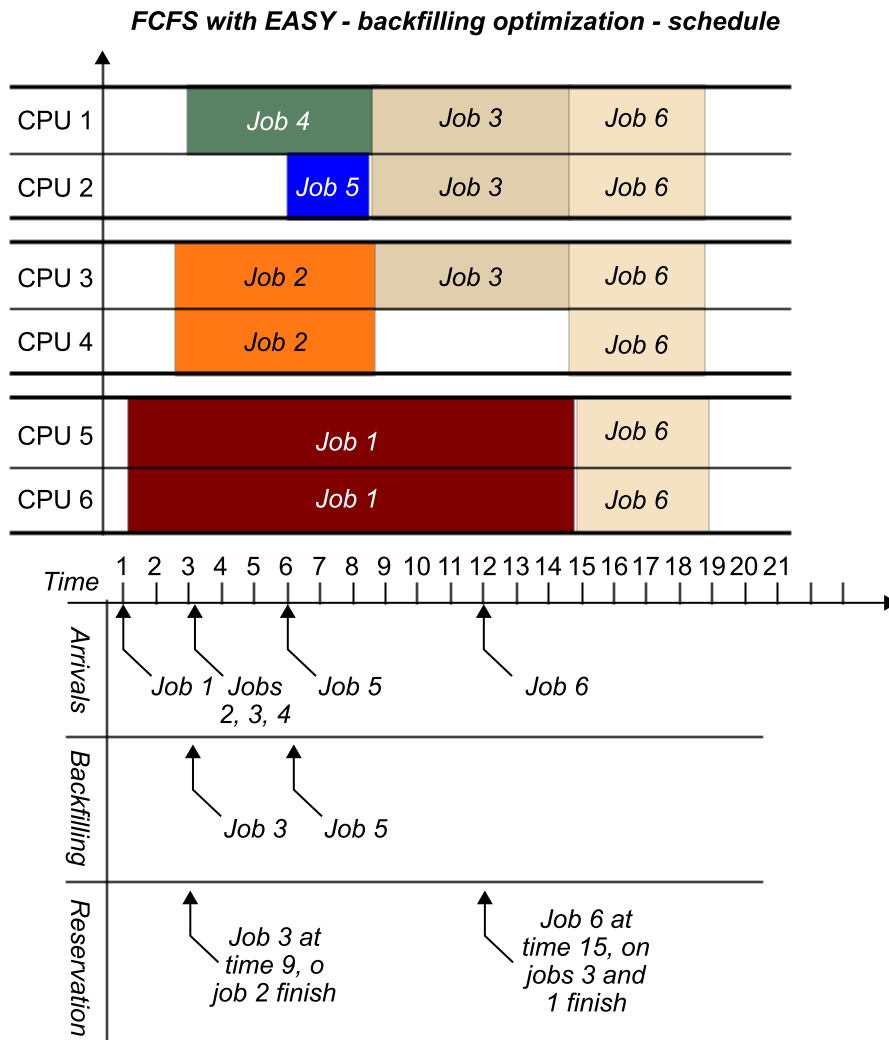
Figura 48. Exemple de planificació basada en la política FCFS



Amb la política de planificació *backfilling* es permet executar treballs que s'han posat a la cua posteriorment a d'altres que ja són en la cua mentre que el temps estimat d'execució dels treballs que ja hi són no s'endarrereixi. Noteu que per a aplicar aquesta optimització, cal saber el temps d'execució aproximat dels treballs quan s'envien (per tant, abans de l'execució), ja que sinó és així no hi ha prou informació per a assegurar que el temps d'execució del treball del principi de la cua no s'endarrereixi. La figura 49 presenta una possible planificació de l'exemple anterior amb optimització mitjançant una política *backfilling*. En aquest exemple, als treballs *job 4* i *job 5* se'ls ha aplicat *backfilling* perquè els treballs *job 3* i *job 6* no poden començar, ja que no hi ha prou processadors

disponibles per a ells. Gràcies a aquesta optimització, la utilització del sistema es pot millorar molt significativament. Molts dels processadors que quedarien aturats sense feina amb una política FCFS poden avançar a partir de *backfilling*.

Figura 49. Exemple de planificació FCFS optimitzada amb EASY *backfilling*

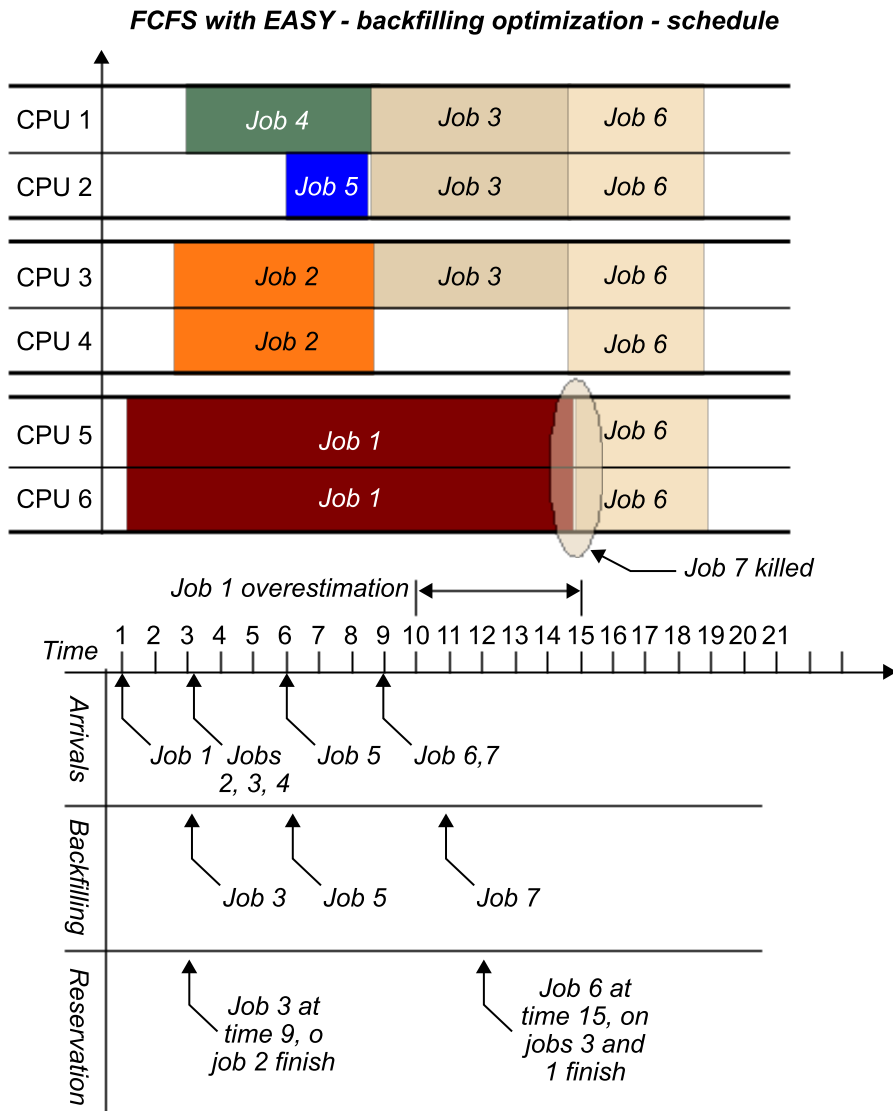


L'inconvenient principal d'aquest algoritme és que l'usuari ha de proporcionar el temps d'execució estimat per al treball quan aquest s'envia al sistema de cues. En situacions en què el temps d'execució que l'usuari ha proporcionat és inferior al temps real d'execució, el treball és matat pel planificador quan detecta que ha excedit el temps d'execució sol·licitat. L'altra situació que es dona freqüentment en aquest escenari és que el temps d'execució estimat és substancialment superior al real.

En l'exemple de la figura 50 es mostren les dues situacions que s'han descrit anteriorment. En aquesta planificació el treball *job 1* ha estat sobreestimat. Un cop el planificador ha detectat aquesta situació i utilitzant l'estimació de temps d'execució proporcionat per l'usuari, el planificador aplica *backfilling* al

treball *job 7*. Tot i això, el treball *job 7* és matat més endavant, ja que excedeix el temps d'execució especificat. Noteu que en cas que el treball *job 7* no fos matat, el temps d'inici del treball *job 6* podria endarrerir-se.

Figura 50. Exemple de sobreestimació i subestimació de recursos en planificació FCFS optimitzada amb *EASY backfilling*



La política *backfilling* que hem presentat anteriorment és la política basada en *backfilling* més senzilla, que va proposar per Lifka i altres, que s'anomena *EASY-backfilling*. S'han proposat moltes variants d'aquesta política, de les quals podem destacar les que tenen en compte consideracions com, per exemple, l'ordre en què es poden aplicar *backfillings* als treballs. Un exemple és la variant Shortest-Job-Backfilled-First.

Activitat

Busqueu informació sobre diferents polítiques i tècniques de planificació de treballs present com a referència el treball del doctor Feitelson.

Lectura complementària

J. Skovira; W. Chan; H. Zhou; D. A. Lifka (1996). "The easy-loadleveler api project. proceedings of the workshop on job scheduling strategies for parallel processing". *Lecture Notes In Computer Science* (vol. 1162, pàg. 41-47).

Bibliografia

- Gaster, B.; Howes, L.; Kaeli, D. R.; Mistry, P.; Schaa, D.** (2011). *Heterogeneous computing with OpenCL*. Burlington, Massachusetts: Morgan Kaufmann.
- Grama, A.; Karypis, G.; Kumar, V.; Gupta, A.** (2003). *Introduction to parallel computing*. Reading, Massachusetts: Addison-Wesley.
- Hwang, K.; Fox, G. C.; Dongarra, J. J.** (2012). *Distributed and cloud computing: from parallel processing to the Internet of things*. Burlington, Massachusetts: Morgan Kaufmann.
- Jain, R.** (1991). *The Art of computer system performance analysis: techniques for experimental design, measurement, simulation, and modeling*. Wiley-Interscience.
- Kirk, D. B.; Hwu, W. W.** (2010). *programming massively parallel processors: a hands-on approach*. Burlington, Massachusetts: Morgan Kaufmann.
- Lin, C.; Snyder, L.** (2008). *Principles of parallel programming*. Reading, Massachusetts: Addison Wesley.
- Munshi, A.** (2009). *The OpenCL specification*. Khronos OpenCL Working Group.
- Munshi, A.; Gaster, B.; Mattson, T. G.; Fung, J.; Ginsburg, D.** (2011). *OpenCL programming guide*. Reading, Massachusetts: Addison-Wesley Professional.
- Nvidia** (2007). *Nvidia CUDA compute unified device architecture*. Technical Report.
- Nvidia** (2007). *The CUDA compiler driver NVCC*.
- Pacheco, P.** (2011). *An introduction to parallel programming*. Burlington, Massachusetts: Morgan Kaufmann.
- Sanders, J.; Kandrot, E.** (2010). *CUDA by example: an introduction to general-purpose GPU programming*. Reading, Massachusetts: Addison-Wesley Professional.
- Skovira, J.; Chan, W.; Zhou, H.; Lifka, D. A.** (1996). "The easy-loadleveler api project. Proceedings of the workshop on job scheduling strategies for parallel processing". *Lecture Notes In Computer Science* (vol. 1162, pàg. 41-47).

