



MÉMOIRE DE STAGE MASTER 2

Stage du 1^{er} février 2023 au 31 juillet 2023

Intitulé : Interruptions en espace utilisateur pour le réseau
BXI

Charles GOEDEFROIT

Encadré par Alexandre DENIS (Inria), Grégoire PICHON (Atos) et
Mathieu BARBE (Atos)

17 juillet 2023

Table des matières

1	Introduction	2
1.1	Presentation Inria	2
1.2	Presentation Atos (Eviden)	2
1.3	Environnement de travail	2
1.4	Le cadre	2
1.5	Presentation du plan	3
2	Context	3
2.1	HPC	3
2.2	OS bypass	3
2.3	Polling	4
2.3.1	Inconvénients du polling	5
2.4	BXI	5
2.5	MPI	5
2.6	NewMadeleine	5
2.7	Travaux antérieur	6
3	Problématiques / Objectifs	6
3.1	Sujet	6
3.1.1	Nouveau mécanisme d'interruption en espace utilisateur	6
3.1.2	contentions...	6
3.2	Ne plus faire de Polling	6
3.3	Objectifs	6
3.3.1	Objectifs global	6
3.3.2	Objectif du stage	6
4	Exploration de uintr	7
4.1	Prérequis est accès	7
4.2	Fonctionnement des uintr	7
4.2.1	Détail	7
4.2.2	Capacité présente et futur	7
4.2.3	Exemple	8
4.2.4	Partage du FD	8
4.2.5	+	8
4.3	Tests du mécanisme	8
4.4	Correction du patch pour l'appel système uintr_alt_stack . .	8
4.5	Mesure de la latence (Mathieu)	8
4.6	Performances	8

5	Intégration dans NewMadeleine	9
5.1	Présenté NewMadeleine details	9
5.2	Modification	9
5.3	Suite de tests	10
5.4	Performances	10
5.4.1	Résultats avec attente active	10
5.4.2	Résultats recouvrement communication par du calcule	10
6	Remerciements	10
7	Annexes	11

1 Introduction

J'ai effectué mon stage de 6 mois dans l'équipe-projet TADaaM d'Inria. Le sujet du stage à été proposé par Alexandre Alexandre Denis (Inria) et Grégoire Pichon (Atos). Ils ont aussi encadré le stage aussi que Mathieu BARBE (Atos).

1.1 Présentation Inria

Inria est TODO : here

1.2 Présentation Atos (Eviden)

en restructuration je parle donc d'Atos dans ce document

1.3 Environnement de travail

- labo de recherche
- un bureau dans l'open space de l'équipe projet TADaaM.
- Salle de visio pour la réunion de suivie de stage toute la semaine
- Baby foot ?
- Cafeteria ?
- ... ?

1.4 Le cadre

Inria (TADaaM) :

- objectifs du service : faire de la recherche...
- type de personnel : chercheur, post-doc, ingénieur, doctorant et stagiaire
- Culture informatique : HPC, système... ?
- Matériel et logiciels mis à disposition :
 - un ordinateur portable
 - un doc avec écran, clavier et souris
 - service inria (email, mattermost, webex...)
- contraintes : je suis libre pour la démarche, les outils...

Atos (BXI-LL) :

- objectifs du service : développé le soft ba niveau pour la NIC BXL... (entre MPI et la carte) donc driver, lustre, portal...
- type de personnel : ingénieurs... ?
- Culture informatique : HPC, système... ?
- Matériel et logiciels mis à disposition :
 - une machine avec 2 CPU SPR accessible via VPN
- contraintes : je suis libre

1.5 Présentation du plan

Contexte dans lequel le stage se place puis les objectifs à long terme et les objectifs du stage. Le déroulement du stage avec les résultats. Pour finir un bilan et la suite. TODO :

2 Contexte

2.1 HPC

Le stage s'est déroulé dans le contexte du Calcul Haute Performance (HPC, High-Performance Computing en anglais). Le HPC utilise des supercalculateurs pour la simulation numérique et le pré-apprentissage d'intelligences artificielles. Ces simulations simulent la dynamique des fluides, la résistance structurelle, les interactions moléculaires, les flux d'air...

Elles couvrent différents domaines :

- L'industrie : les médicales, l'automobile, l'aviation, les constructions, l'aérospatiale...
- La défense : simulation atomique
- La recherche scientifique : la création de galaxies, la fusion nucléaire, le climat...
- La météo

Aujourd'hui, les supercalculateurs sont composés de plusieurs ordinateurs que l'on appelle nœuds de calcul. Ces nœuds sont regroupés en grappe (cluster), ils sont tous vus et utilisés comme une seule grande machine. Ce fonctionnement pose des questions sur la répartition du calcul, la distribution de la mémoire, la communication entre les différents nœuds... Dans le contexte du stage, nous nous sommes concentrés sur les communications entre nœuds. Les nœuds sont connectés entre eux par un réseau haute performance. Ce réseau est dédié aux communications et n'est pas forcément de type Ethernet. La gestion des nœuds est généralement effectuée par un second réseau plus traditionnel (Ethernet ; TCP/IP). Ce genre de réseau a une latence autour de quelques microsecondes. Chaque nœud possède une ou plusieurs cartes réseau et il faut les programmer.

2.2 OS bypass

Le stage se passe donc aussi dans un contexte système. Chaque nœud a son propre système d'exploitation qui permet la gestion des ressources, des processus, des fichiers, des périphériques. Pour cela, le système a 2 espaces :

- un espace noyau où seul le code du système peut s'exécuter. Le code du système peut donc être modifié n'importe où dans la mémoire, exécuté n'importe où...

- un espace utilisateur où le code de l'utilisateur est exécuté. Cette espace est limité par le noyau qui controle ce que fait l'utilisateur.

En temps normale les périphériques sont programmé directement depuis le noyau du système d'exploitation ceci pour des questions de sécurité, de standardisation des accès... Lors que l'on passe par le noyau (kernel en anglais) on à un surcoût qui n'est pas négligeable dans notre cas. Pour passer par le noyau on utilise généralement des appels système qui se présente sous la forme d'une fonction. Un appel système va effectuer un changement de contexte (context switch) pour passer de l'espace utilisateur à l'espace noyau. Ce changement de contexte est coûteux car il sauvegarde les états du code de l'utilisateur avant d'exécuter celui du noyau. Une fois le context switch effectué le code noyau de l'appel système s'exécute avant de refaire un context switch pour cette fois passer de l'espace noyau à l'espace utilisateur, et donc restorer l'états du code de l'utilisateur. Donc le fait de passer par un appel système coûte plusieurs microseconde. On voit donc qu'on ne peut pas utiliser d'appel système car un appel système est déjà du même ordre de temps qu'une communication. En HPC on programme donc directement la carte réseau à partir de l'espace utilisateur (OS bypass en anglais). Pour cela on initialise toujours la carte à partir du noyau mais on fait une projection de la mémoire et des registre de la carte réseau dans l'espace d'adressage virtuelle du processus utilisateur.

Pour transmettre des événements à l'utilisateur les périphériques utilisent généralement les interruptions, qui passe par le système donc on ne les utilise pas en HPC. En HPC on fait donc du polling.

2.3 Polling

Le polling consiste à scruter (poll) régulièrement si on a reçu un événement. Concrètement cela consiste à lire une zone mémoire modifiée par la carte réseau et voir si un bit est passé à 1. Pour cela il est possible de dédié un thread qui va faire de l'attente active, scruter sans cesse si un événement a été reçu. Mais on perd une unité de calcul lorsque ce que le thread est ordonner donc de la puissance de calcul, donc cette technique est peu utilisée. Il est aussi possible d'entrelacer le code de l'utilisateur avec des scrutations, c'est fréquemment utilisé mais ça oblige à l'utilisateur de prendre en compte la progression. Une autre solution qui est utilisée par Pioman dans NewMadelaine (j'en parlerai plus tard) consiste à faire ces scrutations de manière opportuniste dans les threads qui ont fini leur calcul, mais pour cela il faut déjà utiliser plusieurs threads et avoir une application qui a des calculs d'une durée hétérogène.

Il faut donc, peu importe la technique, régulièrement scruter pour faire progresser les communications.

2.3.1 Inconvénients du polling

Dans le cas d'un thread dédié qui fait de l'attente active la réactivité est excellente hormis quand il y a plus de thread que d'unité de calcul est que le thread n'est pas ordonnancer.

Quand on entrelace le calcul et les scrutations la réactivité est moins bonne car il faut attendre que le calcul soit fini pour faire un poll. C'est ce qui est fait habituellement.

Quand on utilise les threads de façon opportuniste pour faire des scrutations la réactivité est moins bonne car il faut qu'il y ait un thread disponible.

Donc un choix doit être fait entre perdre de la capacité de calcul ou perdre en réactivité. En plus on perd un peu de temps de calcul à effectuer du polling.

2.4 BXI

BXI est un type de réseau dédié pour les communications entre nœuds...

- les cartes réseaux BXI (Bull eXascale Interconnect) sont développées par Atos
- elles sont capables de faire des communications sans intervention du CPU
- le CPU a juste à soumettre une commande et la carte s'occupe du reste
- pour savoir si la communication est finie l'application peut poll ou recevoir une interruption en provenance de la carte
- portal4
-

2.5 MPI

- présenté MPI
- les communications asynchrones
- la progression se fait :
 - au moment des appels à la biblio (MPI_Isend, MPI_Irecv, MPI_test, MPI_wait).
 - grâce à un thread dédié
 - grâce à une politique d'utilisation des threads de façon opportuniste (quand ils ne font pas de calcul).
- ...
-

2.6 NewMadeleine

- NewMadeleine est une bibliothèque de communication...
- elle est basée sur un système de progression asynchrone...

- supporte différent type de communication grâce a un système de driver (portal4 pour BXI, ibverb pour IB, shm...)
- ...

2.7 Travaux antérieur

- Travaux de Mathieu Barbe durant un stage en 2019 sur l'utilisation d'interruption pour transmettre des événements réseau.
- Ces travaux vis à diminué la latence du au fait de passer par un driver noyaux.
- Il parle dans les perspective de directement traitement des interruptions depuis l'espace utilisateur ce qui mène à mon stage.

3 Problématiques / Objectifs

3.1 Sujet

3.1.1 Nouveau mécanisme d'interruption en espace utilisateur

- très récent...
- sur les dernier CPU Inter Sapphire Rapide...
-

3.1.2 contentions...

3.2 Ne plus faire de Polling

- utiliser des interruption en HPC comme dans les autre domain
- quand un périphériques a besoin de remonté une information il n'y a plus besoin de polling

3.3 Objectifs

3.3.1 Objectifs global

- L'objectifs à terme est de faire progressé les communications entre 2 noeuds sans utiliser du Polling
- Donc permettre au carte réseau BXI de d'utiliser les interruption en espace utilisateur.
- Réduire globalement le temps de calcul. (Mathieu)
-

3.3.2 Objectif du stage

- Comprendre le fonctionnement de uintr (à partir du manuel du CPU, des manuel du patch et des quelque document disponible)

- Connaître leurs propriétés.
- Mesurer leur performance.
- intégré le mécanisme dans le driver de mémoire partagé (shm) de NewMadeleine, pour des communication entre processus (ipc)
- faire progresser les communication à partir d'un handler d'interruption
- évité les problèmes du polling :
 - réactivité
 - mobilisation d'une ressource ou truffé le code de MPI_Test (simplifie la vie à l'utilisateur)
 - ...

4 Exploration de uintr

4.1 Prérequis est accès

- difficulté d'accès au processeur...
- J'ai eu accès au processeur 2 mois 1/2 après le début du stage
- J'ai donc eu accès à une machine avec 2 CPU SPR... VPN...
- il faut une version patché du noyaux linux pour utilisé c'est nouvelle interruption Ce patch n'est pas disponible en mainline linux. J'ai donc télécharger le patch fourni par intel, je l'ai compiler puis je l'ai installer sur la machine
- Pour utiliser les nouvelle instructions il faut une version récente de GCC... (version). Pas disponible avec LLVM-Clang

4.2 Fonctionnement des uintr

TODO : comparaison avec les interruption de base

4.2.1 Détail

- 5 instructions
- registre MSR pour l'init

4.2.2 Capacité présente et futur

- process => process
- kernel => process
- device => process
- au niveaux des thread...
- masquage
- le thread dois être en ring-3
- peut fonctionné si le thread est endormi mais j'ai pas testé car ...

4.2.3 Example

- exemple avec schéma

4.2.4 Partage du FD

- Pour mes testes "jouer" j'ai utiliser `uintr_register_self()`
- pipe / socket / URL pour NewMadeleine (on en reparle en section XXX)
- `pidfd_getfd`

4.2.5 +

- comparé au signaux ?
- TODO ...

4.3 Tests du mécanisme

- entre process
- entre threads
- avec alt stack
- test d'écrasement des interruptions
- avec binding
- avec turbo boost

4.4 Correction du patch pour l'appel système `uintr__alt__stack`

TODO

4.5 Mesure de la latence (Mathieu)

TODO : expliquer que le binding est important

4.6 Performances

TODO

- Dans un premier temps j'ai pas bind les thread par pu donc j'avais de mauvais résultats.
- J'ai fait différents bind (on vois que c'est similaire sauf quand on passe entre 2 NUMA)
- si on monte la frequence c'est mieux (turbo boost)
-

5 Intégration dans NewMadeleine

5.1 Présenté NewMadeleine details

- liste de progression recv
- liste de progression send
- p_pw
- post
- poll
- driver
- nm_schedule (appelé par nm_wait, ...)
- pioman
- existe : progression du nm_schedule ou de pioman vers le core_task

5.2 Modification

TODO : on désactive le poll avec les driver qui on des handler. Mais on fait toujours un poll qui fonctionne du premier coup.

- ajout d'un driver sig_shm
- ajout d'un driver uintr_shm
- progression à partir du handler du driver vers les core_task
- problématiques de gestion des interruptions
 - on ne peut pas faire d'attente dans les handler Les fonction doivent être async safe
 - on ne peut pas utiliser d'allocateur donc il faut utiliser des p_pw déjà alloué
 - pour avoir un p_pw disponible il faut faire progresser le communications
 - la progresser à une partie critique qui naissaisite un verrou
 - si on ne peut pas récupéré le core_lock (try_lock) il faut mettre à plus tard le trétement de l'interruption
 - on utilise donc une file lock-free
 - problématique des liste lock-free qui doivent être wait-free
 - le principe d'une lfqueue est d'attente si quelqu'un d'autre modifie la file
 - il faut donc une file wait-free
 - j'ai fait de la biblio
 - décrire les différante solutions
 - solution d'Alexandre
- gestion des multiple interruptions qui s'écrase, (prob_any / pour les large si la progression à traiter le pipeline courant)
- quand est ce qu'on envois des interruption ?
 - au moment ou on poste le premier paquet d'envois, pour indiquer au recepteur que des données sont disponible (le recepteur à toujours un paquet de reception posté)

- au moment ou une progression est fait du coté du recepateur, pour indiquer une reception a l'émetteur
- l'émetteur reçoit une interruption est détermine si l'émission est fini, petit paquet, ou si il faut envoyer la suite, paquet large.

5.3 Suite de tests

Tous les tests ne passe pas ?

5.4 Performances

5.4.1 Résultats avec attente active

TODO : voire le sur coup des interruption

5.4.2 Résultats recouvrement communication par du calcule

TODO : courbe overlap reception

6 Remerciements

7 Annexes