



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)

20 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	4
2.3	Polling	5
2.3.1	Inconvénients du polling	5
2.4	BXI	5
2.5	MPI	6
2.5.1	Les communications point à point	6
2.5.2	Les communications collective	6
2.5.3	La progression asynchrone	7
2.6	NewMadeleine	7
2.7	Travaux antérieur	8
3	Problématiques / Objectifs	8
3.1	Sujet	8
3.1.1	Nouveau mécanisme d'interruption en espace utilisateur	8
3.1.2	contentions...	8
3.2	Ne plus faire de Polling	8
3.3	Objectifs	9
3.3.1	Objectifs global	9
3.3.2	Objectifs du stage	9
4	Exploration de uintr	10
4.1	Prérequis est accès	10
4.2	Fonctionnement des uintr	10
4.2.1	Détail	10
4.2.2	Capacité présente et futur	10
4.2.3	Exemple	10
4.2.4	Partage du FD	10
4.2.5	+	11
4.3	Tests du mécanisme	11
4.4	Correction du patch pour l'appel système uintr_alt_stack . .	11
4.5	Mesure de la latence (Feedback Mathieu)	11
4.6	Performances	11

5	Intégration dans NewMadeleine	11
5.1	Présenté NewMadeleine details	11
5.2	Modification	12
5.3	Suite de tests	12
5.4	Performances	12
5.4.1	Résultats avec attente active	12
5.4.2	Résultats recouvrement communication par du calcule	13
6	Bilan	13
7	Remerciements	13
8	Annexes	14

1 Introduction

J'ai effectué mon stage de 6 mois dans l'équipe-projet TADaaM du Centre Inria de l'université de Bordeaux. 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 l'institut national français de recherche en sciences et technologies du numérique. Compte plus de 3 900 chercheurs et ingénieurs au sein de 215 équipes-projets. La plupart des centres sont communes avec une grande université.

1.2 Présentation Atos (Eviden)

Atos est un des leaders internationaux de la transformation numérique. Atos couvre un champ large d'activité : le cloud, la cybersécurité, les services, le conseil, les supercalculateurs... Atos compte 112 000 collaborateurs. Atos est en restructuration pour se séparer en 2 entités. L'entité qui nous intéresse pour ce stage est Eviden qui englobe les activités autour des supercalculateurs et du HPC. Cette restructuration est récente donc quand je parle d'Atos dans ce document je parle de la partie Eviden.

1.3 Environnement de travail

Mon environnement de travail c'est les locaux du Centre Inria de l'université de Bordeaux. J'ai un bureau dans l'open space de l'équipe-projet TADaaM. Je peux participer et assister à des activités scientifiques intéressantes (séminaire, soutenance de thèse, soutenance HDR, activité diverse...). J'ai accès à la salle de visio conférence, notamment pour les réunions de suivi de stage hebdomadaire. Il y a aussi un Baby-foot, une cafétéria, une petite Médiathèque...

1.4 Le cadre

Du côté de l'équipe TADaaM de Inria l'objectif de l'équipe est de faire de la recherche. Elle est composée de chercheurs, ingénieurs de recherche, post-doc, doctorant et stagiaire. La culture informatique de l'équipe est l'utilisation des environnements linux, des logiciels open source, des clusters de calcul HPC, le traitement des données et le système. L'équipe a mis à ma disposition un ordinateur portable avec une station de travail branchée à un écran, à internet, à un clavier et une souris. Inria donne aussi accès à un ensemble de services, boîte email, service de communication Mattermost, service de visio Webex, intranet avec le menu de la cafétéria... J'ai été libre

d'installer mon système d'exploitation, d'utilisé les outils informatique que je voulais...

Du coté d'Atos l'équipe BXI-LL (BXI Low Level) est une équipe don l'objectif est de fournir le support logiciels bas niveau pour les carte réseaux BXI. Cela consiste à développée les drivers de la carte, le support Lustre¹ ainsi que l'interface logiciels Portal4. Cette interface permet aux implémentations MPI d'utiliser le réseaux BXI. Le personnel est composée d'ingénieurs, de doctorant et de stagiaire. /* TODO : vérifié */ La culture informatique de l'équipe est l'utilisation des cluster HPC, la programmation système et l'utilisation d'un système de sécurité PKI pour accéder aux ordinateurs et aux services interne. Atos nous a donné accès à une machine avec deux processeur Intel® Sapphire Rapids. Pour accéder à cette machine ils nous on fourni un accès ssh et un compte VPN. J'ai pu avoir tous les droits d'accès sur cette machine pour changé le noyaux du système.

1.5 Présentation du plan

Dans ce mémoire je vais commencé par vous présenté le context dans le quelle le stage ce place. Je continuerai par les objectifs à long term et les objectifs du stage. Je vous présenterai le nouveau mécanisme d'interruption et ce que j'ai fait avec. Je présenterai l'intégration de ce mécanisme dans la bibliothèque de communication NewMadeleine et les tests que l'on à effectué. Pour finir je ferai un bilan du stage et je parlerai de la suite.

2 Context

2.1 HPC

Le stage c'est déroulé dans le context du Calcule Haute Performance (HPC, High-Performance Computing en anglais). Le HPC utilise des supercalculateur pour la simulation numérique et le pré-apprentissage d'intelligences artificielles. Ces simulation simule la dynamique des fluides, la résistance structurelle, les interaction moléculaire, les flux d'aire...

Elles couvrent différant domaines :

- L'industrie : le médicales, l'automobiles, l'aviations, la constructions, l'aérospatiales...
- La défense : simulation atomique
- La recherche scientifiques : la création des galaxie, la fusion nucléaire, le climat...
- La météo

1. Lustre est un système de fichier parallèle et distribué. Il est utiliser pour l'I/O est ne fait pas partie de ce stage.

De nos jours les supercalculateurs sont composés de plusieurs ordinateurs que l'on appelle nœuds de calculs. Ceci est regroupé en grappe (cluster), il est tout vu et utilisé 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). Les réseaux haute performance ont 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 modifier n'importe quelle adresse de la mémoire, exécuter n'importe quelle instruction...
- un espace utilisateur où le code de l'utilisateur est exécuté. Cette espace est limitée par le noyau qui contrôle ce que fait l'utilisateur.

En temps normal les périphériques sont programmés 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 a un surcoût qui n'est pas négligeable dans notre cas. Pour passer par le noyau on utilise généralement des appels système qui se présentent 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 restaurer l'état du code de l'utilisateur. Donc le fait de passer par un appel système coûte plusieurs microsecondes. 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 registres 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 passent 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 dédier 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 ordonnancer 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 l'utilisateur à 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, peut 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 threads que d'unités 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 pour Bull eXascale Interconnect est un type de réseau d'interconnexion (réseau haute performance) développé par Atos. Historiquement développé par Bull qui a été racheté par Atos. Ce réseau est dédié aux communications entre nœuds. Il est composé de cartes réseau BXI et de switch BXI. La carte BXI est capable de faire progresser la communication réseau sans aucune intervention du processeur (offload des communications réseau). Le processeur a juste à soumettre une commande dans une file sur la carte et elle s'occupe de tout. Le processeur peut ensuite récupérer une file d'événements pour savoir ce qu'il s'est passé, en somme faire un poll. La carte est

également capable de déclencher des interruptions. L'utilisation de la carte passe donc par une implémentation du protocole Portal4. L'utilisateur utilise donc le protocole pour envoyer et recevoir des paquets réseau.

2.5 MPI

MPI pour Message Passing Interface est un standard pour fournir une interface pour effectuer des communications entre plusieurs processus, qui sont souvent sur des nœuds différents et qu'on appelle *processus MPI*. Cette interface fournit une abstraction pour transmettre des données entre plusieurs processus. L'abstraction masque la complexité des communications. L'interface permet d'envoyer et de recevoir des messages. Pour cela il y a deux modes de communications :

2.5.1 Les communications point à point

C'est à dire entre deux processus MPI, aussi appelé One-to-One. Pour ce faire le processus MPI récepteur va appeler la fonction `MPI_Recv` qui est bloquante et va attendre la réception d'un message. L'émetteur va lui faire un appel à la fonction `MPI_Send` qui est aussi bloquante et va envoyer un message et attendre que la communication soit finie. L'utilisation des fonctions `MPI_Send` / `MPI_Recv` est donc totalement synchrone ce qui bloque le code. La norme propose aussi une version non bloquante de ces fonctions qui sont `MPI_Isend` et `MPI_Irecv`. Cette version se contente de poster la communication et rend immédiatement la main. Pour la progression et vérifier la terminaison il faut donc d'autres fonctions `MPI_Test` qui vérifie la progression, et la fait si nécessaire, et la fonction `MPI_Wait` qui attend activement la terminaison et s'occupe de la progression si nécessaire. Il est important de noter que le standard ne précise pas si la progression se fait en tâche de fond ou non, c'est aux implémentations de la norme MPI de choisir. C'est pour cela que la progression peut se faire au niveau de ces fonctions `MPI_Wait` et `MPI_Test` ou être faite avant est donc l'appel aux fonctions s'occupe juste de la terminaison. L'envoi des messages peut donc être asynchrone.

2.5.2 Les communications collective

Les communications collectives se font entre plusieurs processus. Il en existe de différents types :

- un processus vers plusieurs (One-to-All) par exemple un broadcast d'un message
- de plusieurs processus vers un seul (All-to-One) par exemple une réduction (e.g. un processus reçoit la somme des valeurs des autres processus)

- de plusieurs processus vers plusieurs (All-to-All) par exemple quand tous les processus ont un message pour les autres

Pour les collectives il existe également deux versions, bloquante et non bloquante, qui fonctionnent de la même façon que les communications point à point.

2.5.3 La progression asynchrone

Pour faire progresser les communications de façon asynchrone il est possible de :

- faire des appels à `MPI_test` régulièrement et faire du calcul entre chaque appel. Cela nous permet de recouvrir la latence des communications par du calcul. La progression se fait également dans d'autres appels aux fonctions MPI. Quand il n'y a plus de calcul à faire on repasse à une progression synchrone par un appel à `MPI_Wait` sauf si la communication est déjà finie.
- utiliser un thread dédié aux progressions. Dans ce cas c'est la bibliothèque MPI qui s'occupe de la progression des communications en tâche de fond grâce à un thread dédié. Il faut donc faire attention au placement des threads et à prendre en compte qu'un thread est déjà utilisé par la bibliothèque de communication. Il faut aussi éviter d'appeler trop souvent `MPI_Test` car ça crée de la contention
- utilisation des threads de façon opportuniste c'est à dire quand un des threads à fini son calcul il va faire progresser les communications. C'est ce qui est fait par Pioman dans NewMadeleine.

2.6 NewMadeleine

NewMadeleine est une bibliothèque de communications qui historiquement a été développée pour le RPC (Remote procedure call). Elle a toujours cette capacité mais elle a évolué et implémente une interface MPI. On peut donc la considérer comme une implémentation du standard MPI avec des fonctionnalités supplémentaires. Elle est développée au Centre Inria de l'université de Bordeaux. Elle est basée sur un système de progression événementielle qui lui permet d'être asynchrone. Elle est composée en modules ce qui lui permet de charger dynamiquement des stratégies d'optimisation sur les paquets. Les stratégies sont l'agrégation de paquets, l'utilisation de priorités, l'utilisation d'arbre de décision ?, `split_balance` ?. Elle possède également un système de drivers pour supporter différents types de communications comme des réseaux (e.g. `portals4` pour BXI, `ibverbs` pour InfiniBand, `psm2` pour OmniPath, `ofi` (Open Fabrics Interfaces), TCP) et localement (`shm` (mémoire partagée), `socket`, `self`).

2.7 Travaux antérieur

TODO :

- système de progression dans NewMadeleine
- progression avec pioman
- overlap
- 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

Le sujet est *Interruptions en espace utilisateur pour le réseau BXI*.

3.1.1 Nouveau mécanisme d'interruption en espace utilisateur

Ce nouveau mécanisme qui permet de dérouler une interruption à partir de l'espace utilisateur et très récent. Il est seulement disponible sur les processeur Intel® Sapphire Rapids qui sont officiellement sortis le 10 janvier 2023. La plupart des processeur ont été disponibles à la vente le 14 mars. AMD n'a pas encore annoncé de de support pour les interruption en espace utilisateur.

3.1.2 contentions...

Le fonctionnement actuelle de la progression des communications amène à des problèmes de contention mémoire car plusieurs threads peuvent chercher à lire / modifier la même zone mémoire. Utiliser les interruptions permettrai de ne plus avoir ce problèmes car seul le thread concerné par une zone mémoire serai prévenu.

TODO :

3.2 Ne plus faire de Polling

Dans la plupart des autre domain les périphériques envoi une interruption à l'application pour avertir d'un changement. Souvent on passe par les signaux ou un appel système bloquante. L'idée du stage est de ne plus avoir à faire du polling et utiliser les interruptions comme dans les autre domains.

3.3 Objectifs

Il y a plusieurs objects. Les principaux sont :

- La réduction du temps de calcul.
- utiliser des interruptions en espace utilisateur pour replacé le polling.
- Simplifier pour l'utilisateur le recouvrement des communications par du calcul pour qu'il n'est plus à ajouter des `MPI_Test` en pleins milieu de ces calcul.
- Amélioré la réactivité des communications sans qu'il y est besoin d'une unité de calcul dédié à l'attente active.

3.3.1 Objectifs global

Le stage s'inscrit dans un objectif globale qui est de faire progresser les communications entre plusieurs noeuds du réseau BXI sans faire de polling et en utilisant les interruptions en espace utilisateur. Cette objectif permet de réduire globalement le temps de calcul d'une application. Pour ce faire la carte réseau BXI devra être capable de levé des interruptions en espace utilisateur. Le fait de supprimer le polling et de réduire le temps de calcul permettra de diminué la consommation électrique.

3.3.2 Objectifs du stage

Le premier objectif est de défricher le fonctionnement des interruption en espace utilisateur à partir du manuel Intel® de l'architectures 64 et IA-32 pour les développeurs logiciels, de la présentation^{2 3} du mécanisme de Sohil Mehta⁴ et du patch du noyau linux avec ces manuels. Le second objectif est de connaître les propriétés du mécanisme et de mesuré ça performance. Le troisième objectif est d'envisager l'intégration de ces interruption dans le driver de mémoire partagé (shm) de NewMadeleine. L'objectif est de testé dans un premier temps le fonctionnement avec des communications entre processus (IPC⁵ en anglais). Pour intégrer les interruption dans les drivers NewMadeleine il faut également permettre la progression des communications à partir d'un handler d'interruption. Le dernier objectif est de montré que l'utilisation d'interruption permet bien d'amélioré le recouvrement des communications par du calcul.

2. Qui est une diapositive associer à des discussion sur LWN.net

3. <https://lwn.net/Articles/869140/>

4. Ingénieur Intel® qui à développé le patch noyau

5. Inter Process Communication

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... détails Saphir Rapid et 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 (Feedback Mathieu)

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 (plus pourquoi avec binding que comparaison avec et sans)
- avec turbo boost

4.4 Correction du patch pour l'appel système uintr _alt _stack

TODO

4.5 Mesure de la latence (Feedback 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 voit que c'est similaire sauf quand on passe entre 2 NUMA)
- si on monte la fréquence 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 recepteur, pour indiquer une reception a l'émetteur
 - l'émetteur reçoit une iterruption 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 Bilan

7 Remerciements

8 Annexes