



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)

4 août 2023

Table des matières

1	Introduction	3
1.1	Présentation Inria	3
1.2	Présentation Atos (Eviden)	3
1.3	Environnement de travail	3
1.4	Le cadre	4
1.5	Présentation du plan	4
2	Contexte	5
2.1	HPC	5
2.2	OS bypass	5
2.3	Les interruptions matérielles	6
2.4	Polling	7
2.4.1	Inconvénients du polling	7
2.5	BXI	8
2.6	MPI	8
2.6.1	Les communications point à point	8
2.6.2	Les communications collective	9
2.6.3	La progression en tâche de fond	9
2.7	NewMadeleine	10
2.8	Les signaux	10
2.9	Travaux antérieur	10
2.10	Définitions	11
3	Problématiques / Objectifs	11
3.1	Sujet	11
3.1.1	Nouveau mécanisme d'interruption en espace utilisateur	11
3.2	Projet global	11
3.2.1	Les objectifs	12
3.3	Objectifs du stage	12
3.4	La suite	13
4	Exploration des interruptions en espace utilisateur.	13
4.1	Prérequis et accès	13
4.2	Fonctionnement des interruptions en espace utilisateur	14
4.2.1	Les interruptions	14
4.2.2	Les uintr	16
4.2.3	Capacités présentes et futures	18
4.2.4	Exemple de fonctionnement	19
4.2.5	Partage du descripteur de fichier	22
4.3	Tests du mécanisme	22
4.4	Correction d'un bug dans le patch noyau	24
4.5	Mesure de la latence	24

4.6	Performances	26
5	Intégration dans NewMadeleine	31
5.1	Présenté NewMadeleine details	31
5.2	Modification	32
5.3	Suite de testes	33
5.4	Performances	33
5.4.1	Résultats avec attente active	33
5.4.2	Résultats recouvrement communication par du calcule	33
6	Bilan	33
7	Remerciements	33
8	Annexes	34
9	Bibliographie	40

1 Introduction

J'ai effectué mon stage de fin de master de six mois dans l'équipe-projet *TADaaM* du Centre Inria de l'université de Bordeaux. Le sujet du stage a été proposé par Alexandre Denis (Inria) et Grégoire Pichon (Atos). Ils ont aussi encadré le stage avec Mathieu BARBE (Atos).

1.1 Présentation Inria

Inria est l'institut national français de recherche en sciences et technologies du numérique. Il compte plus de 3 900 chercheurs et ingénieurs au sein de 215 équipes-projets. La plupart des centres sont commun avec de grandes universités.

1.2 Présentation Atos (Eviden)

Atos est l'un des leaders internationaux de la transformation numérique. Elle couvre un large éventail d'activités, notamment le cloud, la cybersécurité, les services transactionnels, le conseil, l'infogérance, le Big Data, les supercalculateurs, etc. *Atos* compte 112 000 collaborateurs. Actuellement, *Atos* est en restructuration afin de se séparer en deux entités. L'entité qui nous intéresse pour ce stage est Eviden, qui englobe notamment les activités liées aux supercalculateurs et au HPC. Cette restructuration est récente, donc lorsque nous évoquons *Atos* dans ce document, nous parlons de la partie Eviden.

Atos est le seul grand constructeur européen de supercalculateurs. Elle est donc le leader européen et est bien positionnée mondialement. Elle est présente dans le secteur du HPC depuis l'acquisition de la société *Bull* en 2014. *Bull* était déjà impliquée dans le HPC depuis les années 2000 et en avait fait son cœur de métier.

La machine *Leonardo* en Italie, construite par *Atos*, se classe quatrième au dernier Top500. Ce classement évalue les machines les plus puissantes pour le HPC et est publié deux fois par an.

1.3 Environnement de travail

Mon environnement de travail se trouve dans les locaux du Centre Inria de l'université de Bordeaux. Un bureau dans l'open space de l'équipe-projet *TADaaM* a été mis à ma disposition. Je peux participer et assister à des activités scientifiques intéressantes, telles que des séminaires, des soutenances de thèse, une soutenance HDR et d'autres activités diverses. Nous avons également accès aux salles de visioconférence, notamment pour les réunions de suivi de stage hebdomadaires. En plus, il y a un Baby-foot, une cafétéria, une petite médiathèque, etc.

1.4 Le cadre

Du côté de l'équipe *TADaaM* d'Inria, l'objectif de l'équipe est de mener des recherches sur les sujets suivants :

- Gestion des I/O (ordonnancement, bande passante...)
- Placement de processus
- Partitionnement de maillage (i.e. SCOTCH)
- Localité matérielle (i.e. hwloc)
- Optimisation des communications pour les réseaux haute performance, MPI (i.e. NewMadeleine)

Elle est composée de chercheurs, d'ingénieurs de recherche, de post-doc, de doctorants et de stagiaires. Sa culture informatique est l'utilisation d'environnements Linux, de logiciels open source, de clusters de calcul HPC, le traitement des données et le systèmes. L'équipe a mis à ma disposition un ordinateur portable avec une station de travail reliée à un écran, à Internet, à un clavier et une souris. Inria donne également l'accès à un ensemble de services, notamment une boîte e-mail, un service de communication (Mattermost), un service de visio (Webex) et un intranet avec le menu de la cafétéria, entre autres. Nous sommes libres d'installer le système d'exploitation et d'utiliser les outils informatiques de notre choix.

Du côté d'*Atos* l'équipe BXI-LL (BXI Low Level) a pour objectif de fournir le support logiciel bas niveau pour les cartes réseaux *BXI*. Cela comprend le développement des pilotes de la carte, le support Lustre (un système de fichiers parallèle et distribué utilisé pour l'I/O, qui ne fait pas partie de ce stage), ainsi qu'une implémentation de l'interface logicielle *Portal4*. *Portal4* est une interface logicielle qui permet une abstraction des opérations possibles sur les réseaux HPC. Son but est de faciliter l'utilisation de ces opérations réseau avec les meilleures performances possibles pour les implémentations MPI. /* je suis pas sur de ma présentation de Portal4 */ L'implémentation d'*Atos* permet aux implémentations MPI d'utiliser le réseau *BXI*. Le personnel de l'équipe est composé d'ingénieurs, de doctorants et de stagiaires. La culture informatique de l'équipe est l'utilisation de clusters HPC, la programmation système et l'utilisation d'un système de sécurité PKI pour accéder aux ordinateurs et aux services internes. *Atos* nous a donné accès à une machine avec deux CPU Intel® Sapphire Rapids. Pour accéder à cette machine, ils nous ont fourni un accès SSH et un compte VPN. Nous avons eu tous les droits d'accès sur cette machine pour changer le noyau du système.

1.5 Présentation du plan

Dans ce mémoire, nous commencerons par une présentation du contexte dans lequel le stage se place. Ensuite, nous aborderons les objectifs à long terme ainsi que les objectifs du stage. Nous continuerons en décrivant le

nouveau mécanisme d'interruption en espace utilisateur et ce que nous en avons fait. Ensuite, nous verrons l'intégration de ce mécanisme dans la bibliothèque de communication NewMadeleine, ainsi que les tests que nous avons effectués. Pour finir, nous ferons un bilan du stage et évoquerons les travaux qui suivront.

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 : médical, automobile, aviation, construction, aérospatiale...
- La défense : simulation atomique...
- La recherche scientifique : création de galaxies, fusion nucléaire, climat...
- La météo

De nos jours, les supercalculateurs sont composés de plusieurs ordinateurs que l'on appelle noeuds de calcul. Ceux-ci sont regroupés en grappes (clusters) et sont tous vus et utilisés comme une seule grande machine. Ce fonctionnement soulève des questions concernant la répartition du calcul, la distribution de la mémoire et la communication entre les différents noeuds. Dans le contexte du stage, nous nous sommes concentrés sur les communications entre noeuds. Les noeuds sont connectés entre eux par un réseau haute performance dédié aux communications, qui n'est pas forcément de type Ethernet. La gestion des noeuds est généralement effectuée par un second réseau plus traditionnel (Ethernet ; TCP/IP). Les réseaux haute performance ont une latence de l'ordre de quelques microsecondes ce qui n'est pas le cas des réseaux classiques qui ont une plus grande latence. Avec ce type de réseaux, le débit est aussi bien plus élevé. Chaque noeud 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 noeud 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 (kernel en anglais) où seul le code du système peut s'exécuter. Le code du système peut donc modifier n'importe quel endroit de la mémoire, exécuter n'importe quelle instruction...
- Un espace utilisateur où le code de l'utilisateur est exécuté. Cet espace est limité 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, cela est fait pour des raisons de sécurité et de standardisation des accès. Cependant, lorsque l'on passe par le noyau, il y a un surcoût qui n'est pas négligeable dans notre cas. En effet, pour passer par le noyau, on utilise généralement des appels système qui prennent la forme d'une fonction. Un appel système effectue 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. save états du process (registres, registre vectoriel, le pc, esp, flags...), Une fois le changement de contexte effectué, le code noyau de l'appel système s'exécute avant de refaire un changement de contexte pour, cette fois, passer de l'espace noyau à l'espace utilisateur et ainsi restaurer l'état du code de l'utilisateur. Le noyau peut aussi être amené à mettre à jour ses structures internes, exécuter des tâches qu'il avait mises à faire plus tard, etc. Donc, le fait de passer par un appel système coûte plusieurs microsecondes. On voit donc qu'il n'est pas préférable d'utiliser les appels système, car leur durée est du même ordre de grandeur qu'une communication entre noeuds. 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 virtuel du processus utilisateur.

Pour transmettre des événements à l'utilisateur, les périphériques utilisent généralement les interruptions ordinaires, qui passent par le système donc on ne les utilise très peu en HPC. En HPC on privilégie donc le polling que nous verrons en se 2.4.

2.3 Les interruptions matérielles

Les interruptions matérielles, également appelées IRQ pour "Interrupt ReQuest", existent depuis longtemps et servent à remonter des événements extérieurs, par exemple d'un clavier, d'un port de communication, d'un port IDE, etc. Elles ont ensuite été utilisées pour signaler des événements en provenance de périphériques plus variés, comme des cartes réseau. Dans ce document, nous utiliserons le terme "interruption ordinaire" ou "IRQ" pour les désigner.

2.4 Polling

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

Le *polling* consiste à scruter (poll) régulièrement si un événement a été reçu. Concrètement, cela consiste à lire une zone mémoire modifiée par la carte réseau et vérifier si un bit est passé à 1. Pour cela, il est possible de dédier un thread qui effectuera une attente active, scrutant sans cesse si un événement a été reçu. Cependant, cette technique entraîne la perte d'une unité de calcul quand le thread est ordonnancé, ce qui réduit la puissance de calcul, donc elle est utilisée dans certains cas comme expliqué dans cet article [4]. Une autre approche consiste à entrelacer le code de l'utilisateur avec les scrutations, ce qui est fréquemment utilisé, mais cela oblige l'utilisateur à prendre en compte la progression. Une troisième solution, utilisée par *Pioman* dans *NewMadeleine* (nous en parlerons plus tard), consiste à effectuer 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 avec des calculs de durée hétérogène.

2.4.1 Inconvénients du polling

Dans le cas d'un thread dédié qui effectue une attente active, la réactivité est excellente, sauf lorsque le nombre de threads dépasse le nombre d'unités de calcul et lorsque le thread n'est pas ordonnancé.

Dans le cas où l'on entrelace le calcul et les scrutations, la réactivité est moins bonne car il faut attendre que le calcul soit terminé pour effectuer un poll. C'est ce qui est habituellement fait.

Dans le cas où l'on utilise les threads de façon opportuniste pour effectuer des scrutations, la réactivité est moins bonne car il faut qu'il y ait un thread disponible.

Un choix doit donc être fait entre perdre de la capacité de calcul ou perdre en réactivité. De plus, un peu de temps de calcul est perdu à effectuer du polling.

Dans le cas où l'on effectue trop de scrutations, on peut rencontrer des problèmes de contention mémoire. En effet, le CPU manipule la mémoire par lignes de 8 mots mémoire appelées lignes de cache. Ces lignes de cache font 512 bits sur un CPU 64 bits (car $8 * 64 = 512$), et 256 bits pour des CPU 32 bits (car $8 * 32 = 256$). Pour garantir la cohérence d'accès à la mémoire, ces lignes de cache sont invalidées lorsque quelqu'un d'autre les manipule. Ainsi, si un thread scrute trop souvent une zone mémoire ou que plusieurs threads cherchent à scruter la même zone mémoire, on aboutit à des problématiques de contention de la mémoire, car la ligne de cache est invalidée plus souvent, ce qui entraîne un surcoût. Avec le polling nous sommes donc contraints à respecter une granularité liée à la taille d'une ligne

de cache. Utiliser les interruptions permettrait de ne plus avoir ce problème, car seul le thread concerné par une zone mémoire serait prévenu donc plus de contention à ce niveau.

En plus, utiliser des interruptions permettrait de ne plus faire des scrutations pour rien, mais seulement scruter lorsqu'il y a un événement à traiter. Le nombre de scrutations serait donc grandement réduit.

2.5 BXI

BXI [6] 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 noeuds. Il est composé de cartes réseau *BXI* et de switches *BXI*. Les cartes *BXI* sont capables de faire progresser les communications réseau sans aucune intervention du CPU (offload des communications réseau). Le CPU a juste à soumettre une commande dans une file sur la carte, et elle s'occupe de tout. Le CPU peut ensuite récupérer une file d'événements pour savoir ce qu'il s'est passé, en somme faire un poll. Les cartes sont également capables de déclencher des interruptions. Les cartes implémentent donc l'API de communication *Portal4* [13]. Les utilisateurs utilisent donc cette API pour envoyer et recevoir des paquets réseau.

2.6 MPI

MPI, pour "Message Passing Interface", est un standard permettant de fournir une interface pour effectuer des communications entre plusieurs processus, souvent situés sur des noeuds différents, que l'on appelle *processus MPI*. Cette interface fournit une abstraction pour transmettre des données entre ces différents processus, masquant ainsi la complexité des communications. L'interface permet d'envoyer et de recevoir des messages. Pour cela, il existe deux modes de communications :

2.6.1 Les communications point à point

C'est-à-dire entre deux processus MPI, également 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 également bloquante et va envoyer un message, puis attendre que la communication soit terminée. L'utilisation des fonctions `MPI_Send` / `MPI_Recv` bloque le code, ce qui nous fait perdre du temps à attendre. Le standard MPI propose également 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 utiliser d'autres fonctions comme `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 MPI 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 des fonctions `MPI_Wait` et `MPI_Test`, ou être faite avant, et donc l'appel aux fonctions s'occupe juste de la terminaison. L'envoi des messages peut donc être asynchrone. /* gardé cette phrase ? */

2.6.2 Les communications collective

Les communications collectives se font entre plusieurs processus. Il en existe 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 lorsque tous les processus ont un message pour les autres.

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

2.6.3 La progression en tâche de fond

Pour faire progresser les communications en tâche de fond, il est possible de :

- Faire régulièrement des appels à `MPI_test` et effectuer des calculs entre chaque appel. Cela permet de recouvrir la latence des communications par du calcul. La progression peut également se faire dans d'autres appels aux fonctions MPI. Lorsqu'il n'y a plus de calcul à effectuer, on repasse à une progression bloquante par un appel à `MPI_Wait`, sauf si la communication est déjà terminée.
- Utiliser un thread dédié aux progressions. Dans ce cas, c'est la bibliothèque MPI ou l'application 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 prendre en compte qu'un thread est déjà utilisé pour les communications. Il faut aussi éviter d'appeler trop souvent `MPI_Test` car cela peut créer de la contention.
- Utilisation des threads de façon opportuniste, c'est-à-dire qu'un des threads, une fois son calcul terminé, fera progresser les communications. C'est ce qui est fait par *Pioman* dans *NewMadeleine*.

2.7 NewMadeleine

NewMadeleine est une bibliothèque de communications qui prend en charge le RPC pour "Remote Procedure Call" et implémente également 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, ce qui lui permet d'être asynchrone. Elle est composée de modules, ce qui lui permet de charger dynamiquement des stratégies d'optimisation sur les paquets. Ces stratégies sont l'agrégation de paquets, l'utilisation de priorités et le multi-rail. /* TODO : précision multi-rail/split_balance */. Elle possède également un système de drivers pour supporter différents types de communications, tels que des réseaux (e.g. portals4 pour BXI, ibverbs pour InfiniBand, psm2 pour OmniPath, ofi pour "Open Fabrics Interfaces", TCP) et des communications locales (shm pour "mémoire partagée", socket, self).

2.8 Les signaux

Les signaux[12] sont l'un des mécanismes permettant l'interface entre un processus et le système. Pour cela, le processus utilisateur effectue des appels système pour définir les signaux pour lesquels il souhaite être notifié. Pour être notifié, l'utilisateur peut définir un handler pour chaque signal auprès du système. Ces handler sont déclenchés par le noyau qui prend en charge toutes les opérations nécessaires (sauvegarde de l'état du processus, changement de contexte, etc.).

Les signaux que l'utilisateur peut recevoir correspondent à des exceptions système (e.g. *SIGFPE* pour une division par zéro ou *SIGSEGV* pour une erreur de segmentation), à des informations (e.g. *SIGTERM* pour demander au processus de se fermer) ou à des notifications d'autres processus (i.e. *SIGUSR1* et *SIGUSR2*). Il est possible de masquer la réception de certains signaux en utilisant des appels système de masquage.

2.9 Travaux antérieur

TODO : Amélioré? le faire pour tous ce qui existe

Les origines de *NewMadeleine* sont décrites dans cet article [1]. Le système de progression dans *NewMadeleine* fait l'objet de plusieurs travaux. le fonctionnement du système de progression opportuniste *Pioman* est décrit dans cet article[3]. Pour l'impact de l'overlap et la manière de le mesurer, l'article [5] explique tout en détail. Concernant l'utilisation des interruptions pour transmettre des événements réseau, les travaux de Mathieu Barbe pendant son stage en 2019 sont disponibles ici [2]. Ces travaux visent à réduire la latence due au passage par un driver noyau. Ils abordent également les perspectives de traitement direct des interruptions depuis l'espace utilisateur, ce

qui a conduit à ce stage.

2.10 Définitions

Il est important d’avoir les définitions suivantes en tête :

- *CPU* désigne la puce dans sa totalité.
- *core* ou *processeur* désigne l’un des coeur du *CPU*.
- *core logique*, *processeur logique* ou *unité de calcul* désigne une unité de calcul au sein d’un coeur. Il y en a deux par core dans les CPUs Intel® Sapphire Rapids fourni par *Atos*.

/* TODO : besoin d’expliquer ? : Un processus peut avoir plusieurs threads, un thread correspond à un fil d’exécution de code. Et tous les threads au sien d’un même processus partage le même espace d’adressage virtuel. */

3 Problématiques / Objectifs

3.1 Sujet

Le sujet est *Interruptions en espace utilisateur pour le réseau BXI* [11].

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 est très récent. Il est seulement disponible sur les CPU Intel® Sapphire Rapids qui sont officiellement sortis le 10 janvier 2023.¹ La plupart des CPU ont été disponibles à la vente le 14 mars.² AMD n’a pas encore annoncé de support pour les interruptions en espace utilisateur.

3.2 Projet global

Dans la plupart des autres domaines, les périphériques envoient une interruption ordinaire à l’application, par le biais des signaux ou d’un appel système bloquant, pour avertir d’un changement. L’idée serait de faire la même chose en HPC grâce aux interruptions en espace utilisateur. Le projet global vise donc à faire progresser les communications entre plusieurs nœuds du réseau *BXI* sans faire de polling et en utilisant les interruptions en espace utilisateur. Cela permettrait de réduire globalement le temps de calcul d’une application. Pour ce faire, la carte réseau *BXI* devra être capable de lever des interruptions en espace utilisateur. Le fait de supprimer le polling

1. <https://www.datacenterknowledge.com/intel/intel-launches-sapphire-rapids-after-4-delays-it-worth>

2. <https://www.intel.com/content/www/us/en/newsroom/news/4th-gen-intel-xeon-sprints-into-market.html>

et de réduire le temps de calcul permettra de diminuer la consommation électrique.

3.2.1 Les objectifs

Il y a plusieurs objectifs, les principaux sont les suivants :

- La réduction du temps de calcul.
- Utiliser des interruptions en espace utilisateur pour remplacer le polling.
- Simplifier pour l'utilisateur le recouvrement des communications par du calcul, afin qu'il n'ait plus besoin d'ajouter des `MPI_Test` en plein milieu des calculs.
- Améliorer la réactivité des communications sans avoir besoin d'une unité de calcul dédiée à l'attente active.

Ce stage est donc une première étape de ce projet global.

3.3 Objectifs du stage

Le premier objectif est de défricher le fonctionnement des interruptions en espace utilisateur à partir des éléments suivants :

- Le manuel Intel® de l'architecture 64 et IA-32 pour les développeurs logiciels [8].
- La présentation du mécanisme de *Sohil Mehta*, ingénieur chez Intel®, qui a développé le patch noyau. Cette présentation est une diapositive associée à des discussions sur LWN.net [15].
- Le patch du noyau Linux [9] avec ses manuels [10].

Le second objectif est de connaître les propriétés du mécanisme et de mesurer sa performance. Le troisième objectif est de ne plus avoir à faire du polling, que ce soit en dédiant un thread ou en utilisant les threads de façon opportuniste, afin de ne plus perdre de temps de calcul à poll et de résoudre les problèmes de réactivité.

Le troisième objectif est de montrer que l'utilisation des *Uintr* dans les communications HPC est possible. Pour cela, on se place dans un cadre simplifié :

- On se concentre sur les communications entre processus (IPC, Inter-Process Communication en anglais). Pour cela, on va utiliser la mémoire partagée (shm) pour communiquer au lieu du réseau *BXI*.
- On se limite aux cas généraux des communications (donc pas de gros messages, pas de communications multi-thread, etc).

On envisage donc l'intégration de ces interruptions dans le driver de mémoire partagée (shm) de *NewMadeleine*, ce qui nous permet une utilisation réelle pour voir ce que donnerait une version complète.

Pour intégrer les interruptions dans les drivers de *NewMadeleine*, il faut également permettre la progression des communications à partir d'un handler d'interruption. Le dernier objectif est de montrer que l'utilisation d'interruption permet bien d'améliorer le recouvrement des communications par du calcul.

3.4 La suite

Les objectifs suivants du projet global seront traités dans une thèse qui fait suite au stage.

4 Exploration des interruptions en espace utilisateur.

Pour cette partie, j'ai utilisé mes connaissances personnelles autour du système, du développement noyau, de Linux, ainsi que ce que j'ai appris en cours de *programmation système*, de *système d'exploitation*, d'*architecture des ordinateurs* et *Programmation des Architectures Parallèles*.

4.1 Prérequis et accès

Pour utiliser le mécanisme d'interruption en espace utilisateur, que nous allons abréger en *Uintr* dans la suite de ce document, il est nécessaire d'avoir accès à un CPU Intel® Sapphire Rapids. Du fait qu'ils étaient sortis récemment, ils étaient assez difficiles d'accès. *Atos* nous a donné accès, non sans difficulté, à une machine qui possède 2 CPU Intel® Sapphire Rapids, lesquels sont des Intel® Xeon® Platinum 8470. Les difficultés étaient liées à la disponibilité d'une machine, à trouver un endroit où l'installer et à s'assurer qu'elle soit connectée à un réseau accessible depuis l'Inria. Nous avons donc obtenu un accès VPN qui utilise l'ancien système VPN, car le nouveau ne fonctionnait pas. Nous avons eu accès à la machine environ deux mois et demi après le début du stage. La machine était déjà configurée avec le système d'exploitation Red Hat Enterprise Linux (RHEL) version 9.1 avec un noyau Linux version *5.14.0-162.6.1*. Elle possède également au moins une carte *BXI* v2 que nous n'avons pas utilisées pendant le stage.

Il faut également avoir une version patchée du noyau Linux prenant en charge le nouveau mécanisme. Cette version patchée n'est pas encore disponible dans la branche principale du noyau. Cependant, elle est accessible sur le GitHub d'Intel® [9]. Elle est basée sur la version *6.0.0*. Nous avons donc téléchargé cette version patchée, puis nous l'avons compilée et installée sur la machine, toutes les commandes nécessaires pour cela sont disponibles en

annexe. Lors de la compilation, il est nécessaire d'activer le support des *Uintr* (Voir la figure 8 en annexe). Il est également possible d'activer le support permettant à un thread bloqué, c'est-à-dire non ordonnancé ou dans un appel système interruptible, de recevoir une *Uintr*.

Le mécanisme utilise de nouvelles instructions, donc une version récente du compilateur *GCC* est nécessaire pour compiler les programmes utilisateurs qui utiliseront les *Uintr*. Il faut donc la version *11.3.0* ou plus récente de *GCC*, et sur RHEL il faut la version *12.1.1* ou supérieure. Le support n'est pas encore disponible dans d'autres compilateurs comme *LLVM-Clang* ou *ICC*. Pour compiler un programme utilisateur, il faut spécifier le flag de compilation `-muintr` pour les fichiers qui définissent un handler d'interruption ou qui utilisent les nouvelles instructions.

4.2 Fonctionnement des interruptions en espace utilisateur

Dans le cadre de ce stage, nous avons étudié en détail le fonctionnement des interruptions ordinaires ainsi que le fonctionnement des interruptions en espace utilisateur. Les *Uintr* n'étaient pas connus des équipes *Inria* et *Atos*.

4.2.1 Les interruptions

Pour commencer, nous allons voir le fonctionnement des interruptions matérielles. Nous nous concentrerons sur l'envoi d'interruptions entre deux processus fixés sur deux unités de calcul. Pour ces interruptions les CPU disposent d'une unité dédiée à leur traitement, l'*APIC* pour "Advanced Programmable Interrupt Controller". Cette *APIC* permet au système d'enregistrer un handler pour chaque interruption. Le noyau définit un tableau appelé *IDT* pour "Interrupt Descriptor Table" qui contient 256 entrées, correspondant à 256 interruptions possibles. Les indices de l'*IDT* sont des valeurs comprises entre 0 et 255, que l'on appelle également **vecteurs** d'interruption. Parmi eux, les vecteurs entre 0 et 31 sont réservés pour les exceptions et les interruptions système, les vecteurs entre 32 et 127 sont dédiés aux interruptions liées aux périphériques, le vecteur 128 est réservé pour les appels système, et les vecteurs entre 129 et 255 sont destinés à des utilisations diverses.

Il est important de savoir que chaque unité de calcul (processeur logique) possède un *APIC ID* physique. À titre informatif, le *core ID* est un sous-ensemble de l'*APIC ID*.

Pour déclencher une interruption, il y a quatre possibilités :

1. Une exception déclenchée par un processeur (e.g. une division par zéro, un défaut de segmentation...).
2. Une instruction comme `INT80 numSysCall` pour déclencher un appel système, ou bien `INT3` pour définir un point d'arrêt, ou encore `INT0`, `BOUND` et `INT n`.

3. Des broches du CPU dédiées à la réception d'interruptions lancées à partir d'un périphérique.
4. Demander à l'*APIC* elle-même grâce à un registre *ICR* pour "Interrupt Command Register". Il existe un *ICR* par vecteur, donc il faut écrire l'*APIC ID* du destinataire dans le *ICR* du vecteur que l'on veut déclencher. Seul le CPU et le noyau peuvent modifier les *ICR*.

On voit bien que les *IRQ* fonctionnent au niveau du noyau et du CPU.

Nous allons voir un exemple d'envoi d'*IRQ* entre deux unités de calcul en cours d'exécution. Tout d'abord, l'initialisation des *IRQ* se fait au démarrage du système et consiste principalement à définir les handlers noyau dans l'*IDT*. Il faut aussi définir quel vecteur nous voulons utiliser, le code du handler noyau qui sera invoqué, comment l'utilisateur va contacter le système et comment faire l'identification du récepteur. Pour cela, on peut appliquer un patch au noyau ou faire un module noyau.

Dans notre exemple, nous allons supposer que nous avons déjà patché le noyau en ajoutant des appels système et que nous avons choisi un vecteur.

Nous allons maintenant voir les étapes de l'envoi d'une *IRQ*, illustrées sur la figure 1 :

- ① Le récepteur fait un appel système pour indiquer au noyau comment il veut être averti d'une interruption (e.g. un descripteur de fichiers qu'il va lire, un appel système bloquant, une zone mémoire où lire, un handler utilisateur, etc.).
- ② L'émetteur peut donc avertir le noyau qu'il faut envoyer une interruption. Pour cela, il peut utiliser un appel système ou écrire dans un descripteur de fichiers, par exemple.
- ③ Le noyau détermine l'unité de calcul où se trouve le récepteur. Pour cela, il peut utiliser par exemple un *PID* pour "Processus ID" donné par l'émetteur ou autre. Ainsi, il peut déterminer l'*APIC ID* de l'unité de calcul à interrompre.
- ④ Le noyau écrit donc l'*APIC ID* dans le *ICR* d'un vecteur déterminé à l'avance. L'émetteur reprend alors la main après un autre changement de contexte.
- ⑤ L'*APIC* va donc interrompre le récepteur qui va alors stopper son exécution et passer dans le noyau. Une fois dans le noyau, le handler va se déclencher et exécuter le code prévu au préalable (écriture dans un descripteur de fichiers, écriture dans une zone mémoire, déclenchement d'un handler utilisateur, etc.).

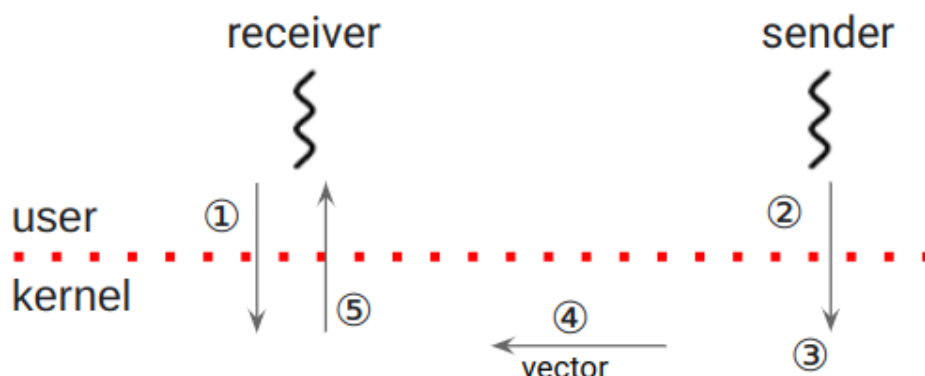


FIGURE 1 – L'envoi d'une interruption ordinaire

Lors du déclenchement du handler noyau, certains registres actuels sont sauvegardés, tels que le pointeur de pile `RSP`, le registre d'états `RFLAGS`, le registre `CS` et le registre de pointeur d'instruction `RIP`. Cette sauvegarde est réalisée en les empilant dans une nouvelle pile. Le vecteur de l'interruption est aussi empilé en tant que code erreur (*errorCode*).

Une fois que le handler noyau a fini de s'exécuter, il doit exécuter l'instruction `iret` qui a pour effet de dépiler les registres sauvegardés et de les restaurer.

Il est possible de masquer les interruptions grâce à deux instructions utilisables seulement par le noyau, qui sont `cli` et `sti`. Ces instructions permettent de modifier le flag *IF* pour "Interrupt Flag" qui se trouve dans le registre d'états de l'unité de calcul, `RFLAGS` (aussi nommé `EFLAGS` sur les architectures 32 bits). La liste des instructions pour les *IRQ* se trouve dans le tableau 1.

Comme nous l'avons vu, ce mécanisme fonctionne totalement dans le noyau du système. Dans notre exemple, il faut au minimum deux changements de contexte (context switch) pour le récepteur et l'émetteur, et peut-être même plus si le récepteur doit déclencher un handler côté utilisateur.

4.2.2 Les `uintr`

Le mécanisme d'interruption en espace utilisateur utilise cinq nouvelles instructions qui sont listées dans le tableau suivant 1. Les deux premières, `clui` et `stui` qui sont analogues à `cli` et `sti`, permettent le masquage des interruptions. En effet, tout comme les interruptions ordinaires qui ont un flag *IF* pour activer ou désactiver les interruptions, les *Uintr* ont un flag *UIF* pour "User Interrupt Flag". L'instruction suivante, `testui`, permet à l'utilisateur de savoir si les interruptions sont masquées ou non. Cette instruction existe car l'utilisateur n'a pas accès directement au *UIF*, contrai-

rement aux interruptions ordinaires qui ont un accès direct à *IF* puisqu'elles fonctionnent dans le noyau.

L'instruction suivante, **uiret**, fonctionne de manière similaire à celle des interruptions ordinaires (**iret**), mais elle n'utilise pas les mêmes registres et surtout elle est utilisable en espace utilisateur. Enfin, la dernière instruction permet d'envoyer une *Uintr* en utilisant un indice, que nous allons voir dans la section 4.2.4.

Interruption (IRQ)	Interruption utilisateur (<i>Uintr</i>)
cli (CL ear IF)	clui (CL ear UIF)
sti (SeT IF)	stui (SeT UIF)
	testui (Read UIF)
iret (Interrupt RE Turn)	uiret (User Interrupt RE Turn)
APIC pins, APIC ICR, INT n , INT3 , INT0 , BOUND et INT80 n	sendipi <uipi_index>

TABLE 1 – Instructions des interruptions et des interruptions en espace utilisateur

Le mécanisme est également accompagné de six nouveaux *registres d'états*, appelés registres *MSR* pour "Model-Specific Registers". Ces registres sont modifiés par le noyau grâce à des appels système que l'utilisateur effectue pour initialiser les *Uintr* et sont utilisés par le CPU. Le tableau 2 décrit ces registres, et nous expliquerons leur utilité par la suite.

Nom du registre	Description
IA32_UINTR_STACKADJUST	Utilisé par le récepteur pour définir l'adresse de la pile alternative
IA32_UINTR_HANDLER	Utilisé par le récepteur pour définir l'adresse du handler <i>Uintr</i>
IA32_UINTR_MISC	Utilisé par l'émetteur pour définir la taille du <i>UITT</i> et par le récepteur pour que l' <i>APIC</i> connaisse le vecteur d'interruption ordinaire qu'il doit reconnaître pour déclencher le handler <i>Uintr</i> et le dernier bit, pour le flag de masquage <i>UIF</i>
IA32_UINTR_PD	Utilisé par le récepteur pour définir l'adresse du <i>UPID</i>
IA32_UINTR_RR	Utilisé par l' <i>APIC</i> pour <i>lister</i> les vecteurs <i>Uintr</i> qu'il doit envoyer au récepteur. Correspond aux derniers 64 bits du <i>UPID</i>
IA32_UINTR_TT	Utilisé par l'émetteur pour définir l'adresse du <i>UITT</i>

TABLE 2 – Liste des six registres d'états des *Uintr*

4.2.3 Capacités présentes et futures

Le mécanisme d'interruption en espace utilisateur a une interface pour l'utilisateur similaire aux signaux 2.8. Nous allons donc voir les capacités des *Uintr* en les comparant à celles des signaux. Tout d'abord, le fonctionnement des *Uintr* se fait au niveau des threads, tandis que celui des signaux se fait au niveau du processus. Il est possible d'avoir un fonctionnement qui se rapproche d'un fonctionnement par threads en utilisant plusieurs options.

Avec les *Uintr*, il est possible de définir un handler différent pour chaque thread d'un processus, alors que pour les signaux, on peut définir un seul handler pour tous les threads d'un processus. Par contre, avec les signaux, il est possible de définir un handler différent pour chaque signal, ce qui n'est pas possible avec les *Uintr*. Pour les *Uintr*, il faut gérer cette différenciation manuellement en appelant la fonction correspondant au vecteur reçu.

Il existe 64 signaux possibles, parmi lesquels les 32 premiers ont une signification particulière. En revanche, pour les *Uintr*, il existe également 64 vecteurs possibles, entre 0 à 63, qui n'ont aucune signification particulière par défaut.

Pour les *Uintr*, le masquage se fait via une instruction, tandis que pour les signaux, il faut effectuer un appel système pour les masquer.

L'envoi d'un signal se fait par le noyau suite à une exception, une décision du noyau ou la demande d'un processus grâce à un appel système (`kill(signum)` ou `tgkill(signum)`). Pour les *Uintr*, l'envoi peut se faire depuis un autre processus ou depuis le noyau, et dans le futur, il pourra également se faire depuis un périphérique.

Avec les signaux, le handler peut être déclenché que le processus cible soit endormi ou non, tandis que pour les *Uintr*, c'est différent. Il faut que le thread soit en espace utilisateur pour recevoir une interruption. Sinon, l'interruption sera reçue quand le thread revient en espace utilisateur. On a vu précédemment, dans les prérequis, section 4.1, qu'une fonctionnalité existe lors de la compilation du noyau pour autoriser l'interruption d'un thread bloqué. Si la fonctionnalité est activée, il est donc possible d'interrompre un thread qui n'est pas ordonnancé ou qui est en train de faire un appel système interruptible, et ainsi le passer en espace utilisateur pour qu'il puisse recevoir l'interruption. Pour utiliser cette fonctionnalité, l'utilisateur doit renseigner un flag au moment de définir le handler.

Il existe donc trois flags :

- `UINTR_HANDLER_FLAG_WAITING_ANY` qui active la fonctionnalité,
- `UINTR_HANDLER_FLAG_WAITING_RECEIVER` et `UINTR_HANDLER_FLAG_WAITING_SENDER` qui s'ajoutent au précédent pour préciser si c'est l'émetteur ou le récepteur qui prendra en charge le surcoût du passage par le noyau.

Pour les signaux, le déclenchement du handler est géré par le noyau qui sauvegarde l'état du processus, définit une pile alternative si besoin, change

de contexte et appelle le handler utilisateur. Pour les *Uintr*, le déclenchement du handler est effectué par le CPU, il est donc très sommaire :

- changer la pile si une pile alternative est disponible dans le registre `IA32_UINTR_STACKADJUST`,
- empiler l'ancien pointeur de pile, le registre d'états de l'unité de calcul, le registre de pointeur d'instruction `RIP` et le vecteur *Uintr*
- aller à l'adresse du handler utilisateur disponible dans le registre `IA32_UINTR_HANDLER`.

C'est donc à l'utilisateur qu'appartient la responsabilité de sauvegarder les registres généraux, les registres vectoriels (SIMD)... et de les restaurer à la sortie du handler. Le compilateur permet déjà de sauvegarder les registres généraux avec le flag `general-regs-only`. Cependant, pour les registres vectoriels, il faut les sauvegarder avant de les utiliser. Il faut faire attention avec les opérations sur les chaînes de caractères de la *libc* car les fonctions `memcpy`, `memmove`, `memset` et `memcpy` utilisent des registres vectoriels par défaut. Le compilateur fournit le flag `-minline-all-stringops` qui permet de les remplacer par des versions inline de ces opérations afin de ne plus utiliser de registres vectoriels.

Une fois que le handler a fini son exécution, il faut s'occuper du retour. Pour les signaux, c'est le noyau qui s'en occupe, tandis que pour les *Uintr*, il incombe à l'utilisateur d'utiliser l'instruction `uiret`. Donc, l'unité de calcul va dépiler le vecteur et les registres qui suivent pour les restaurer, ce qui permettra au code de continuer là où il en était.

Que ce soit dans un handler de signal ou dans un handler d'interruption utilisateur, on a la même contrainte : on ne peut pas faire d'attente, donc on peut seulement appeler des fonctions et des appels système dits "async safe".

4.2.4 Exemple de fonctionnement

Nous allons maintenant voir un exemple d'initialisation des *Uintr* illustré sur la figure 2 :

- ① Le récepteur enregistre auprès du noyau un handler d'interruption qu'il a défini à l'aide de l'appel système `uintr_register_handler(ui_handler)`. Le noyau va enregistrer ce handler dans le registre `IA32_UINTR_HANDLER` et va initialiser une zone mémoire nommée *UPID* pour "User Posted Interrupt Descriptor". Ce *UPID* permet au mécanisme des *Uintr* de manipuler des informations propres à ce thread, essentielles pour l'envoi d'*Uintr*. L'adresse du *UPID* est enregistré dans le registre `IA32_UINTR_PD`.
- ② Le récepteur donne au noyau un vecteur *Uintr* entre 0 et 63 qu'il souhaite recevoir (8 dans la figure), à l'aide de l'appel système

`uvec_fd <- uintr_vector_fd(8)` . Le noyau lui renvoie un descripteur de fichier qui pointe vers une structure contenant à la fois le vecteur et l'adresse du *UPID*.

- ③ Le récepteur envoie ce descripteur de fichier aux émetteurs potentiels, un seul dans notre cas. Nous verrons comment partager ce descripteur de fichier dans la section 4.2.5.

- ④ L'émetteur va s'enregistrer auprès du noyau grâce au descripteur de fichier en appelant l'appel système

`uipi_index <- uintr_register_sender(uvec_fd)` . Pour ce faire, le noyau possède un tableau *UITT* pour "User Interrupt Target Table" qui fait une taille de 256 entrées par défaut. L'adresse de ce tableau doit être enregistrée dans le registre `IA32_UINTR_TT` , et sa taille doit être placée dans les quatre premiers octets du registre `IA32_UINTR_MISC` . Ainsi, la taille du *UITT* peut varier en fonction des besoins. Chaque entrée de ce tableau *UITT* correspond à une zone mémoire nommée *UITTE* pour "User Interrupt Target Table Entry". Le noyau va rechercher une entrée libre dans le *UITT* et remplir l'*UITTE* avec le vecteur *Uintr* et l'adresse du *UPID* obtenus grâce au descripteur de fichier. Il va également mettre un vecteur d'interruption ordinaire dans le cinquième octet du registre `IA32_UINTR_MISC` , ce vecteur "ordinaire" étant dédié aux *Uintr* et a pour valeur 236. Pour finir il retourne l'indice de l'entrée à l'émetteur.

- ⑤ L'utilisateur peut maintenant envoyer autant d'interruptions que nécessaire, totalement depuis l'espace utilisateur.

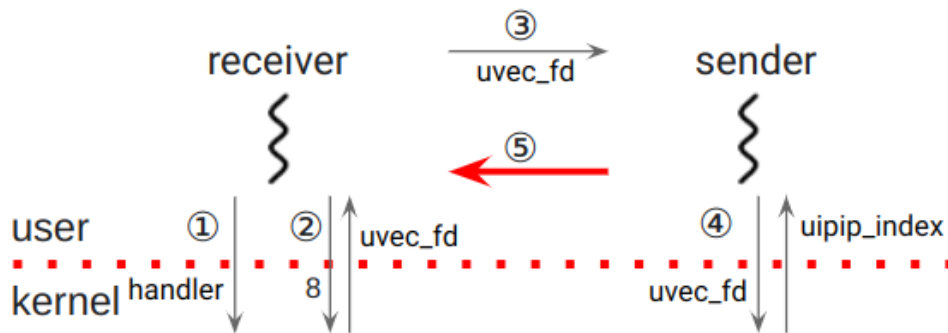


FIGURE 2 – Phase d'initialisation des *Uintr*

Maintenant que nous avons vu l'initialisation, nous allons pouvoir voir comment l'envoi d'interruption est réalisé, comme illustré dans la figure 3 :

- ① L'émetteur utilise l'instruction `senduipi` avec l'indice récupéré lors de l'initialisation. L'unité de calcul peut accéder au tableau *UITT* qui se trouve dans le registre `IA32_UINTR_TT` . Elle peut vérifier que

l'indice se trouve bien dans le tableau, en utilisant la taille qui se trouve dans le registre `IA32_UINTR_MISC`. De plus, il s'assure que l'entrée dans le *UITT* est valide. À partir de l'indice donné à l'instruction, l'émetteur peut récupérer l'adresse de l'*UPID* et le vecteur *Uintr* à envoyer au récepteur. Dans la zone mémoire de l'*UPID*, il écrit le vecteur *Uintr* à envoyer et détermine s'il faut également envoyer une interruption ordinaire. En effet, les interruptions en espace utilisateur utilisent les interruptions ordinaires, et l'envoi de celles-ci dépend de si une interruption ordinaire n'a pas déjà été envoyée. Dans notre cas, nous considérons que c'est le premier envoi.

- ② Donc, l'émetteur va récupérer le vecteur d'interruption ordinaire dans le registre `IA32_UINTR_MISC` et sélectionner le registre *ICR* correspondant au vecteur "ordinaire". Ensuite, il va récupérer l'*APIC ID* dans l'*UPID* et l'écrire dans le registre *ICR*. Cela permettra d'envoyer l'interruption ordinaire à l'unité de calcul correspondante.
- ③ L'*APIC* va recevoir l'interruption, puis elle va comparer le vecteur d'interruption ordinaire reçu avec celui qui se trouve dans l'*UPID*, qu'il connaît grâce au registre `IA32_UINTR_PD`. Si les vecteurs sont identiques, l'*APIC* va pouvoir déclencher le système de traitement des *Uintr*, sinon, elle va utiliser le mécanisme habituel des *IRQ*.
- ④ Le système de traitement des *Uintr* va donc indiquer dans l'*UPID* que l'interruption ordinaire a déjà été envoyée, puis va commencer le déclenchement des handlers pour tous les vecteurs *Uintr* reçus, en partant du plus grand, 63, jusqu'au plus petit, 0. `IA32_UINTR_RR` est utilisé à cette étape. Dans notre exemple, nous en avons un seul qui est 8. L'*APIC* va donc ordonner à l'unité de calcul de changer de pile si une pile alternative est disponible dans le registre `IA32_UINTR_STACKADJUST`, empiler les registres nécessaires, empiler le vecteur *Uintr* (donc 8 dans l'exemple) et aller à l'adresse du handler qui est disponible dans le registre `IA32_UINTR_HANDLER`.

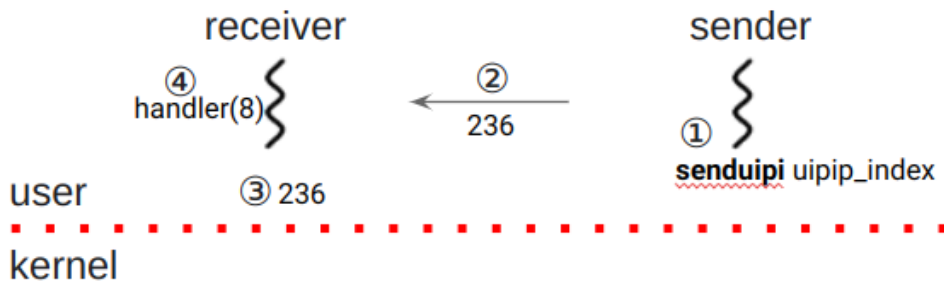


FIGURE 3 – L'envoi d'une *Uintr*

Le mécanisme habituel des *IRQ* peut être utilisé dans le cas où l'*Uintr*

est déclenchée depuis le noyau et dans le cas où le thread destinataire n'est pas ordonnancé ou est dans un appel système interruptible.

4.2.5 Partage du descripteur de fichier

Pour partager le descripteur de fichier, il y a plusieurs façons possibles.

Dans nos tests du mécanisme entre processus, nous avons utilisé l'héritage des descripteurs de fichier. Dans nos tests du mécanisme entre threads, nous avons utilisé un appel système qui permet d'enregistrer tous les threads du processus comme émetteurs, sans utiliser le descripteur de fichier, mais il est possible d'utiliser une variable globale pour y mettre le descripteur de fichier.

Dans le cas où les deux processus sont indépendants, on peut utiliser l'appel système `pidfd_getfd` qui permet de dupliquer le descripteur de fichier d'un autre processus s'il a le même propriétaire et si l'on connaît son *PID* et le numéro du descripteur de fichier. Nous avons donc testé en partageant le numéro de descripteur de fichier avec un pipe. On aurait aussi pu passer par un fichier ou par des sockets.

Par la suite, dans *NewMadeleine*, nous avons utilisé le système d'URL à la connexion que nous verrons dans la section `/* TODO : ref */`.

4.3 Tests du mécanisme

Nous avons donc effectué des tests du mécanisme avec des exemples minimaux de communication entre processus. Nous examinerons les plus pertinents dans cette section.

Tout d'abord, des tests pour mesurer le temps d'envoi d'interruption en espace utilisateur entre deux threads et deux processus.

Pour le test d'envoi entre deux threads, il commence par créer un nouveau thread qui va commencer par se "bind" à une unité de calcul. Nous le verrons en détail dans la section 4.5.

Il va ensuite enregistrer un handler d'interruption, démasquer les interruptions, enregistrer tous les threads du processus comme émetteurs avec l'appel système `uipi_index <- uintr_register_self(vector)`, (`uipi_index` est global au processus) et attendre grâce à une boucle.

Pendant ce temps-là, le thread principal attend une seconde. Ce temps est arbitraire et laisse le temps au thread d'enregistrer un handler d'interruption. Après son attente, il va se "bind" à une unité de calcul puis envoyer une interruption. Pour cela, nous commençons par enregistrer le temps processeur actuel, puis nous utilisons l'instruction `senduipi` avec la variable globale `uipi_index`. Une fois que le handler d'interruption se déclenche, on enregistre le temps actuel du processeur pour ensuite calculer la différence. Ce test est capable de faire cette mesure plusieurs fois. Il finit par imprimer

les différences de temps dans la console et par désallouer et de-enregistrer les *Uintr* et les autres structures.

Pour le teste d'envoi entre deux processus, il commence par créer un pipe pour l'envoi du descripteur de fichier et une zone de mémoire partagée pour stocker les mesures de temps. Puis le processus se fork en deux, le premier devient l'émetteur et le second le récepteur. Comme pour la version avec threads, ils se "bind" à une unité de calcul. Le récepteur enregistre un handler d'interruption et récupère un descripteur de fichier avec l'appel système `uvec_fd <- uintr_vector_fd(vector)` qu'il envoie dans le pipe avant d'attendre grâce à une boucle. L'émetteur reçoit le numéro du descripteur de fichier, il connaît déjà le *PID* du processus grâce au fork, il peut donc utiliser `pidfd_getfd` pour dupliquer le descripteur de fichier. Avec ce descripteur, il s'enregistre en tant qu'émetteur d'*Uintr*. La mesure du temps et l'envoi d'interruption fonctionnent de la même manière que pour la version avec threads. On peut aussi effectuer la mesure plusieurs fois et on imprime et termine proprement.

On verra les résultats de ces mesures dans la section 4.6.

Un test où le thread s'auto-interrompt tout simplement en faisant un `uipi_index <- uintr_register_self(vector)`, puis un `senduipi uipi_index`, et on fait la mesure de la même façon que pour les autres tests.

Pour les tests avec la pile alternative, nous avons un test très simple qui est basé sur celui qui s'auto-interrompt, et nous modifions les tests d'envoi entre deux processus ou threads pour mesurer l'impact les performances. Il est bien sûr possible d'activer ou non l'utilisation de la pile alternative. Pour définir cette nouvelle pile, on le fait juste après avoir enregistré le handler d'interruption et avant le démasquage.

Un test de démasquage des *Uintr* dans un handler d'interruption. Il nous permet de voir que c'est tous à fait possible et cela pose des problématiques similaires à celles des signaux.

Un test consiste à envoyer plusieurs interruptions d'affilée. Il nous permet de voir qu'il y a une différence de comportement par rapport aux signaux. Quand on reçoit une interruption et que le handler d'interruption est déclenché, le comportement est le même, c'est-à-dire que les interruptions vont s'écraiser et le handler d'interruption se déclenchera à nouveau une fois. La différence réside dans le fait que si l'on effectue plusieurs interruptions avant que le handler d'interruption ne se déclenche, les interruptions s'écraient également à ce moment-là. Ainsi, le handler d'interruption ne se déclenchera qu'une seule fois lorsque l'on effectue deux interruptions successives. Avec les signaux, en revanche, le handler de signal se déclencherait deux fois pour deux émissions de signal successives car le handler est déjà déclenché dès le premier signal.

Grâce à un test, nous avons constaté qu'actuellement on peut enregistrer plusieurs fois le même descripteur de fichier, ce qui mène à des doublons

dans le tableau *UITT* avec plusieurs *UITTE* pour le même couple vecteur / *UPID*. Cette limitation de l'implémentation noyau est documentée dans le patch et accompagnée d'un commentaire *TODO*.

4.4 Correction d'un bug dans le patch noyau

En manipulant le mécanisme, nous sommes tombés sur un bug qui concerne l'utilisation d'une pile alternative. Comme pour les signaux, il est possible de définir une pile alternative. Cette nouvelle pile est utilisée au moment où le handler est déclenché. L'interface pour définir la pile alternative est la même pour les signaux et les *Uintr*. Dans les manuels (des signaux et des *Uintr*), il est bien indiqué que l'utilisateur doit lui-même allouer une zone mémoire consacrée à la nouvelle pile, et il a également la responsabilité de libérer la mémoire une fois le handler dé-enregistré. Il est bien indiqué que l'utilisateur doit donner l'adresse de début (*base address*) de la zone mémoire, ainsi que la taille, à l'appel système qui définit la pile alternative. Pour les *Uintr*, l'appel système est `uintr_alt_stack(spAddress, size)`. Il faut noter qu'une pile empile les éléments vers le haut, c'est-à-dire que l'adresse du pointeur de pile décroît à l'ajout d'un élément. Donc, pour utiliser la zone mémoire dédiée à la pile, il faut partir de la fin. Du côté du noyau, le mécanisme des signaux garde en mémoire l'adresse et la taille de la pile pour le moment où le handler doit être déclenché. Le calcul du pointeur de pile de la nouvelle pile se fait donc juste avant de déclencher le handler. Pour les *Uintr*, le noyau se contente juste d'enregistrer l'adresse dans le registre dédié `IA32_UINTR_STACKADJUST`, et le processeur utilise l'adresse telle quelle comme pointeur de pile. On est donc confronté à un bug de débordement mémoire car on part du début de la zone mémoire. Il y a bien un test dans le patch du noyau qui vérifie ce cas, mais mal. Nous sommes tombés sur ce bug au moment d'utiliser les *Uintr* dans *NewMadeleine*, qui manipule bien plus la mémoire qu'un simple test. Nous nous sommes retrouvés avec des problèmes de corruption mémoire, des défauts de segmentation et des "double free detected". Nous avons donc corrigé le test du patch du noyau ainsi que l'appel système. Pour ce faire, on ajoute la taille de la zone mémoire à l'adresse avant de l'écrire dans le registre `IA32_UINTR_STACKADJUST`. Nous avons donc effectué une *Pull request* [14] sur le dépôt GitHub du patch, et à l'heure où nous écrivons ce document, il n'a toujours pas été appliqué par Intel®.

4.5 Mesure de la latence

Pour mesurer la latence entre le moment où on envoie une interruption et où l'interruption est reçue par le handler, on fait deux mesures de temps. Pour faire les mesures de temps, on utilise `clock_gettime` qui utilise l'instruction `rdtsc` et retourne le temps actuel du processeur. On fait une

première mesure juste avant d'envoyer une interruption et une seconde au tout début du handler. Pour obtenir la latence, on a juste à soustraire la première mesure à la seconde.

Nous avons regardé le code assembleur pour nous assurer de la mesure. Pour l'envoi, que l'on peut voir en listing 1, on peut constater qu'entre la mesure et l'instruction d'envoi, il y a seulement une lecture mémoire qui n'est pas très coûteuse. Du côté de la réception, que l'on peut voir en listing 2, on voit la sauvegarde des registres généraux au début du handler. Cette sauvegarde ajoute un petit surcoût, mais il est obligatoire. Le compilateur *GCC* nous force à mettre le flag `general-regs-only` pour compiler un handler d'interruption et donc sauvegarder les registres généraux. On peut voir la déclaration d'un handler avec la mesure du temps en listing 3.

```
1      call    clock_gettime@PLT
2      movq    uipi_index(%rip), %rax
3      senduipi %rax
```

Listing 1 – Code assembleur de l'envoi d'*Uintr*

```
1  handler:
2  .LFB203:
3      .cfi_startproc
4      .cfi_def_cfa_offset 16
5      pushq   %r11
6      .cfi_def_cfa_offset 24
7      .cfi_offset 11, -24
8      pushq   %r10
9      .cfi_def_cfa_offset 32
10     .cfi_offset 10, -32
11     pushq   %r9
12     .cfi_def_cfa_offset 40
13     .cfi_offset 9, -40
14     pushq   %r8
15     .cfi_def_cfa_offset 48
16     .cfi_offset 8, -48
17     pushq   %rdi
18     .cfi_def_cfa_offset 56
19     .cfi_offset 5, -56
20     movl     $1, %edi
21     pushq   %rsi
22     .cfi_def_cfa_offset 64
23     .cfi_offset 4, -64
24     leaq     ts2(%rip), %rsi
25     pushq   %rcx
26     .cfi_def_cfa_offset 72
27     .cfi_offset 2, -72
28     pushq   %rdx
29     .cfi_def_cfa_offset 80
30     .cfi_offset 1, -80
31     pushq   %rax
32     .cfi_def_cfa_offset 88
33     .cfi_offset 0, -88
```

```

34     subq    $8, %rsp
35     .cfi_def_cfa_offset 96
36     cld
37     call    clock_gettime@PLT

```

Listing 2 – Code assembleur de l’handler *Uintr*

```

1  __attribute__((target("general-regs-only")))
2  __attribute__((interrupt))
3  void handler(struct __uintr_frame* ui_frame, u64 vector) {
4      clock_gettime(CLOCK_MONOTONIC, &ts2);

```

Listing 3 – Déclaration du handler *Uintr*

Nous faisons la mesure de la même façon pour les signaux.

Pour ne pas perturber la mesure, il faut "bind" les threads à une unité de calcul. "bind" consiste à demander au noyau de toujours ordonnancer le thread sur la même unité de calcul. Pour ce faire nous utilisons la bibliothèque *hwloc* [7]. Il est donc important de "bind" les threads pendant les mesures car dans le cas contraire, le noyau va changer le thread d’unité de calcul, ce qui va amener à un surcoût non négligeable. En effet, le changement d’unité est un peu coûteux et invalide certains caches. "bind" les threads nous permet aussi de contrôler le placement de ceux-ci, c’est-à-dire si on met le thread émetteur et le thread récepteur sur deux unités de calcul proches ou distantes. Comme nous le verrons dans la section 4.6, le placement a un impact sur la latence des *Uintr*.

Il est aussi important de fixer la fréquence de tous les core du CPU pour avoir des mesures reproductibles.

4.6 Performances

Dans cette section, nous allons voir des mesures de la latence des *Uintr* dans différents contextes. La fréquence des unités de calcul est fixée à 2 GHz lors des mesures. Les mesures faites avec le turbo boost activé montent à 3.8 GHz. Les mesures de latence sont en nanosecondes. Il est intéressant de noter que sur une unité de calcul cadencée à une fréquence de 2 GHz, elle exécute environ deux instructions par nanoseconde. Certaines mesures, que ce soit avec les *Uintr* ou les signaux, sont très élevées mais sont très peu nombreuses. Elles sont certainement liées au système. Dans les graphiques, nous avons donc coupé les valeurs qui dépassent les 8000 nanosecondes. Nous avons défini trois placements à partir de la topologie de la machine fournie par *Atos*. Vous pouvez trouver la topologie de la machine sur la figure 9.

Les trois placements sont les suivants :

- Le placement "proche" consiste à placer les threads sur deux unités de calcul proches mais pas dans le même core.
- Le placement "éloigné" consiste à placer les threads sur deux unités de calcul qui sont dans le même CPU et qui sont éloignées.

- Le placement "très éloigné" consiste à placer les threads sur deux unités de calcul qui se trouvent sur deux CPU différents.

Les mesures que nous allons présenter ont été faites sans que la pile alternative ne soit activée. Nous avons bien fait des mesures avec la pile alternative des *Uintr* activée et nous n'avons vu aucune différence, car l'utilisation de cette pile amène seulement à une copie de la mémoire d'un registre à un autre si le registre contient une adresse, ce qui est très peu coûteux.

Dans nos graphiques, nous appelons une mesure une itération. Nous faisons donc un million d'itérations et nous n'affichons pas la première, car cette itération est énormément perturbée notamment par le coût de chargement des caches.

Dans les graphiques, nous avons représenté la mesure de la latence par des points bleus pour les signaux et par des points rouges pour les *Uintr*. Nous ne cherchons pas à expliquer les mesures des signaux, elles sont juste là pour comparer les *Uintr* avec un mécanisme qui passe par le noyau.

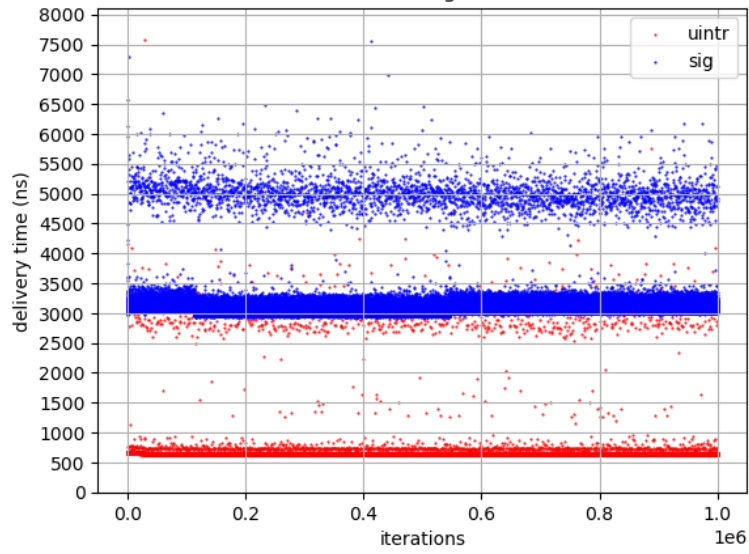
Sur le graphique 4a, les itérations sont représentées sur l'axe des abscisses et la latence sur l'axe des ordonnées. Bien sûr, plus la latence est basse, mieux c'est. Nous observons que le mécanisme a une latence d'environ 642 nanosecondes, ce qui est environ 4.75 fois plus rapide que les signaux.

On voit qu'il y a deux groupes de mesures :

- Un autour des 642 nanosecondes qui comporte la majorité des mesures. On le voit bien dans le tableau 4b qui se trouve juste en dessous du graphique. En effet, entre la mesure minimum et la mesure à 95%, il y a une différence de 19 nanosecondes. On peut voir cette distribution aussi dans l'histogramme sur la figure 5. Cet histogramme porte sur des plages de 300 nanosecondes. On voit bien pour les *Uintr*, en orange, que la majorité se trouve sur la bande entre 300 et 600 nanosecondes.
- Un autre autour de 2700 nanosecondes qui correspond certainement au moment où l'interruption n'a pas pu être reçue car le thread était dans le noyau.

Pour un mécanisme qui fonctionne au niveau des instructions, on pourrait s'attendre à une latence moins grande, mais le mécanisme est bien plus rapide que le fait de passer par le système.

999999 iterations 2 threads with core binding and without the first iteration (Next to)



(a)

	mean	std	min	10%	50%	95%	max
sig	3065	160	2920	3001	3054	3157	68532
uintr	643	90	629	638	642	648	65973

(b)

FIGURE 4 – Mesures de latence entre deux threads avec un placement proche

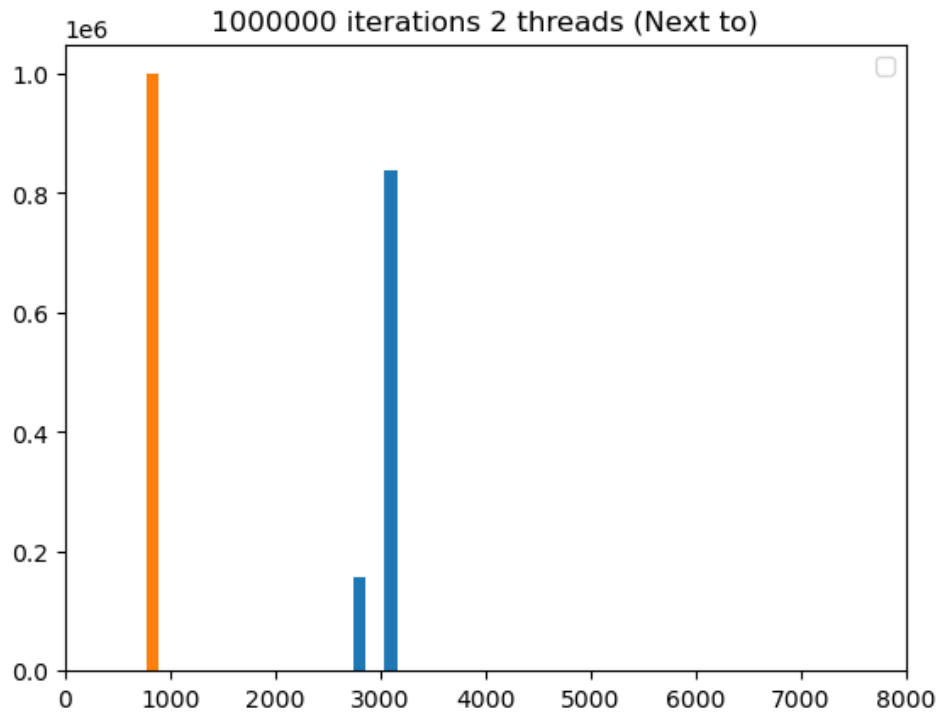


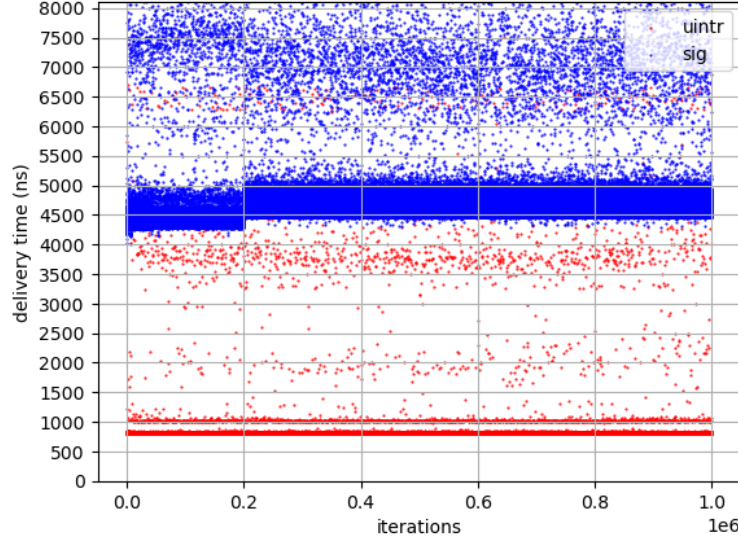
FIGURE 5 – Histogramme de distribution des mesures de la latence entre deux threads avec un placement proche

Nous avons fait les mêmes mesures entre deux processus et les valeurs sont très similaires, peu importe le placement, car le mécanisme fonctionne au niveau des threads. On peut retrouver ces mesures en annexe sur la figure 10.

En comparant les mesures de latences entre le placement proche et éloigné on voit une petite différence de 10 nanosecondes de plus, on peut voir ça sur la figure 11.

Quand on compare entre le placement proche et très éloigné on voit une grande différence qui est due notamment au fait de passer d'un noeud mémoire NUMA à un autre. On a donc une différence d'environ 172 nanosecondes de plus. Le graphique nous montre également une plus grande dispersion de la latence avec la majorité qui est toujours en-dessous de 1000 nanosecondes. On retrouve ces mesures sur la figure 6 juste en-dessous.

999999 iterations 2 threads with core binding and without the first iteration (Very Far)



(a)

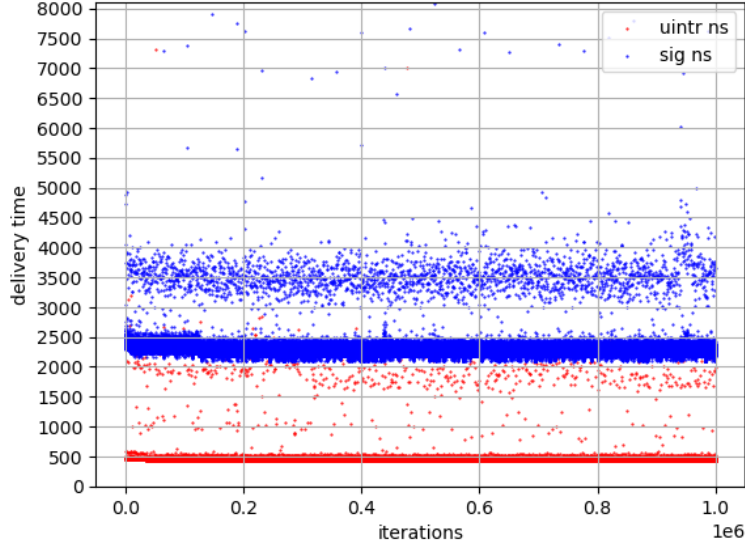
	mean	std	min	10%	50%	95%	max
sig	4601	403	4009	4396	4561	4796	58200
uintr	818	136	801	810	814	820	58072

(b)

FIGURE 6 – Mesures de latence entre deux threads avec un placement très éloigné

Quand on augmente la fréquence des unités de calcul la latence diminue. Pour se faire on active le turbo boost du CPU. On le voit sur les mesures de la figure 7 pour le placement proche mais c'est également le cas pour le placement éloigné et le très éloigné en figure 12 et figure 13.

999999 iterations 2 threads with core binding and without the first iteration (Next to)



(a)

	mean	std	min	10%	50%	95%	max
sig	2262	294	2081	2216	2251	2357	283171
uintr	442	313	426	437	440	450	311486

(b)

FIGURE 7 – Mesures de latence entre deux threads avec un placement proche et le turbo boost activé

Nous n'avons fait aucune mesure de latence pour les threads bloqués et interruptibles.

5 Intégration dans NewMadeleine

TODO : toute cette section n'est pas faite

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 reçois une iterruption est détermine si l'émission est fini, petit paquet, ou si il faut envoyer la suite, paquet large.

5.3 Suite de testes

Tous les testes 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

```
1 git clone git@github.com:intel/uintr-linux-kernel.git
```

Listing 4 – Récupération du code source du noyau Linux patché

```
1 # Configurer le noyau sur la machine cible
2 cp -v /boot/config-$(uname -r) .config
3 make menuconfig
4
5 # Désactiver les clef de signature
6 scripts/config --disable SYSTEM_TRUSTED_KEYS
7 scripts/config --disable SYSTEM_REVOCATION_KEYS
8
9 # Compiler
10 make -j 16
```

Listing 5 – Configurer et compiler le noyau

```
1 sudo make modules_install
2 sudo make install
```

Listing 6 – Installer le noyau

```
1 # À faire la premier fois
2 sudo grub2-mkconfig -o /boot/grub2/grub.cfg
3
4 sudo grubby --set-default /boot/vmlinuz-6.0.0+
```

Listing 7 – Utiliser le noyau installé

```
1 make -C tools/testing/selftests TARGETS=uintr run_tests
```

Listing 8 – Lancer les tests *Uintr* du noyau (le noyau doit être installé)

```
1 sudo grubby --remove-kernel /boot/vmlinuz-6.0.0+
```

Listing 9 – Supprimer le noyau installé

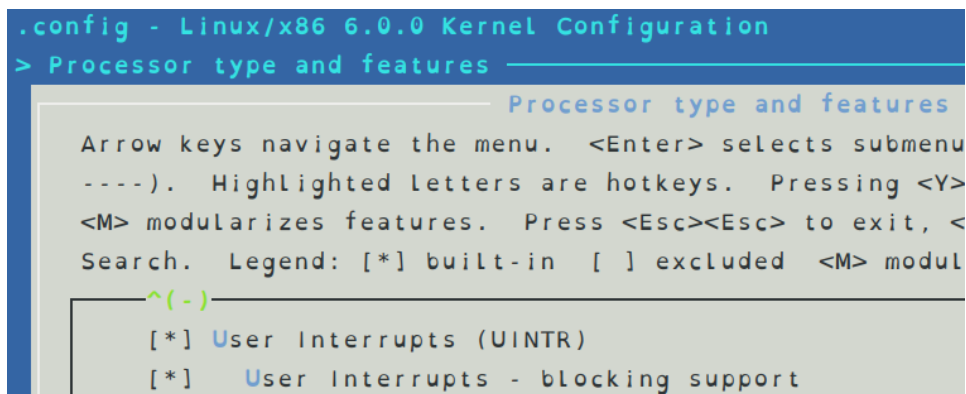


FIGURE 8 – Activer le support des *Uintr* à la compilation du noyau

```

1  #include <x86gprintrin.h>
2
3  #define VECTOR 6
4  bool isOver = 0;
5
6  __attribute__((target("general-regs-only")))
7  __attribute__((interrupt))
8  void ui_handler(struct __uintr_frame* ui_frame, unsigned long
   vector) {
9      /* Second clock_gettime() here */
10     printf("handler invoked (vector: %lld)\n", vector);
11
12     isOver = 1;
13 }
14
15 void receiver() {
16     if (uintr_register_handler(ui_handler, 0)) {
17         perror("interrupt handler register error");
18         exit(1);
19     }
20
21     _stui();
22
23     int uvec_fd = uintr_vector_fd(VECTOR, 0);
24     if (uvec_fd < 0) {
25         perror("vector fd error");
26         exit(1);
27     }
28
29     send_FD_to_sender(uvec_fd);
30
31     while (!isOver) continue;
32
33     if (uintr_unregister_handler(0))
34         perror("interrupt handler unregister error");
35     close(uvec_fd);
36 }

```

Listing 10 – Code récepteur

```

1  #include <x86gprintrin.h>
2
3  void sender(pid_t targetPid) {
4     int uvec_fd = wait_to_receive_FD_from_receiver(targetPid);
5
6     long uipi_index = uintr_register_sender(uvec_fd, 0);
7     if (uipi_index < 0) {
8         perror("sender register error");
9         exit(1);
10    }
11
12    /* First clock_gettime() here */
13    _senduipi(uipi_index);
14

```

```

15     if (uintr_unregister_sender(uipi_index, 0))
16         perror("sender unregister error");
17     close(uvec_fd);
18 }

```

Listing 11 – Code émetteur

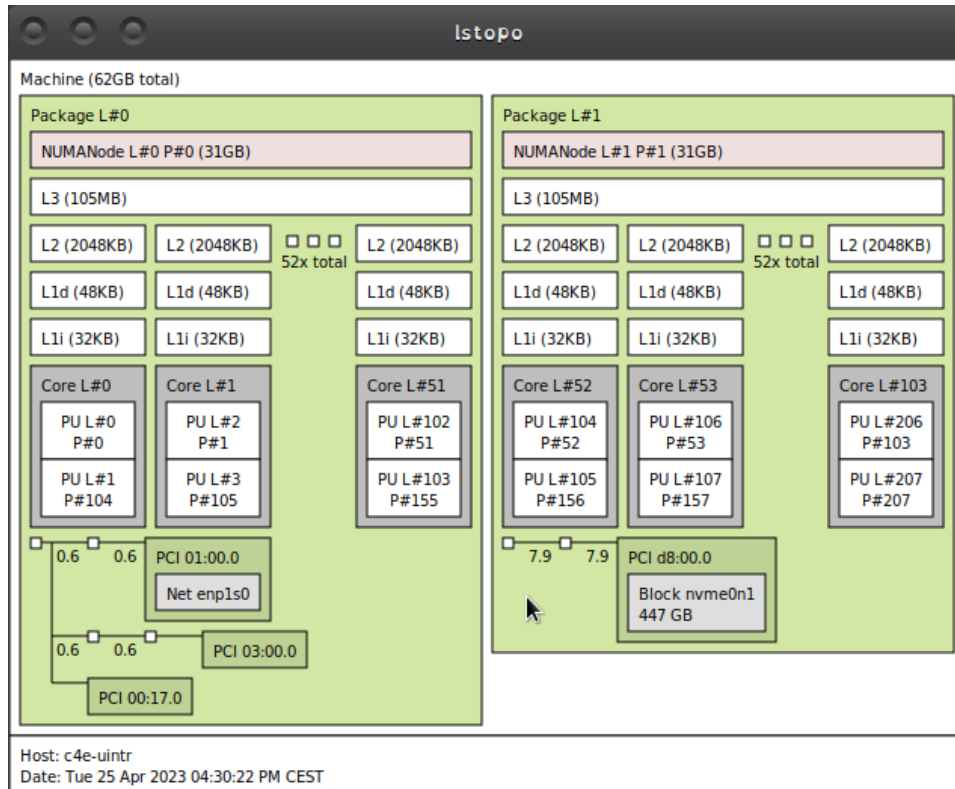
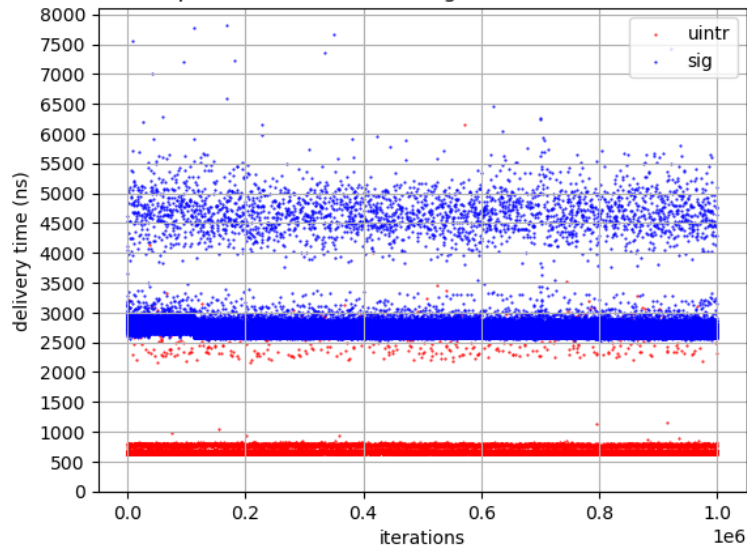


FIGURE 9 – Topologie de la machine fournie par *Atos* obtenue grâce à la commande *lstopo*

999999 iterations 2 process with core binding and without the first iteration (Next to)



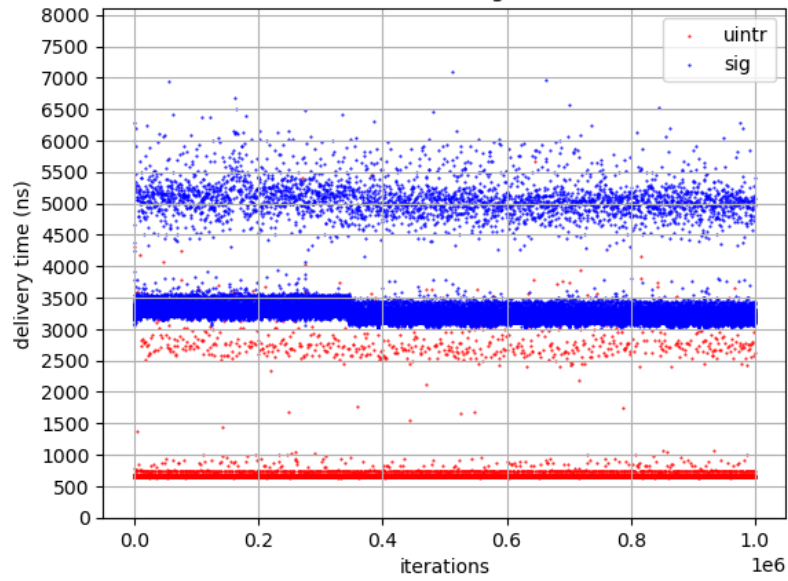
(a)

	mean	std	min	10%	50%	95%	max
sig	2675	147	2520	2620	2658	2781	44311
uintr	651	72	631	642	649	659	46053

(b)

FIGURE 10 – Mesures de latence entre deux processus avec un placement proche

999999 iterations 2 threads with core binding and without the first iteration (Far)



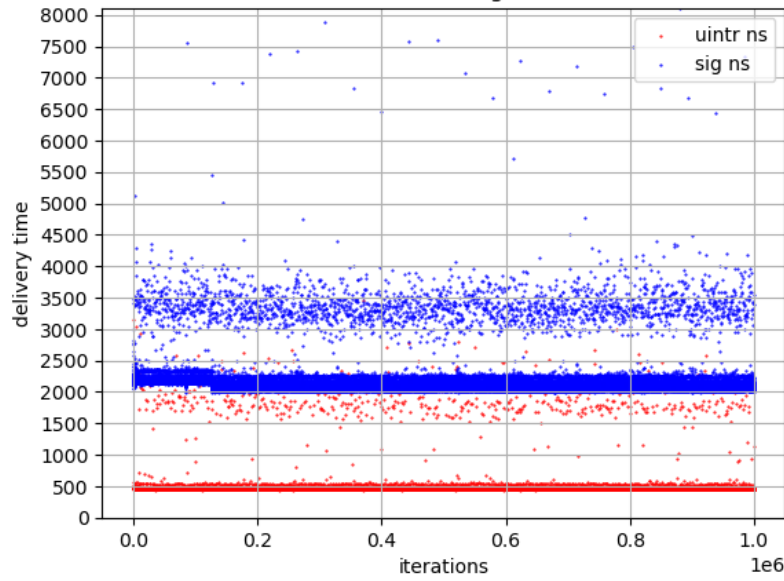
(a)

	mean	std	min	10%	50%	95%	max
sig	3210	156	3014	3141	3177	3302	63690
uintr	654	329	638	648	653	659	325408

(b)

FIGURE 11 – Mesures de latence entre deux threads avec un placement éloigné

999999 iterations 2 threads with core binding and without the first iteration (Far)



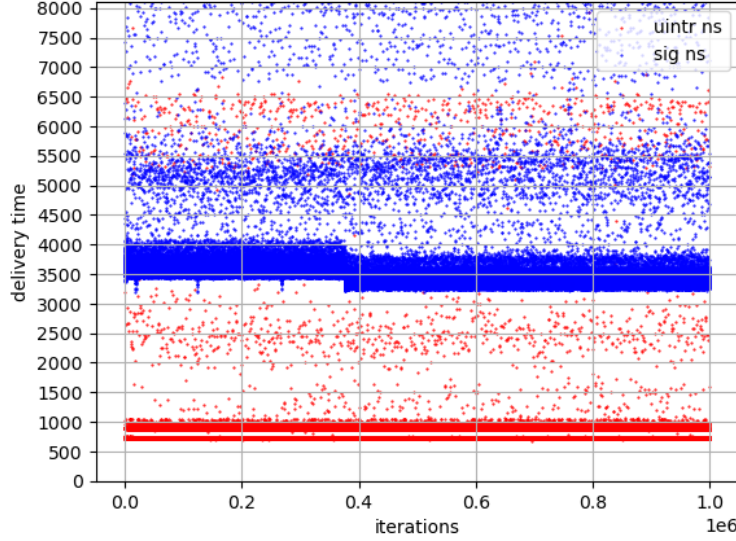
(a)

	mean	std	min	10%	50%	95%	max
sig	2097	282	1981	2064	2082	2178	270625
uintr	455	389	440	450	454	460	388815

(b)

FIGURE 12 – Mesures de latence entre deux threads avec un placement éloigné et le turbo boost activé

999999 iterations 2 threads with core binding and without the first iteration (Very Far)



(a)

	mean	std	min	10%	50%	95%	max
sig	3520	327	3189	3383	3481	3678	51798
uintr	736	159	683	721	725	735	46787

(b)

FIGURE 13 – Mesures de latence entre deux threads avec un placement très éloigné et le turbo boost activé

9 Bibliographie

- [1] Olivier Aumage, Elisabeth Brunet, Nathalie Furmento, and Raymond Namyst. NewMadeleine : a Fast Communication Scheduling Engine for High Performance Networks. In *Workshop on Communication Architecture for Clusters (CAC 2007), workshop held in conjunction with IPDPS 2007*, Long Beach, California, United States, March 2007.
- [2] Mathieu Barbe. Rapport de projet de fin d'études : Design and implement interrupt-handling code of a low-latency network adapter linux driver, 18 mars 2019 - 27 septembre 2019. Chez Bull (Atos). Grenoble INP - ENSIMAG.
- [3] Alexandre Denis. pioman : a pthread-based Multithreaded Communication Engine. In *Euromicro International Conference on Parallel, Distributed and Network-based Processing*, Turku, Finland, March 2015.
- [4] Alexandre Denis, Julien Jaeger, Emmanuel Jeannot, and Florian Reynier. One core dedicated to MPI nonblocking communication progression ? A model to assess whether it is worth it. In *International Sympos-*

sium on Cluster, Cloud and Internet Computing (CCGRID), Taormina, Italy, May 2022.

- [5] Alexandre Denis and François Trahay. MPI Overlap : Benchmark and Analysis. In *International Conference on Parallel Processing*, 45th International Conference on Parallel Processing, Philadelphia, United States, August 2016.
- [6] Saïd Derradji, Thibaut Palfer-Sollier, Jean-Pierre Panziera, Axel Poudes, and François Wellenreiter Atos. The bxi interconnect architecture. In *2015 IEEE 23rd Annual Symposium on High-Performance Interconnects*, pages 18–25, 2015.
- [7] Portable hardware locality (hwloc). <https://www.open-mpi.org/projects/hwloc/>.
- [8] Intel® 64 and ia-32 architectures, software developer’s manual. <https://www.intel.com/content/www/us/en/developer/articles/technical/intel-sdm.html>.
- [9] Intel linux kernel fork with uintr. <https://github.com/intel/uintr-linux-kernel/tree/uintr-next>.
- [10] Intel uintr linux man pages. <https://github.com/intel/uintr-linux-kernel/tree/uintr-next/tools/uintr/manpages>.
- [11] Sujet de stage. https://dept-info.labri.fr/~denis/Enseignement/Sujet_PFE_2023_uintr.html.
- [12] Signal linux man pages. <https://man7.org/linux/man-pages/man2/sigaction.2.html>.
- [13] Portals 4.0 specification. <https://www.sandia.gov/portals/portals-4-0-specification-clone-2/>.
- [14] Pull request for uintr_alt_stack() fix. <https://github.com/intel/uintr-linux-kernel/pull/2>.
- [15] Présentation du support des interruptions en espace utilisateur avec diapositive sur lwn.net. <https://lwn.net/Articles/869140/>.