

Contributions pour l'utilisation des machines parallèles, réparties et hétérogènes

à travers le prisme de la gestion des données partagées

Loïc Cudennec

DGA Maîtrise de l'Information,
BP 7, 35998 Rennes Cedex 9, France
loic.cudennec@intradef.gouv.fr



Carrière scientifique

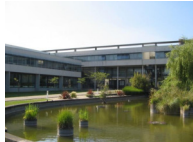
2005
Doctorant

Equipe-projet INRIA
PARIS/KerData

Santa Clara

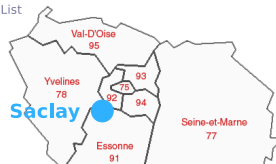


Tsukuba



2009
Ingénieur-Chercheur

CEA List



Saclay

2019
Chargé
d'expertise IA

DGA MI



Activités depuis le doctorat

Publications

- 2 journaux internationaux
- 1 chapitre de livre
- 20 conférences internationales
- 13 workshops internationaux
- 7 conférences francophones
- 1 brevet

Animation scientifique

- 13+ projets ANR, H2020 et industriels
- 35+ relectures (1 jury best-paper)
- 7 organisations de session ICCS (rang A)
- 5 encadrements de thèse (1 direction)
- 14 encadrements de master 2
- 1 organisation de Hackaton étudiant

2009



2012



2019



cea

énergie atomique - énergies alternatives

list

cea tech

MINISTÈRE
DES ARMÉES
*Liberté
Égalité
Patrimoine*

DGA
DIRECTION GÉNÉRALE
DE L'ARMEMENT

Thématiques de recherche

Système

- Cohérence des données
- Déploiement d'applications
- Intergiciels de calcul
- Compilateurs

Recherche opérationnelle

- Optimisation et modèles analytiques
- Maîtrise de la consommation d'énergie
- Adéquation logiciel-matériel pour l'IA
- Apprentissage machine distribué

2005

Grilles de calcul
Calcul haute-performance (HPC)



2009

Processeurs massivement parallèles
Many-coeurs



2015

Calcul hétérogène
HPC embarqué (HPEC)



2019

HP(E)C pour
l'Intelligence Artificielle



Projets de recherche et R&D

- Grille expérimentale Grid'5000
- Accord cadre INRIA-Sun Microsystems
- ANR GDS, LEGO et RESPIRE
- Sakura DISCUSS

- Laboratoire commun CEA-Kalray
- Laboratoire commun CEA-Synopsys

- H2020 M2DC, ExaNode et LEXIS
- EIT Digital Usine du futur

2005

Grilles de calcul
Calcul haute-performance (HPC)



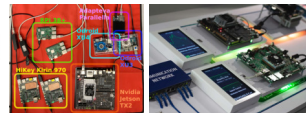
2009

Processeurs massivement parallèles
Many-coeurs



2015

Calcul hétérogène
HPC embarqué (HPEC)



2019

HP(É)C pour
l'Intelligence Artificielle



Encadrement doctorants



Cédric Gernigon

2020-

Suivi de thèse

Dir : Olivier Sentieys (INRIA)



Jussara Marândola Kofuji

2009-2011

Suivi de thèse

Dir : Marcelo Knörich Zuffo (USP, Brésil)



Erwan Lenormand

2018-2022

Direction

Codir : Henri-Pierre Charles (CEA)



Safae Dahmani

2012-2015

Encadrement

Dir : Guy Gogniat (UBS)



Cédric Prigent

2021-

Encadrement

Dir : Gabriel Antoniu (INRIA)

2005

Grilles de calcul

Calcul haute-performance (HPC)



L. Cudennec (DGA)

2009

Processeurs massivement parallèles

Many-coeurs

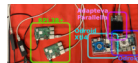


Soutenance d'HDR Informatique U.Rennes 1

2015

Calcul hétérogène

HPC embarqué (HPEC)



2019

HP(E)C pour

l'Intelligence Artificielle



28 septembre 2022

6 / 62

Encadrement doctorants et étudiants M2

Voichita Almasan
2006

Majd Ghareeb
2007

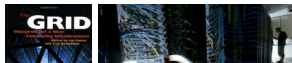
Abdullah Almousa Almaksour
2007

Marius Moldovan
2008

Andre Lage Freitas
2008

2005

Grilles de calcul
Calcul haute-performance (HPC)



L. Cudennec (DGA)

2009

Processeurs massivement parallèles
Many-coeurs



Soutenance d'HDR Informatique U.Rennes 1

Tien Thanh Nguyen
2015

Hayfa Alaya Ben Salem
2015



Jussara Marândola Kofuji
2009-2011

Suivi de thèse
Dir : Marcelo Knörich Zuffo (USP, Brésil)



Safae Dahmani
2012-2015

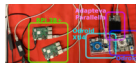
Encadrement
Dir : Guy Gogniat (UBS)

Hamza Chaker
2014

Cédric Maignan
2015

2015

Calcul hétérogène
HPC embarqué (HPEC)



Cédric Gernigon
2020-

Suivi de thèse
Dir : Olivier Sentieys (INRIA)



Erwan Lenormand
2018-2022

Direction
Codir : Henri-Pierre Charles (CEA)



Cédric Prigent
2021-

Encadrement
Dir : Gabriel Antoniu (INRIA)

Louis Syoën
2017

Mathis Valli
2022

Rihab Bennour
2018

2019

HP(É)C pour
l'Intelligence Artificielle



28 septembre 2022

7 / 62

ICCS/Alchemy : Architecture, Languages, Compilation and Hardware support for EMerging and heterogeneous sYstems

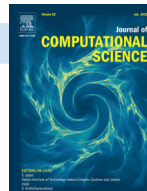
Sessions spéciales ICCS (rang A)

- Co-organisateur avec Stéphane Louise (CEA)
- 7 sessions de 2013 à 2019
- 35 papiers acceptés
- 30 à 40% de taux d'acceptation
- 30+ membres de comité de programme

Orateurs principaux

- 2013 Benoît Dupont de Dinechin (CTO Kalray)
- 2015 Raymond Namyst (Pr Université de Bordeaux)
- 2016 Andrea Olofsson (CEO Adapteva)

2013 2014 2015 2016 2017 2018 2019
 Barcelone Cairns Reykjavik San Diego Zürich Wuxi Faro



Évolution technologique dirigée par la microélectronique

Technologie

- Miniaturisation significative des composants électroniques
- Systèmes de stockage de l'énergie plus performants

Opportunités

- Puissance de calcul nomade
- Liens de communication efficaces
- Autonomie des équipements
- Vers l'**informatique ambiante**

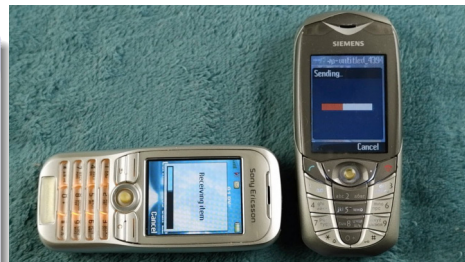


["boomer", 0.99]

Conséquences sur la programmabilité

Sauvegarder une photo le 28 septembre 2004

- Prendre une photo
- Aligner les capteurs infrarouges
- Démarrer la réception
- Sélectionner et envoyer la photo
- Attendre et ne pas désaligner les capteurs !



Sauvegarder une photo le 28 septembre 2022

- Prendre une photo



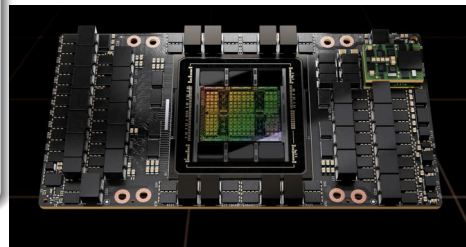
Conséquences sur la programmabilité

Effectuer un calcul le 28 septembre 1992

- Calculer

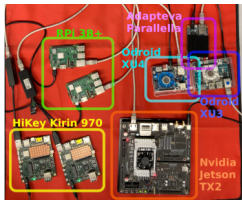
Effectuer un calcul le 28 septembre 2022
(avec un soupçon de mauvaise foi)

- Initialiser l'interface avec le GPU
- Allouer la mémoire distante
- Transférer les données entre le CPU et le GPU
- Lancer le noyau de calcul
- Transférer les résultats entre le GPU et le CPU
- Désallouer la mémoire



L'évolution technologique complexifie (souvent) le modèle de programmation

Des architectures et des applications : des enjeux actuels



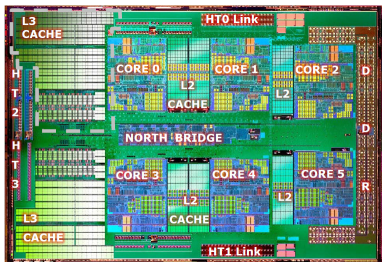
Distribution
des calculs

Parallélisme
d'exécution

Hétérogénéité
matérielle



Organisation en 3 thèmes



Cohérence des caches
 Protocoles de cohérence
 pour many-cœurs



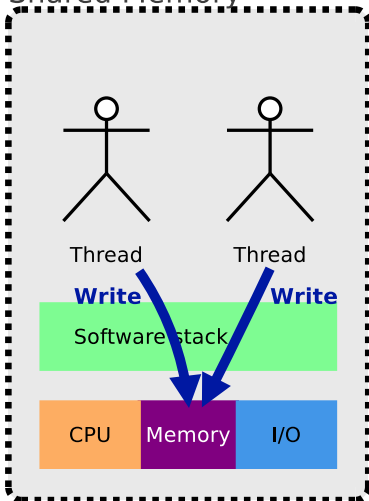
Unification des mémoires
 Mémoire partagée sur
 architectures hétérogènes



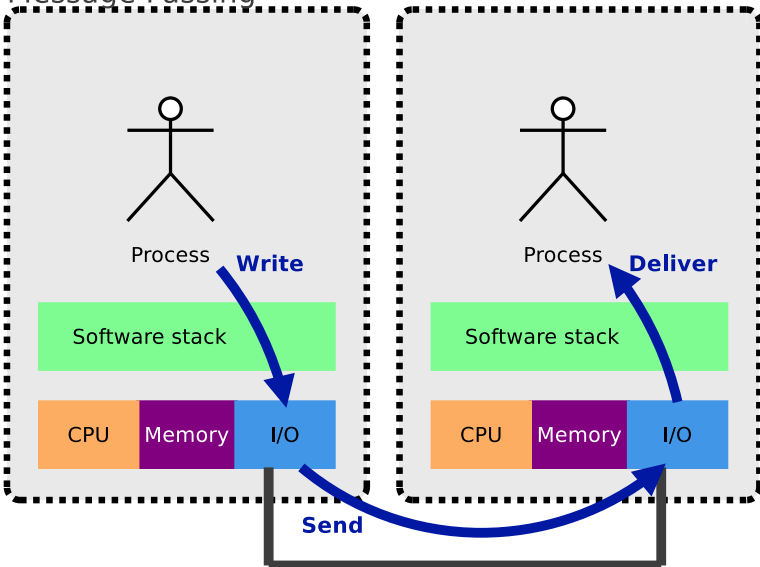
Aide à la décision
 Concilier performance et
 maîtrise de l'énergie

Modèles à mémoire partagée et mémoire distribuée

Shared Memory

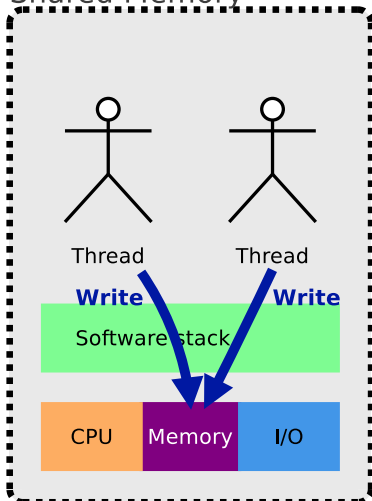


Message Passing



Modèles à mémoire partagée : cohérence des données

Shared Memory



Thread 1

```
a = 3
b = 4
print(a)
```

a = ?

Thread 2

```
a = 7
b = a
print(b)
```

b = ?

- Modèle de cohérence : contrat entre l'utilisateur et le système pour déterminer un ordre partiel sur les évènements
- Protocole de cohérence : protocole de communication implémentant un modèle de cohérence

Augmentation du nombre de cœurs de calcul : exemple du TOP500

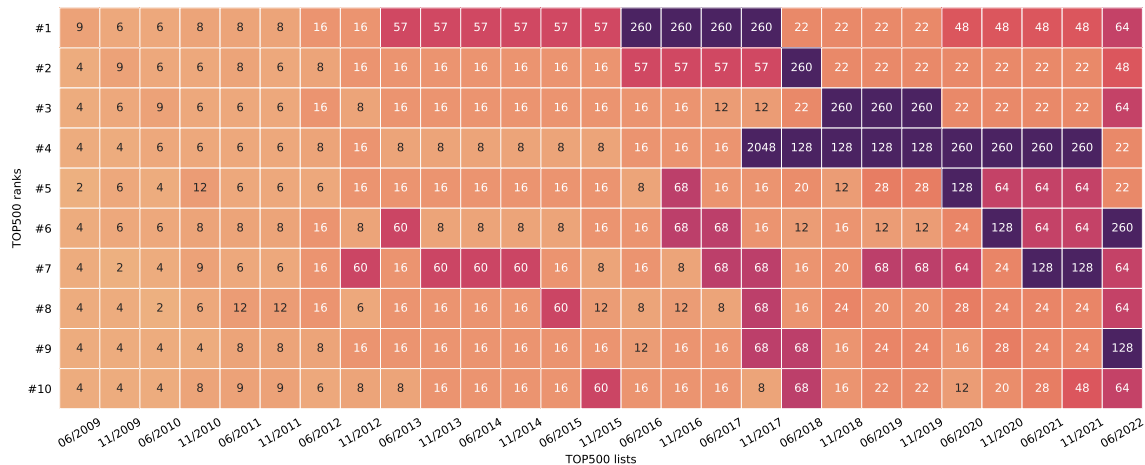


FIGURE – Nombre de cœurs embarqués par le processeur généraliste (multi-cœur ou many-cœur) le plus important présent sur les nœuds de calcul des systèmes classés dans les 10 premières places du TOP500, de juin 2009 à juin 2022.

Kalray, un many-cœur souverain dans la compétition internationale

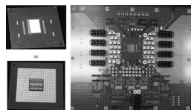


FIGURE – Intel 80-core (2007) [US 80 cœurs 97W]

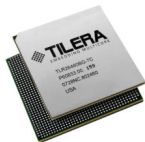


FIGURE – Tilera Tile 64 (2007) [US 64 cœurs 10W]



FIGURE – Intel Xeon Phi (2012) [US 61 cœurs 300W]



FIGURE – PEZY-1 (2012) [JP 512 cœurs 35W]



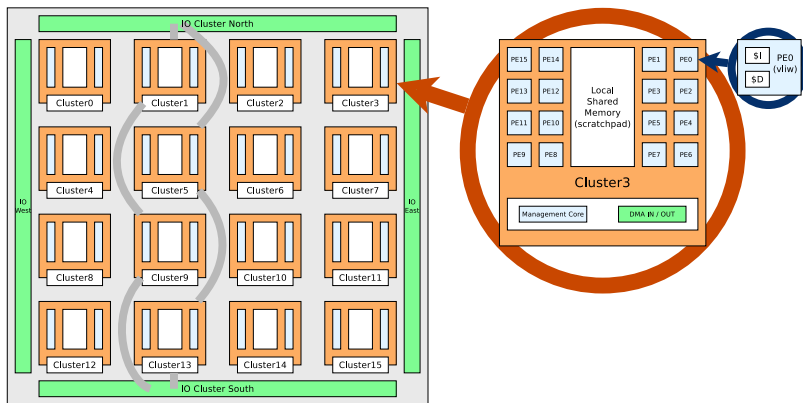
FIGURE – Kalray MPPA-256 (2013) [FR 256 cœurs 12W]

- Laboratoire commun CEA-Kalray, ~10 pax sur 3 années
- Objectif du CEA : travailler sur la programmabilité de la puce
- Proposition d'un langage flot-de-données (Σ -C)
- Livraison d'une chaîne de compilation dans le kit de développement du MPPA
- Performances ~à un Intel i7 et consommation 6 à 20 fois moindre
- Contribution à 6 publications internationales



FIGURE – MPPA AccessCore SDK

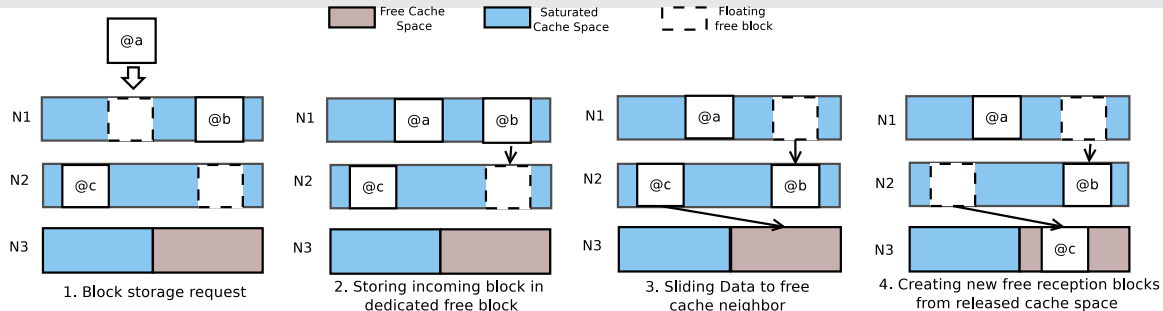
Exemple d'architecture many-cœur : le MPPA 256 de Kalray



Le many-cœur : une grille de calcul sur quelques cm^2

- Peut-on tirer partie de l'accumulation des mémoires locales pour concevoir un protocole de cohérence adapté aux many-cœurs ?

Protocoles de cohérence des caches coopératifs



Glissement de données

- Solliciter un cache voisin pour éviter une éviction dans la mémoire externe
- Maintenir un espace libre flottant pour accélérer les déplacements

HIPS@IPDPS 2013

Introducing a Data Sliding Mechanism for Cooperative Caching in Manycore Architectures. Safae Dahmani, Loïc Cudennec and Guy Gogniat.

Demandes d'accès réussies pour une donnée sur la puce

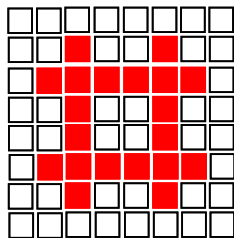


FIGURE – Motif d'accès sur une puce composée de 8x8 cœurs

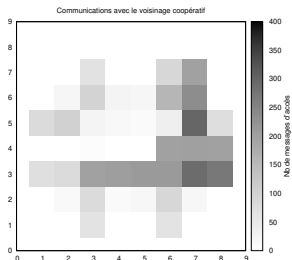


FIGURE – MESI + Cache Elastic [Herrero et al., 2010]

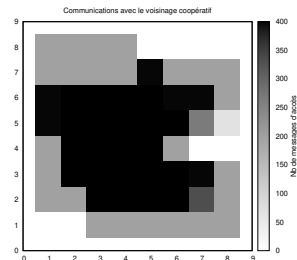
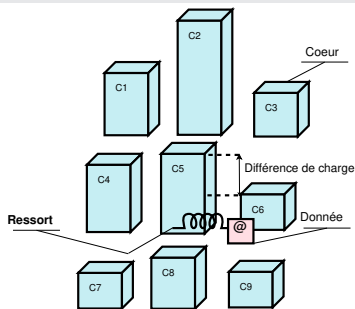


FIGURE – MESI + Cache Glissement [Dahmani et al, 2013]

Comparaison avec le cache Elastic

- Étalement plus important de la surface de coopération
- Conserve les données sur la puce (accès cache distant vs mémoire externe)
- Limitation : éloignement des données sur la puce par rapport au cœur d'origine

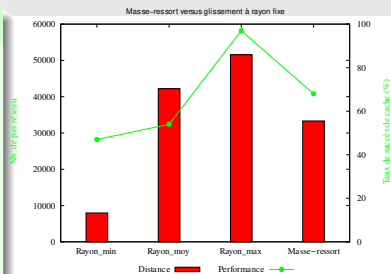
Caches coopératifs avec modèle masse-ressort



Formule discrète

$$D_i = 2AD_{i-1} - A^2D_{i-2} + B\Delta h$$

A et B constantes de raideur du ressort, de viscosité et de masse. Δh la différence de potentiel entre cœurs



Glissement de données avec masse-ressort

- Augmenter le taux de succès de cache et diminuer le nombre de hops
- Obtention d'un meilleur compromis que les approches statiques

MCSoc 2014

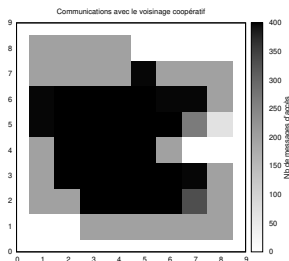
Using the Spring Physical Model to Extend a Cooperative Caching Protocol for Many-Core Processors. Safae Dahmani, Loïc Cudennec, Stéphane Louise and Guy Gogniat.

Comment évaluer un protocole de cohérence des données ?

Comment évaluer un protocole de cohérence des données ?

1. Méthode analytique

- Instrumenter un binaire parallèle
- Journaliser l'entrelacement des accès mémoire
- Rejouer en simulant le nouveau protocole



SOC 2012

Enhancing Cache Coherent Architecture With Access Patterns for Embedded Manycore Systems. Jussara Marandola, Stéphane Louise, Loïc Cudennec, Jean-Thomas Acquaviva and David Bader.

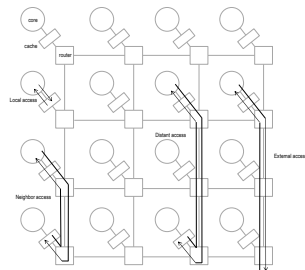
Comment évaluer un protocole de cohérence des données ?

1. Méthode analytique

- Instrumenter un binaire parallèle
- Journaliser l'entrelacement des accès mémoire
- Rejouer en simulant le nouveau protocole

2. Cycles processeur

- Appliquer la méthode 1
- Décomposer chaque type d'accès mémoire (local, distant, externe)
- Déduire le nombre de cycles processeur



COSMIC@CGO 2015

Cycle-based Model to Evaluate Consistency Protocols within a Multi-protocol Compilation Toolchain. Hamza Chaker, Safae Dahmani, Loïc Cudennec, Guy Gogniat and Martha Johanna Sepúlveda.

Comment évaluer un protocole de cohérence des données ?

1. Méthode analytique

- Instrumenter un binaire parallèle
- Journaliser l'entrelacement des accès mémoire
- Rejouer en simulant le nouveau protocole

2. Cycles processeur

- Appliquer la méthode 1
- Décomposer chaque type d'accès mémoire (local, distant, externe)
- Déduire le nombre de cycles processeur

3. Simulateur de NoC

- Appliquer la méthode 2
- Regrouper les accès mémoire concurrents par lots
- Injecter les lots dans le NoC pour créer de la contention réseau

MCSoc 2016

Network Contention-aware Method to Evaluate Data Coherency Protocols within a Compilation Toolchain. Loïc Cudennec, Safae Dahmani, Guy Gogniat, Cédric Maignan and Martha Johanna Sepúlveda.

Comment évaluer un protocole de cohérence des données ?

1. Méthode analytique

- Instrumenter un binaire parallèle
- Journaliser l'entrelacement des accès mémoire
- Rejouer en simulant le nouveau protocole

Objectif

- Évaluer le nombre de messages
- Déterminer le taux d'utilisation des caches

2. Cycles processeur

- Appliquer la méthode 1
- Décomposer chaque type d'accès mémoire (local, distant, externe)
- Déduire le nombre de cycles processeur

Objectif

- Évaluer le coût des accès
- Prendre en compte l'architecture

3. Simulateur de NoC

- Appliquer la méthode 2
- Regrouper les accès mémoire concurrents par lots
- Injecter les lots dans le simulateur de NoC

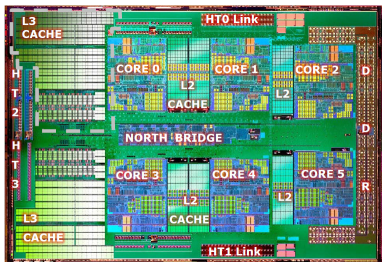
Objectif

- Mesurer la contention réseau
- Faire varier la fréquence temporelle

Évaluer un protocole de cohérence des données reste un problème difficile

Compromis entre précision et vitesse de simulation

Organisation en 3 thèmes



Cohérence des caches
 Protocoles de cohérence
 pour many-cœurs

Unification des mémoires
 Mémoire partagée sur
 architectures hétérogènes

Aide à la décision
 Concilier performance et
 maîtrise de l'énergie

Architectures distribuées hétérogènes : modularité et embarquabilité

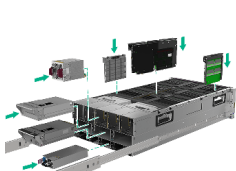


FIGURE – HPE Moonshot

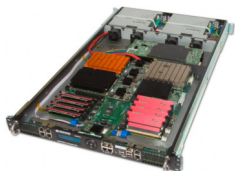


FIGURE – Christmann RECS



FIGURE – Démonstrateur FACE CEA-Renault

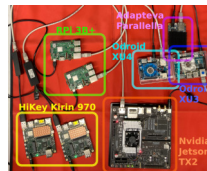


FIGURE – Bac-à-sable cartes embarquées



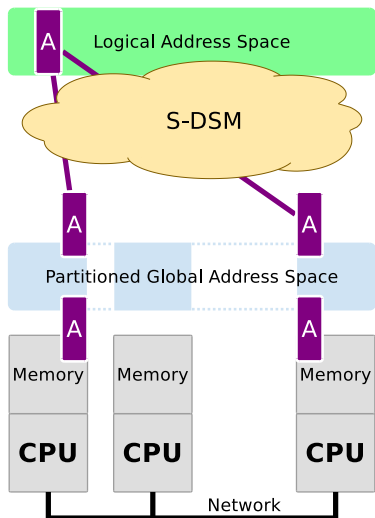
FIGURE – Isybot cobot Industrie 4.0

Caractéristiques communes

- Système distribué et parallèle
- Hétérogénéité des unités de calcul
- Mémoires réparties
- Communications peu performantes

Assemblage de composants sur étagère

Mémoire virtuellement partagée et machines hétérogènes



Mémoire virtuellement partagée (MVP)

- Espace d'adressage logique commun
- Accès concurrents aux données
- Localisation et transfert transparents pour l'utilisateur
- Abstraction des mémoires physiques et des moyens de communication

Positionnement

Faciliter la programmation des machines hétérogènes à l'aide d'une MVP

Comparaison MVP et passage de message (MP)

Référence	Système MP	Système MVP	#	Performance de la MVP
[Lu et al., 1995]	PVM	TreadMarks	8	15% of PVM speedup to 5% slower
[Scales and Gharachorloo, 1997]	Oracle database*	Shasta	2	2-4 times slower
[Bader and JáJá, 1999]	MPICH	CVM	8	10 times slower
[Werstein et al., 2003]	PVM, MPI	TreadMarks	32	slower, except for Mandelbrot kernel
[Antoniu et al., 2007]	GridRPC*	JuxMem	53	3 times faster
[Gelado et al., 2010]	CUDA*	GMAC ADSM	2	match
[Dimakopoulos and Hadjidoukas, 2011]	MPI	MOME, MOCHA	16	6-8 times slower
[Lin et al., 2012]	SysLink	Reflex	3	83% less energy, 2.5% slower
[Nelson et al., 2015]	MapReduce*, GraphLab*	Grappa	128	1.3 times faster
[Kaxiras et al., 2015]	MPI	Argo	128	match or exceed
[Beri et al., 2017]	StarPu*, SUMMA*	Unicorn	10	performs quite close

Tendances

- Les MVP deviennent aussi performantes que les systèmes par passage de message
- Les performances réseau autorisent des motifs de communication plus complexes
- Évolution favorable pour les systèmes de stockage basés sur les caches distribués

Conception et implémentation d'une MVP

MVP SAT - Share Among Things (CEA)

- Démonstrateur pour projets Européen (M2DC, LEXIS) à ISC et Teratec
- 14 publications de 2017 à 2021
- Support technique pour la thèse d'Erwan Lenormand
- ANSI C, 10k+ LoC, 1,5 h.an



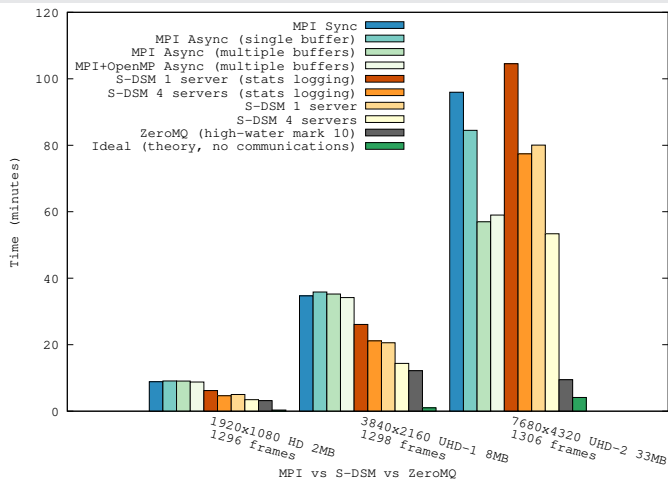
Thématiques de recherche

- Modèles de programmation (événementiel)
- Sécurité (chiffrement par attributs)
- Hétérogénéité des clients (FPGA)
- Optimisation du déploiement et de la consommation d'énergie

HeteroPar@Euro-Par 2017

Software-Distributed Shared Memory over Heterogeneous Micro-Server Architecture. Loïc Cudennec.

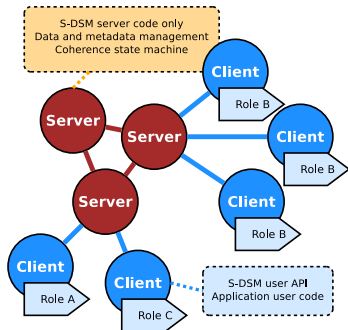
Comparaison MVP et passage de message (MP)



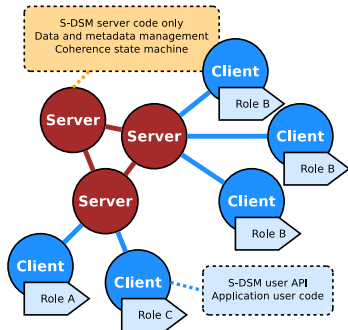
ParaMo@Euro-Par 2020

Experiments using a Software-Distributed Shared Memory, MPI and 0MQ over Heterogeneous Computing Resources. Loïc Cudennec and Kods Trabelsi.

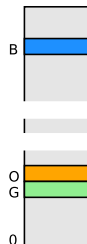
Quelques choix de conception de SAT



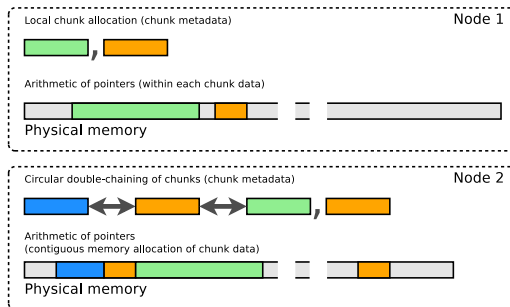
Quelques choix de conception de SAT



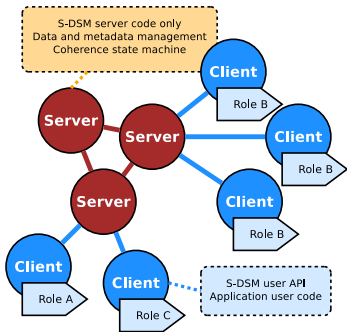
4294967295



S-DSM logical
address space
(unsigned long)



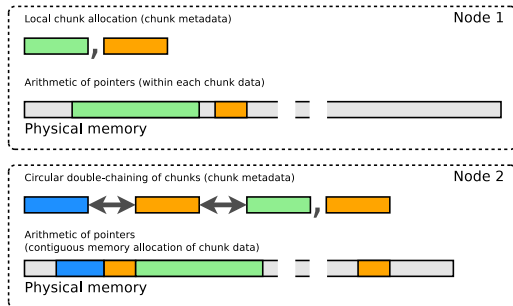
Quelques choix de conception de SAT



4294967295



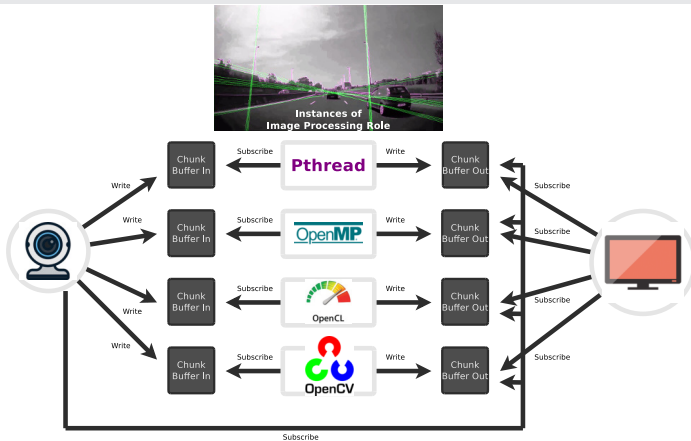
S-DSM logical address space (unsigned long)



Allocation et accès dans SAT

```
Consistency_t * consistency = newHomeBaseMESI() ;
Chunk_t * chunk = MALLOC(consistency, chunkid, N * sizeof(unsigned int));
unsigned int i = 0;
WRITE(chunk)
for(i = 0; i < N; i++) {
    chunk->data[i] = i;
}
RELEASE(chunk)
```

Modèle de programmation évènementiel



Principe

Un chunk est une zone mémoire partagée sur laquelle il est possible de souscrire aux notifications de modification

Apports du modèle évènementiel

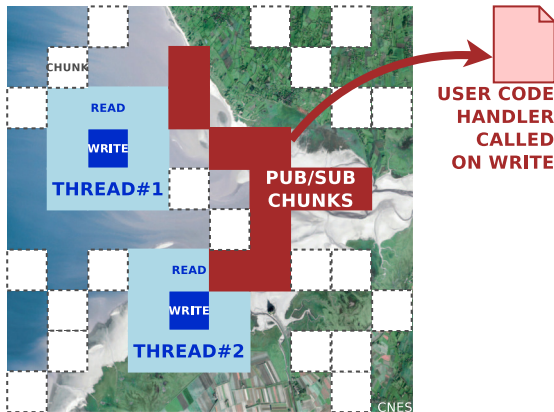
- Ordonnancement FIFO par construction
- Implémente un programme flot-de-données

HeteroPar@Euro-Par 2018

Merging the Publish-Subscribe Pattern with the Shared Memory Paradigm. Loïc Cudennec.

Modèle de programmation évènementiel

SHARED DATA



Principe

Un chunk est une zone mémoire partagée sur laquelle il est possible de souscrire aux notifications de modification

Apports du modèle évènementiel

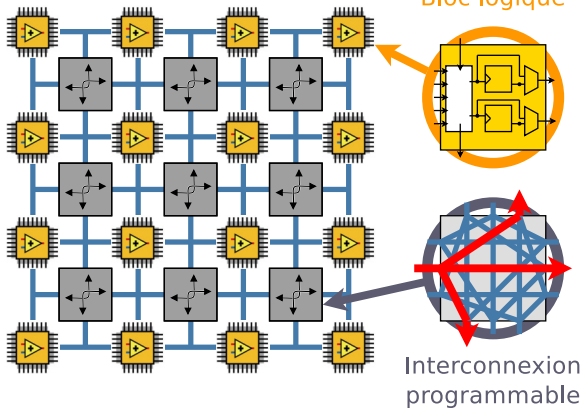
- Processus de surveillance non-intrusif
- Couplage de code HPC et évènementiel transparent

HeteroPar@Euro-Par 2018

Merging the Publish-Subscribe Pattern with the Shared Memory Paradigm. Loïc Cudennec.

Intégrer des FPGA dans la mémoire partagée

Matrice FPGA

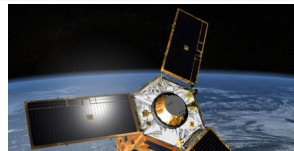
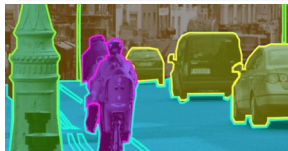
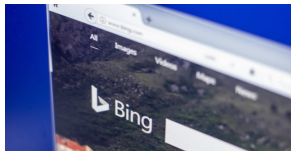


Processeurs à la logique programmable

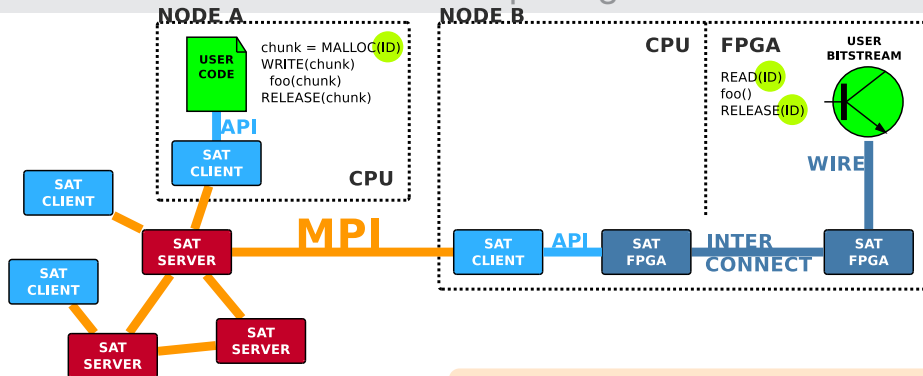
- Exprimer un algorithme sous forme de circuit logique
- Performance calculatoire (temps de propagation des signaux)
- Performance énergétique (/2 GPU, /40-400 CPU)

Problématiques

- Intégration difficile dans un logiciel
- Modèle de calcul flux-de-données



Intégrer des FPGA dans la mémoire partagée



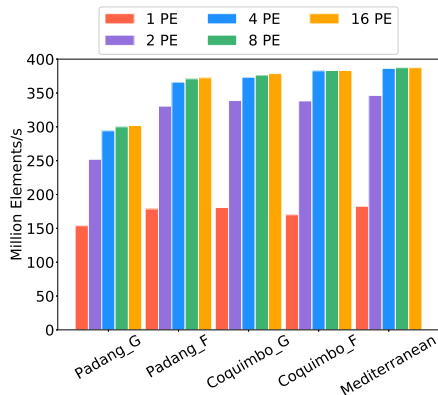
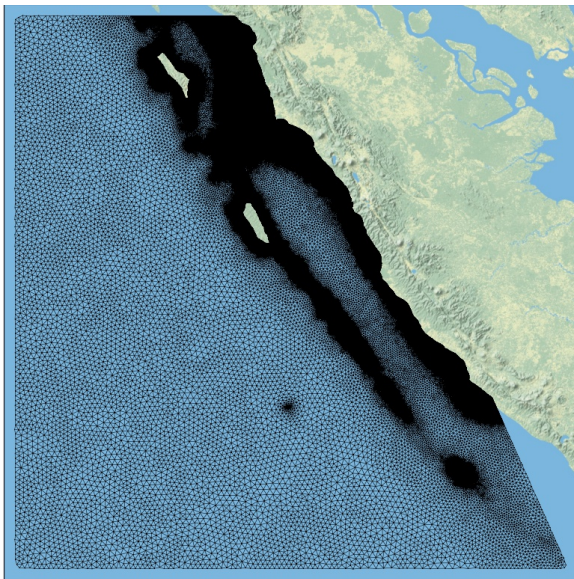
Principe

- Accéder à l'espace partagé depuis le circuit logique utilisateur
- Rendre le FPGA capable d'initiative pour accéder aux données

RSP 2020

A combined fast/cycle accurate simulation tool for reconfigurable accelerator evaluation : application to distributed data management. Erwan Lenormand, Thierry Goubier, Loïc Cudennec and Henri-Pierre Charles.

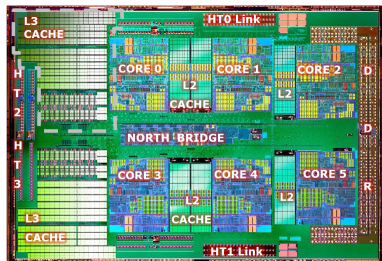
Intégrer des FPGA dans la mémoire partagée



HeteroPar@Euro-Par 2021

Data management model to program irregular compute kernels on FPGA : application to heterogeneous distributed system. Erwan Lenormand, Thierry Goubier, Loïc Cudennec and Henri-Pierre Charles.

Organisation en 3 thèmes

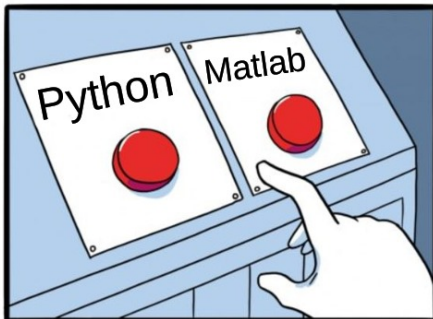


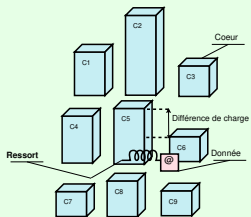
Cohérence des caches
Protocoles de cohérence
pour many-cœurs

Unification des mémoires
Mémoire partagée sur
architectures hétérogènes

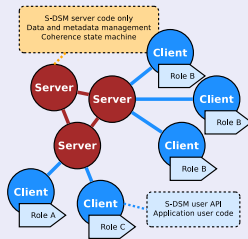
Aide à la décision
Concilier performance et
maîtrise de l'énergie

Fournir des outils considérés performants n'est pas suffisant

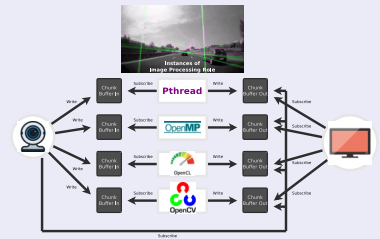




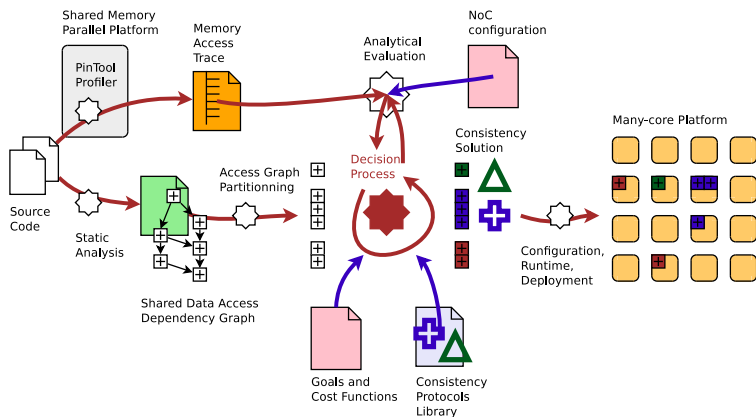
Choix du protocole de cohérence de cache et paramétrage



Choix de la topologie applicative et projection sur les ressources de calcul



Optimisation de la consommation énergétique des communications



Principe

- Associer à chaque donnée partagée un protocole de cohérence
- Déterminer les paramètres pour chaque accès
- Chercher à satisfaire des objectifs de performance et de sobriété énergétique

ROADEF 2015

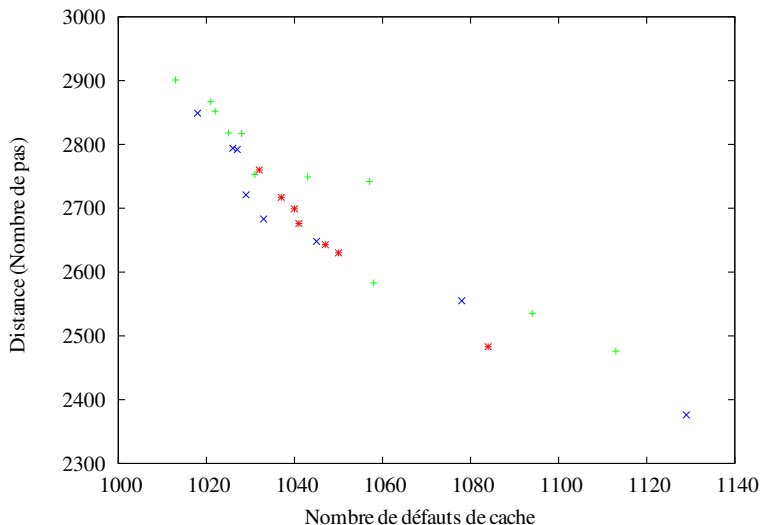
Aide à la décision pour le choix et le paramétrage de protocoles de cohérence des données. Safae Dahmani, Loïc Cudennec et Guy Gogniat.

Brique de décision

- FastPGA [Eskandari et al., 2006] *Pareto Genetic Algorithm*
- Multi-objectifs : nombre de défauts de cache et distance sur la puce

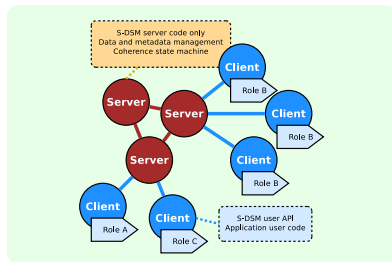
Expérimentation

- Trace de 2080 accès
- Variation du facteur de diversité de FastPGA

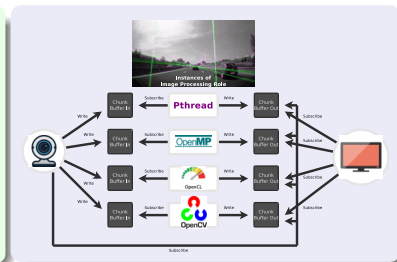


EMO@MCDM 2015

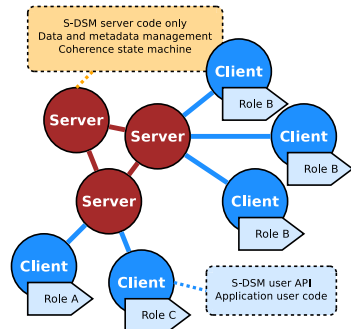
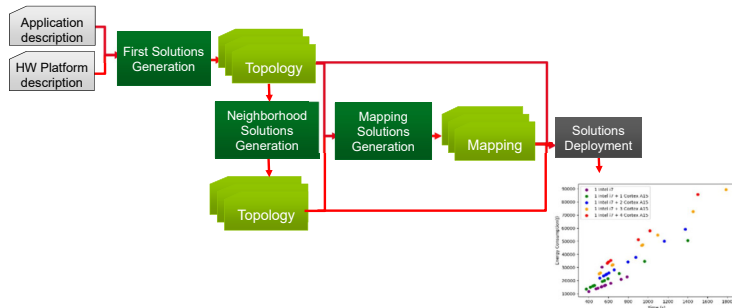
A Multiobjective Consistency Decision Engine for a Manycore Compilation Platform using an Evolutionary Approach. Safae Dahmani, Sergiu Carpov, Loïc Cudennec and Guy Gogniat.



Choix de la topologie
applicative et projection sur
les ressources de calcul



Optimisation de la
consommation énergétique
des communications



Principe

- Déterminer la topologie de l'application (# serveurs, clients, interconnexions)
- Déterminer une projection sur les ressources de calcul
- Chercher à satisfaire des objectifs de performance et de sobriété énergétique

HeteroPar@Euro-Par 2019

Application Topology Definition and Tasks Mapping for Efficient Use of Heterogeneous Resources. Kods Trabelsi, Loïc Cudennec and Rihab Bennour.

Brique de décision

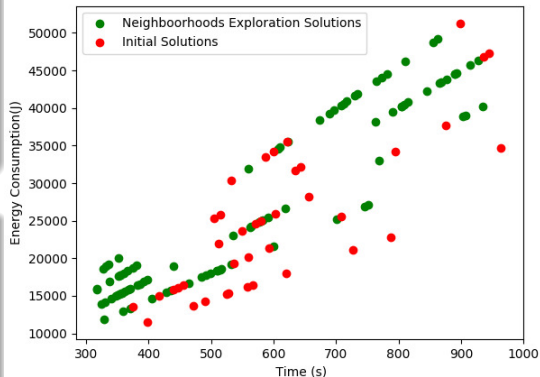
- RVV [Hansen and Mladenović 1999] : Recherche à voisinage variable
- Multi-objectifs : temps d'exécution et consommation d'énergie

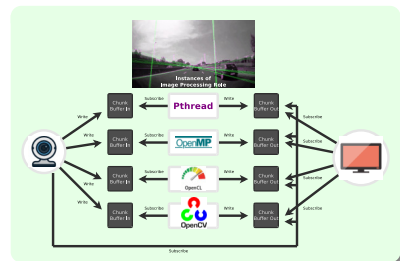
Expérimentation

- Matériel : Christmann RECS, 2 Intel Core i7, 8 ARM Cortex A15
- Application : traitement de flux vidéo sur mémoire partagée

Résultat

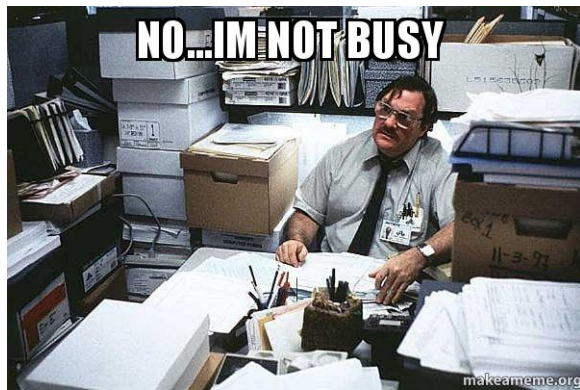
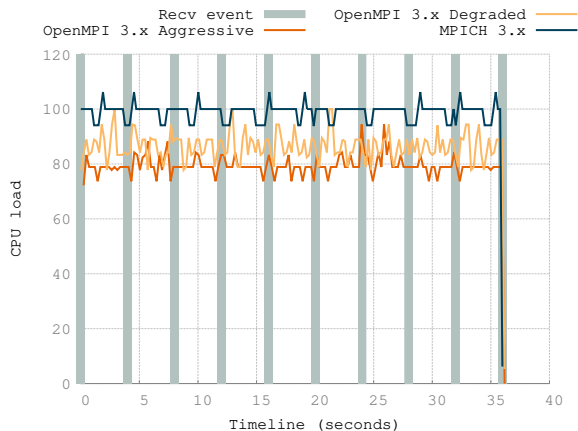
- Solutions initiales : solutions métier et aléatoires
- Exploration : meilleures solutions souvent contre-intuitives





Optimisation de la
consommation énergétique
des communications

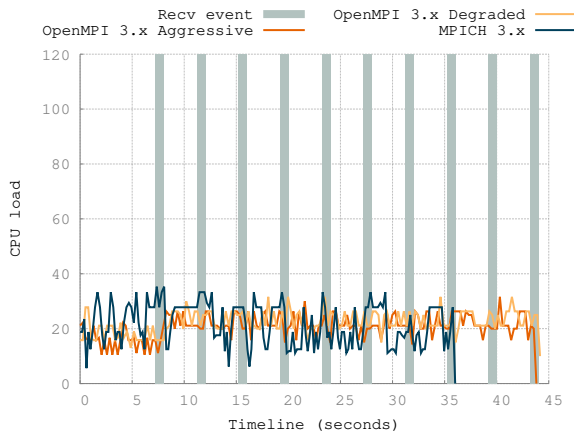
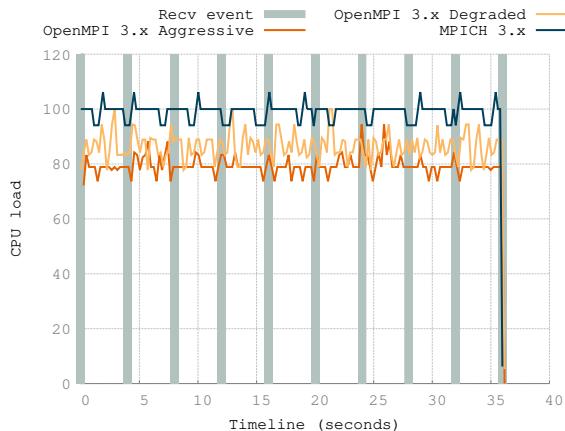
MPI et les enjeux du HPC embarqué (HPEC)



Comportement d'un processus en attente de message

- Attente active adaptée pour la réactivité des applications HPC
- Utilisation déraisonnée du processeur → consommation d'énergie

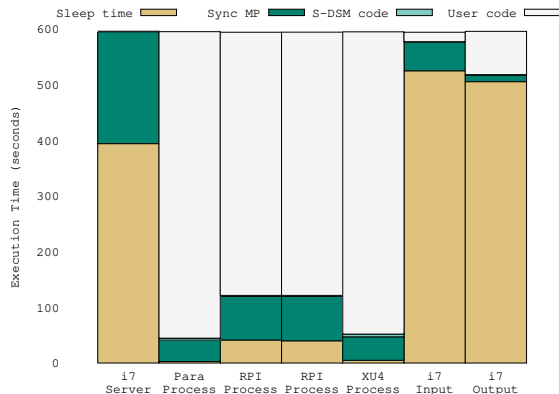
MPI et les enjeux du HPC embarqué (HPEC)



Mécanisme de micro-sieste paramétrable

- Interrompt succinctement l'attente active pour redonner la main à l'OS
- Implémentation transparente en espace utilisateur entre l'application et MPI

MPI et les enjeux du HPC embarqué (HPEC)



	TIME	FPS	W	kJ	FPS/W
	12'52	2.23	60	44.39	0.037
	12'40	2.26	66	47.68	0.034
	12'13	2.35	52	36.98	0.045
	12'42	2.26	52	38.16	0.043

Résultats (micro-sieste en violet)

- Gain en économie d'énergie (20%)
- Gain en performance (5%) lorsque des processus sont co-localisés sur un même processeur

Concurr. Comput. Pract. Exp. 32(24) (2020)

Adaptive Message Passing Polling for Energy Efficiency : Application to Software-Distributed Shared Memory over Heterogeneous Computing Resources. Loïc Cudennec.

Synthèse

Panorama des travaux

Many-cœurs

Langage flot-de-données

- Langage et modèle de calcul
- Conception d'une chaîne de compilation
- Transfert industriel

Cohérence des caches

- Protocoles coopératifs
- Outils d'évaluation
- Outils pour le choix des protocoles

Machines hétérogènes

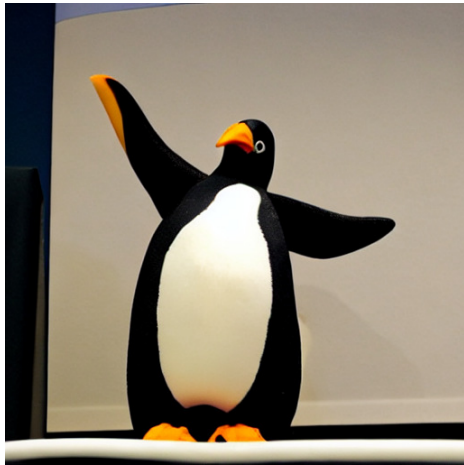
Mémoire partagée

- Conception d'une MVP
- Programmation événementielle
- Intégration des FPGA

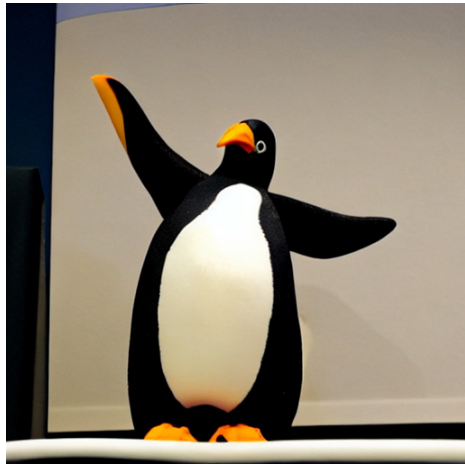
Performance

- Étude mémoire partagée vs passage de message
- Exploration des topologies applicatives
- Optimisation de l'énergie dans MPI

Perspectives



"penguin delivering a keynote speech"
Stable Diffusion [Rombach et al. 2022]



Emad
@EMostaque



Replying to [@KennethCassel](#)

We actually used 256 A100s for this per the model card, 150k hours in total so at market price \$600k

2:46 PM · Aug 28, 2022 · Twitter for iPhone

114 Retweets 65 Quote Tweets 1,039 Likes

Enjeux calculatoires en intelligence artificielle

Apprentissage

Centres de calcul

- Grands volumes de données
- Accumulation d'accélérateurs spécifiques
- Distribuer les calculs

Continuum informatique

- Apprentissage fédéré
- Migration des calculs entre cloud et edge
- Volatilité et reconfiguration

Inférence

Systemes embarqués

- Embarquabilité des modèles : compromis précision et empreinte mémoire
- Frugalité énergétique : vers l'internet des très petits objets
- Résilience des systèmes : le retour de l'informatique autonome

Contributions pour l'utilisation des machines parallèles, réparties et hétérogènes

à travers le prisme de la gestion des données partagées

Loïc Cudennec

DGA Maîtrise de l'Information,
BP 7, 35998 Rennes Cedex 9, France
loic.cudennec@intradef.gouv.fr

