


Historique.

Un peu d'histoire pour expliquer où nous en sommes.

En 1999, le 22 septembre, avec la thèse de François Galilée, la spécification d'Athapascan était posée: Fork + Shared + les modes et droits d'accès r, w, rw, cw ainsi que les modes postponed (rp,wp, rpwp, cwp). Une implémentation C++ complète a existé (sisi). Le coût à l'exécution était énorme : environ 10000 appels de fonction C pour 1 fork sans paramètre. Ceux-ci s'expliquaient de deux manières: 1/ un algorithme de terminaison distribuée était utilisé pour calcul la fin d'accès à une version pour déclencher les calculs sur la version suivante ; 2/ les choix des structures de données ainsi que leur allocation dans le tas ; 3/ le mauvais couplage entre la couche de communication thread safe Athapascan-0 et l'implémentation d'Athapascan-1 au dessus. Bien qu'en théorie, les algorithmes de cette implémentation distribuée devaient permettre d'avoir du bon speedup, en pratique il fallait avoir des calculs à très gros grain. L'ordonnancement était basé sur des variations autour du vol de travail avec des heuristiques pour permettre de contrôler l'utilisation mémoire, ou bien par ordonnancement de type statique d'un graphe de flot de données mais aucune perf n'a pu être observée.

Quelques semaines plus tard, le 14 décembre 1999, Mathias Doreille montrait dans sa thèse qu'il était possible d'avoir de très bon speedup sur des programmes ordonnancés par une variation d'un algorithme de type **ETF** et en utilisant une implémentation adhoc très légère directement au dessus de MPI mais incomplète vis-à-vis des spécifications d'Athapascan-1. L'ordonnancement local des messages et du calcul ne permettait pas d'avoir des performances raisonnables sur des programmes du type itératifs (jacobi, poisson) normaux (sans augmentation complexe / artificielle du grain). 

En 2004, Rémi Revire a montré la faisabilité d'avoir une implémentation plus efficace permettant à la fois un ordonnancement par vol de travail et par ordonnancement de graphe en se basant sur une implémentation qui, grâce à l'ordre d'exécution quasi séquentiel d'un programme Athapascan, permet de gérer le cycle de vie des tâches et objets partagés en rapport à leurs portées de déclaration : la gestion par pile de l'exécution des programmes séquentiels était retrouvée ! Néanmoins, l'algorithme de terminaison du calcul distribué (non plus de la fin d'utilisation d'une version) était encore mal conçu, mais heureusement n'intervenait que pour déterminer la terminaison globale de l'exécution des processus.

Du point de vu des performances, l'implémentation se basait sur l'utilisation de lock pour verrouiller l'accès à certaines données lors des opérations de vols. Le partitionnement statique du graphe utilisait Scotch et était à l'état de prototype instable. L'ensemble était basé sur une couche de communication appelée Inuktitut qui permettait de communiquer par envoi/réception de messages actifs.

Laurent Pigeon avait implanté un ensemble d'algorithmes pour la diffusion parallèle de message durant son master.

Lors d'un stage d'été, Xavier Besseron avait implanté le côté communication sur architecture hétérogène en utilisant soit un protocole de type Xdr (eXternal data representation), soit de type ASN.

Lors de son stage, Everton Hermann a effectué le portage d'Inuktitut sur Myrinet.

Le coût de cette implantation était d'environ 1000 appels de fonction C pour 1 fork

(toujours sans paramètre).

En 2005, Kaapi est née des cendres d'Athapascan & Inuktitut : une évaluation synthétique aurait été "ça marche pas trop mal, code de niveau prototype de recherche, mais peut mieux faire"

Plus précisément:

- peu de support du côté d'Inuktitut qui avait été conçu pour servir à la conception de Taktuk (1ère version hélas très bogguée en C++) ou des outils ka-XXX pour la diffusion de fichiers.
- 3 protocoles de communication dans Inuktitut: active message, write & signal (idem active message mais avec une allocation par la couche de comm des données reçues dans l'espace mémoire utilisateur) et une variante allocate & write & signal bien adaptée à un portage de la couche de comm directement au dessus de CORBA.
- la portée des structures de données dans Athapascan n'était pas encore bien comprise, leur gestion nécessitait encore l'utilisation de locks.
- l'algorithme de vol était décomposé en trois endroits : soit le vol effectif d'une tâche lors d'une requête local, soit lors vol suite à une requête à distance, soit le réveil d'un thread ce qui posait pas mal de problème de "stabilité" de celui-ci en cas de modification (...) En pratique toute tentative de modification d'une partie du vol de tâches ou du choix du reveil d'un thread était source de longues heures de debug...

Beaucoup de concepts avait été validés (couplage ordo statique / vol de travail) mais l'ensemble des sources étaient difficilement abordables par les étudiants arrivant sur le projet. En 2005, il fut donc fait les choix suivants:

- 1- récupération d'Inuktitut et suppression du code non essentiel: protocol message actif
- 2- lors du stage d'Everton Hermann il a été vu qu'un fonctionnement de la couche de message par vol de travail permettait simplement d'agréger des messages vers le même destinataire, de plus elle était original dans le sens ou la fenêtre d'agrégation dépendait de l'activité réseau...
- 3- une même opération pour le vol de tâche ou le réveil de thread: la fonction de vol prend en entrée un thread et retourne un thread ! Cela a l'avantage d'être simple et uniforme : dans le premier cas, un objet thread est volé et le thread résultat est rempli par la tâche volée, dans le second cas un objet processeur est volé et il retourne le thread à réveiller.
- 4- la notion de graphe de flot de données (application) doit être découpé du vol de travail (ordonnancement), l'ordonnancement de Kaapi se base donc sur un ensemble de threads qui peuvent être volés, les threads pouvaient être directement du code applicatif ou bien structurés par pile de tâches. En pratique, seul la deuxième possibilité a été utilisé, même à travers les algorithmes adaptatifs de Daouda. Cette structuration a permis de réduire l'utilisation de verrou dans l'implémentation.
- 5- une conception uniforme des objets pouvant être sérialisés et qui sont représentés par leur "format" : les data ou les threads possèdent un format. La couche API Athapascan génère automatique avec un ensemble d'expression template les objets formats en fonction du type des objets à sérialiser.
- 6- la capacité d'interagir avec l'exécution par l'envoi de signaux : SIGSTOP &

SIGKILL.

7- et enfin, un prototype pour la collection de statistique d'exécution des processus et l'enregistrement d'événement dans des fichiers de trace.

Les derniers points 5/ et 6/ ont facilité en partie l'implémentation des protocoles de tolérance aux pannes : arrêt des activités et sérialisation *de facto* des threads.

L'interface de programmation de la couche de communication n'a pas bougé depuis 2005, en revanche un développement a été nécessaire pour permettre de faire communiquer un ensemble de processus dynamiques et rendre robuste la panne d'un ou plusieurs processus sans provoquer l'arrêt les autres de manière non contrôlée. De plus, cette couche est "multi-réseau" avec une sélection du réseau à la source en fonction du destinataire qui a l'avantage d'être simple, mais n'est pas forcément très performante.

Cette version de Kaapi a permis de reporter successivement les 3 plugtest 2006 (non officiel dans le sens que nous n'avons pas été classé comme les autres car non java-compliant, mais les résultats étaient bien devant les autres), 2007 et 2008.

Pour 2009 : le ETSI Grid@Work plugtest s'oriente vers les aspects uniquement déploiement... on ne participera pas.

Mes critiques sont :

1- code gros devenant gros ne facilitant pas la gestion de l' "histoire" pour les nouveaux arrivants.

2- les perfs de l'aspect purement multi-core / SMP / many-core peuvent être meilleur

3- les choix des routes / des interfaces dans la couche communication doit être amélioré. N'ayant pas de travaux de recherche (thèse) sur ce point, nous n'arrivons pas à nous reposer sur une couche externe "portable" et robuste (grid, cluster, ...) ayant des capacités de dynamique / robustesse en cas de panne qui soient suffisantes.

4- l'absence d'une interface permettant de gérer la transparence référentielle des objets partagés, i.e. être capable d'accéder en lecture ou écriture aux objets quel que soit le site d'exécution. Ce point est délicat et doit être guidé par certaines décisions d'ordonnancement.

5- l'absence de la notion de collection d'objets distribués (utilisant un support comme décrit dans le point 4/ ci-dessus) permettant la description de la plupart des algorithmes en calcul scientifique pour lesquels il serait intéressant d'avoir à manipuler des tableaux multi-dimensionnels d'objets "shared".

6- absence d'une couche "serveurs de stockage fiable" pour les protocoles de tolérance aux pannes... De même que pour 3, il est difficile de trouver du source dans ce domaine qui permettent de bonne performance. Néanmoins une évaluation des outils récents seraient à refaire (la dernière date de 2006-2007).

6- mauvais couplage des algorithmes adaptatifs.

Début 2009: naissance de X-Kaapi pour répondre aux points 2/ et 6/ tout en facilitant le portage de X-Kaapi sur des systèmes embarqués dans le cadre de nos collaborations avec ST, i.e. en utilisant du C !

Objectifs X-Kaapi & Problématiques.

Les objectifs de X-Kaapi sont de capitaliser le savoir faire que nous avons eu en vu de cibler de nouvelles architectures (pour nous) et de nouvelles applications.

Le contexte technologique est induit par des architectures à mémoire hiérarchique composées par des processeurs de type multi-core ou many-core ou du domaine de l'embarqué (MpSoC). Ces processeurs communiquent par de la mémoire ou un réseau. La latence d'accès à la mémoire varie fortement en fonction de la distance entre la données et le cœur qui y accède. L'environnement d'exécution disponible n'est peut-être pas aussi complet que celui d'un PC actuel. En particulier en ce qui concerne l'ordonnancement de threads, la gestion mémoire (pas de mémoire virtuelle, taille limité). Le nombre total de cœur peut être très important.

Les deux architectures cibles initiales sont :

- une bluegene : un nœud est composé d'un ensemble de cœurs partageant une mémoire local, partageant ou non une hiérarchie de caches. Le réseau d'interconnexion est de type grille 3D torique (sous contrainte de routage au bord en fonction d'une réservation partielle). Ce réseau n'est utilisé que pour les communication bi-point. La machine dispose aussi d'un réseau de diffusion (broadcast) et d'un réseau de synchronisation. La taille mémoire par nœud est limité (256MBytes ou 512MBytes). Les possibilités de l'OS sont à découvrir. La couche de communication DCMF disponible est de bas niveau et offre une communication de type messages actifs et lecture/écriture à distance. Les questions qui se posent sont liés à la capacité à traiter des handlers de communication en concurrence avec du calcul, la capacité du système à faire progresser les communications et à ordonnancement plusieurs threads.
- MpSoC : la board ST composé de 4 cœurs ST 231 qui communiquent par des envois de message ou du partage d'objet (dans un sens très particulier) vers, et uniquement vers, le cœur principal ST40. Très peu de mémoire (quelques MBytes), un OS nouveau pour nous (OS21) mais qui possède la notion de task ~ thread et qui permet leur ordonnancement (préemptif ?) sur un cœur ST231. A noter que l'OS du ST40 est un Linux.

En pratique et pour commencer le développement, une architecture du type Idkoiff sera utilisé.

Les problématiques associées sont :

1- API pour le calcul

Dans Kaapi/Athapascan le modèle imposé est de type graphe de flot de données sur lequel nous avons construit beaucoup d'outils (partitionnement de graphe, protocol de tolérance aux pannes, ...). En revanche, ce modèle est très contraignant pour la programmation et l'absence d'abstraction pour décrire et manipuler des objets partagés contenant d'autres objets partagés (ex: un array of shared lui même shared) est contraignant. De plus, dans le cas des algorithmes à grain fin, la construction d'une représentation abstraite est un facteur pénalisant.

Les questions qui se posent sont :

1- doit-on conserver cette interface de programmation fork/shared ?

Comme élément de réponse, notons qu'il est possible d'exécuter des programmes série-parallèle récursif avec un surcoût de l'ordre de 1 (à 2 ?) par rapport au même code séquentiel : cela nécessite éventuellement l'écriture d'un compilateur C -> C pour faciliter la génération de code. Cf présentation log_prog que j'ai faite avant de partir.

En pratique cette interface est aussi utilisée pour la construction d'un graphe de tâche en exécutant que les instructions fork d'une tâche, sans exécuter le corps des tâches créées :

2- comment transformer la manière d'interpréter les programmes pour permettre de construire le futur d'une exécution afin d'améliorer l'ordonnancement par un algorithme "statique" ?

Enfin, remarquons que dans le cas des algorithmes adaptatifs, la construction de tâches indépendamment des ressources disponibles est source de surcoût, voir par exemple l'exemple du prefix de Jean-Louis/Daouda. La solution proposée et retenue dans le prototype X-Kaapi est de donner la capacité au programmeur d'explicitement les points où seront fait les extractions de parallélisme. A la charge du système, i.e. de l'algorithme de vol, de faire les bons choix...

Une autre solution serait peut-être d'utiliser une approche à la TBB (Daouda pourra me corriger) : dans TBB dans le cas d'une exécution séquentiel, la création de 2 tâches à un niveau de récursion ne se traduit pas forcément par l'exécution des deux tâches. Dans le cas où la première tâche se termine et s'il n'y a pas eu de requête de vol, la seconde est tout simplement annulée afin que le calcul séquentiel puisse se poursuivre.

Pour terminer, et au vu des très bons résultats expérimentaux sur les algorithmes de la STL parallélisés avec X-Kaapi par Daouda :

3- faut-il utiliser seulement l'interface X-Kaapi pour la programmation parallèle ?

En conclusion il serait important que **TOUS** nous ayons une discussion sur une éventuelle interface de programmation (...).

2- Runtime pour l'ordonnancement local

Dans Kaapi, le runtime abstrait les ressources de calcul en associant un K-processeur par ressource de calcul. Celui-ci étant en charge de gérer l'ordonnancement des calculs sur cette ressource. Le calcul est ensuite représenté par un K-thread qui peut s'exécuter sur l'un des K-processeurs en respectant éventuellement des contraintes de placement. Un K-thread peut être volé afin de lui extraire du travail. De même un K-processeur peut être volé pour lui extraire un K-thread.

Le voleur exécute le code du traitement de vol en concurrence à l'exécution de la victime.

Dans le prototype X-Kaapi, un K-processeur est lancé sur chacun des cœurs choisis lors du lancement de l'application. Le K-thread déroule les instructions du calcul dans lequel des points de vol ont été insérés. Si un ou plusieurs K-processeurs ont posté une requête de vol vers un même autre K-processeur, le K-

thread actif sur le K-processeur victime répond à la requête. La limitation actuelle est qu'il n'est pas possible de définir plusieurs "niveaux" (appelés aussi contexte) de vol correspondant au cas d'un algorithme récursif...

Le voleur et la victime coopèrent pour le traitement des requêtes de vol.

Sachant que ces deux types de vols (concurrent: Kaapi ; coopératif: X-Kaapi) ont leurs avantages/désavantages en fonction de la classe de problème, la question naturelle qui se pose est :

1- comment coupler efficacement un vol concurrent avec un vol coopératif ?

De plus, "qui peut le plus, peut le moins", la question préliminaire est :

2- faut-il conserver un vol coopératif ? Sachant qu'il induit une très forte dégradation des performances dans le cas d'algorithmes à gros grain et que le vol concurrent possède un surcoût raisonnable.

Enfin pour terminer cette section, nous avons observé que nous avons de bien meilleurs performances si les données des algorithmes parallèles sont distribués sur les différents bancs mémoire.

3- quel description hiérarchique de la machine faut-il donner ?

Une description par domaine de partage mémoire hiérarchique tel qu'utilisé dans X-Kaapi (mais hard-codé) ou telle que sortie par les outils libtopology (??) semble naturelle. De plus il est simple d'associer une certaine sémantique pour les déplacements mémoire : pas de contrôle évident pour les caches si ce n'est par localité temporelle ; contrôle possible sur les bancs des architectures NUMA ; contrôle de placement sur les architectures à mémoire distribuée.

Une sous question que l'on doit se poser et qui dépend aussi des algorithmes de la question 5 :

4- faut-il voir la topologie du réseau autrement que de manière hiérarchique ? i.e. doit-on voir la topologie du réseau hypertransport d'Idkoiff ?

5- quel algorithme de vol permet d'exploiter au mieux cette hiérarchie ?

A Jean-Noël de nous donner la réponse théorique et expérimentale : les choix fait dans Kaapi et X-Kaapi de toujours favoriser un vol local avant un vol distant sont à justifier et à améliorer.

3- Communication inter "memory domain"

Enfin pour terminer, il convient de voir comment traiter les aspects communications inter domaine de partage.

Le modèle d'architecture dans X-Kaapi est actuellement le suivant :

- une hiérarchie de "domaines de mémoire" correspondant à la hiérarchie mémoire que l'on retrouve sur les architectures actuelles: cache L1 -> cache L2 -> cache L3 -> banc mémoire -> mémoire partagée -> mémoire distribuée. Certains niveaux faisant apparaître les organisations en cluster ou grid pourrait être ajoutés.

- un ensemble de (K-)processeurs associé à chaque niveau : pour un niveau donné, les processeurs associés à un niveau peuvent communiquer en utilisant la

mémoire de ce niveau.

Une approche similaire a été choisie pour la définition des clusters dans Kaapi, sans descendre au niveau local (mémoire partagée et en dessous).

Pour chaque niveau, il existe une certaine manière de (liste non exhaustive, par exemple ajout d'un broadcast ?) :

- lire une donnée,
- écrire une donnée.

Outre cet aspect transport de données, les threads ont aussi besoin de se synchroniser.

Dans X-Kaapi, la communication & synchronisation entre 2 thread est gérée par une zone tampon (buffer) associée à un état. Seul l'écriture à distance est autorisée (cf papier renpar). Un thread voulant envoyer une donnée à un autre doit :

- écrire dans la zone tampon du destinataire,
- changer l'état associé.

A la charge du destinataire de vérifier le changement d'état. Si ce coût de scrutation est important il convient d'ajouter la capacité à exécuter un traitement (~ interruption chez le destinataire).

Dans Kaapi la communication est gérée 1/ soit à travers la mémoire partagée 2/ soit par message actif.

Il s'agit ici de se poser la question d'une interface uniforme pour la communication entre threads. Il me semble qu'il faille s'orienter vers une interface de type écriture à distance + signalisation et ma dernière question :

1- quelle interface précise pour la communication ?

Specs.

1/ Multi-core

2/ Communication

3/ Algo adaptatif

4/ Programmation des applis de calcul scientifique

5/ Scheduling

Deadlines.