



CENTRE NATIONAL DE LA
RECHERCHE SCIENTIFIQUE



OBSERVATOIRE DE PARIS

STAGE DE FIN D'ÉTUDE

Développements logiciels autour de GDL

Auteur :
Nodar KASRADZE

Responsable :
M. Alain COULAIS



UNIVERSITÉ PARIS 8

Année scolaire 2012-2013

Table des matières

Remerciements	3
Résumé	4
Introduction	5
1 Le cadre de travail	6
1.1 Observatoire de Paris	6
1.2 Matériels utilisés	8
1.3 Logiciels utilisés	10
2 Détails du sujet de stage	12
2.1 Présentation d'IDL	12
2.2 Présentation de GDL	13
2.3 Logiciel libre	15
3 Travail réalisé	17
3.1 Les premiers pas	17
3.2 Fonction d'inversion de matrice	18
3.2.1 Contexte	18
3.2.2 Réalisation	19
3.2.3 Tests	19
3.3 Nombre de threads	20
3.3.1 Contexte	20
3.3.2 Réalisation	20
3.3.3 Étude des performances	21
3.4 Factorisation de Cholesky	23
3.4.1 Description de factorisation de Cholesky	23
3.4.2 La réalisation de CHOLDC, CHOLSOL	24
3.4.3 La réalisation de LA_CHOLDC, LA_CHOLSOL	25

3.5	Convolution	26
3.5.1	Contexte	26
3.5.2	Réalisation	26
3.5.3	Tests	27
3.6	Correction de bugs	28
Conclusion		29
Références		31
Table des figures		32
Abréviations		34
Annexes		i
A	Code de <i>set_num_threads</i>	i
B	Fonction de convolution	iii

Remerciements

Je tiens tout d'abord à remercier mon tuteur, monsieur Alain Coulais, sans qui ce stage n'aurait jamais eu lieu, pour sa gentillesse, ses conseils toujours plus qu'utiles, pour sa présence durant toute la durée de ce stage et enfin, pour tout l'enseignement qu'il a su m'apporter.

Je tiens à remercier monsieur Michel Caillat pour sa présence, sa gentillesse et pour toutes les discussions passionnantes partagées.

Enfin je remercie l'ensemble du personnel du LERMA pour son accueil, sa gentillesse et sa bonne humeur, qui ont rendu mon séjour pendant le stage dans l'Observatoire de Paris des plus agréables.

Résumé

Mon stage s'est déroulé au sein du Laboratoire d'Etudes du Rayonnement et de la Matière en Astrophysique (LERMA) – Observatoire de Paris. Son but a été de contribuer au développement de GDL (clone libre d'IDL - logiciel largement utilisé en astronomie professionnelle) pour atteindre un niveau de maturité et de facilité d'utilisation garantissant la préservation à long terme des capacités de traitements et d'analyse pour de nombreuses expériences en astronomie au sol ou spatiale.

Ce stage était principalement destiné à la mise en place de :

- 1) fonctionnalités manquantes,
- 2) correction des bugs connus,
- 3) tests de régression et performance.

Ce stage nécessitait une connaissance solide du langage C++, car le projet GDL contient plus de 120 000 lignes écrites en ce langage. J'ai aussi eu à apprendre la syntaxe IDL/GDL pour pouvoir écrire des tests qui vont assurer la qualité des nouvelles fonctionnalités.

Summary

I've done my internship at LERMA - Observatoire de Paris. The main goal was to contribute to GDL (free clone of IDL - very popular software in the domain of astronomy) to reach a level of maturity and usability ensuring long term preservation of analysis capabilities for numerous ground experiments and spaces missions based on IDL.

The course of internship was to implement :

- 1) missing functionalities,
- 2) correction of existing bugs,
- 3) performance and regression tests.

Solid knowledge of C++ language was necessary for this internship, because the GDL project contains more than 120 000 lines of code in that programming language. I've learned IDL/GDL syntaxe as well, that I've used to code the tests for ensuring the quality of recently implemented functionalities.

Introduction

La première année de master informatique à l'Université Paris 8 Vincennes-Saint-Denis se termine par un stage d'une durée minimum de 3 mois. J'ai choisi de réaliser ce stage au sein du Laboratoire d'Études du Rayonnement et de la Matière en Astrophysique (LERMA) – Observatoire de Paris, situé à Paris 14ème. Les chercheurs de ce laboratoire travaillent avant tout à faire avancer le front des connaissances dans plusieurs axes de l'astronomie, notamment : la formation et l'évolution des galaxies, la formation de certaines étoiles, la chimie de la poussière interstellaire. Du logiciel libre est utilisé largement et est aussi produit : par exemple pour l'infrastructure de l'interféromètre ALMA et pour GDL destiné à être en libre accès pour tout le monde. C'est ainsi que du 13 mai 2013 au 14 août 2013, je travaille sur un logiciel nommé GDL. Ce logiciel permet de faciliter le traitement des données. Ce projet utilise plusieurs bibliothèques open source, ce qui diminue le travail à réaliser et préserve l'utilisation de code testé pendant plusieurs années par une communauté ouverte.

Ma première approche sur ce projet a été de consacrer une première partie brève du stage à prendre en main les librairies utilisées tout en faisant les diverses installations qui me seront nécessaires pour la conception sur mon poste (GSL, Eigen, PLPlot, GraphicsMagick, ...). Ensuite je prévois une première ébauche de conception me permettant mieux d'assimiler les bibliothèques, suivi d'une étude approfondie de besoin qui précède la conception spécifique de GDL.

Le travail effectué fut aussi réalisé à trois, avec mon tuteur de stage Alain Coulais et Mouadh Ayad, étudiant à l'Université Paris-Sud, IUT de Cachan et lui aussi stagiaire au LERMA travaillant sur un sujet complémentaire au mien, notamment l'élaboration d'une fonction SMOOTH et suite de tests de régression. En plus j'ai été en contact par courriel électronique avec Marc Schellens (fondateur et leader du projet GDL) et d'autres membres de la communauté. Le fait d'avoir travaillé en collaboration avec ces personnes a renforcé le cadre professionnel de ce stage car il est aujourd'hui impossible de travailler seul sur des projets de grande ou même petite envergure.

Dans la première partie de ce rapport je vais décrire le lieu du stage et l'environnement de travail. Après je vais présenter le cahier de charge et le but du stage. Dans la troisième et dernière partie je vais expliquer le travail que j'ai réalisé pendant ces trois mois, les difficultés que j'ai rencontrées, comment j'ai choisi de les résoudre et les résultats obtenus. Pour conclure le rapport, je résume les apports de ce stage, aussi bien personnels que professionnels.

1 Le cadre de travail

1.1 Observatoire de Paris

L'Observatoire de Paris est un Grand Établissement sous tutelle du ministère de l'Enseignement supérieur et de la Recherche. Il a le statut d'EPSCP (Établissement public à caractère scientifique, culturel et professionnel). Ses trois missions principales sont : 1) la recherche, en contribuant au progrès de la connaissance de l'Univers par l'acquisition de données d'observation, le développement et l'exploitation des moyens appropriés, l'élaboration des outils théoriques nécessaires, 2) la formation initiale et continue, 3) la diffusion des connaissances. L'Observatoire de Paris abrite l'École Doctorale Astronomie et Astrophysique d'Île-de-France.

L'Observatoire de Paris compte un grand total d'environ 600 emplois permanents. Du côté MESR - 89 CNAP, 10 EC, 2 PRAG et 232 personnels de soutien ; 248 titulaires du CNRS y sont affectés. L'Observatoire a le statut d'université et dispense un enseignement universitaire de haut niveau en astronomie et astrophysique allant du master au doctorat, avec possibilité de formation à distance, 245 étudiants y sont inscrits (ce chiffre intègre les étudiants des universités partenaires) ; par ailleurs 35 enseignants-chercheurs d'autres établissements de l'enseignement supérieur et 9 personnels divers travaillent dans les laboratoires de l'établissement ; l'établissement emploie 42 contractuels dont 11 sur postes vacants, 19 sur contrat CNRS, le reste sur budget de l'établissement ; son budget annuel hors salaires est de 20 M€ et sa masse salariale est estimée à 35 M€.

		GEPI	LERMA	LESIA	LUTH	SYRTE	IMCCE	Présidence	USN	HORS LES MURS	TOTAL
CNAP	Astronomes	10	4	17	5	4	3	1	0	2	46
	Astronomes-adjoints	9	5	15	1	5	5	0	1	2	43
	Sous total CNAP	19	9	32	6	9	8	1	1	4	89
CNRS	Directeurs de recherche	2	8	13	8	3	1	0	0	0	35
	Chargés de recherche	3	3	11	9	11	4	0	0	0	41
	Sous total CNRS	5	11	24	17	14	5	0	0	0	76
CNU	Professeurs	1	1	4	2	0	2	0	0	0	10
	Maîtres de conférence	2	7	8	3	2	5	0	0	0	27
	Sous total CNU	3	8	12	5	2	7	0	0	0	37
TOTAL		27	28	68	28	25	20	1	1	4	202

FIGURE 1 – Répartition des effectifs (chercheurs et EC) par composante et par corps

Cette institution représente le plus grand pôle national de recherche en astronomie. Il a été fondé en 1667. Près d'un tiers (30 %) des astronomes français y poursuivent leurs recherches au sein de sept laboratoires situés sur trois campus (Paris fondé en 1667 par Louis XIV, Meudon fondé en 1876 par Jules Janssen dans des bâtiments du château de Meudon et Nançay construit en 1953 sur un terrain de l'ENS).

Il existe un partenariat fort entre l'Observatoire et le CNRS/INSU car les laboratoires de recherches sont tous des "Unités mixtes de recherche" et, souvent, en rattachement secondaire à d'autres universités scientifiques.

Les activités de recherche de l'Observatoire de Paris sont structurées autour de :

- 5 UMR associées au CNRS :
 - GEPI : Galaxies, Étoiles, Physique et Instrumentation
 - LERMA : Laboratoire d'Études du Rayonnement et de la Matière en Astrophysique
 - LESIA : Laboratoire d'Études Spatiales et d'Instrumentation en Astrophysique
 - LUTH : Laboratoire Univers et Théories
 - SYRTE : SYstèmes de Référence Temps Espace
- 1 institut qui est aussi une UMR CNRS :
 - IMCCE : Institut de mécanique céleste et de calcul des éphémérides
- 1 unité de service et de recherche :
 - La station radioastronomique de Nançay
- 1 laboratoire de recherche associé ayant pour tutelle principale l'Université de Paris-Diderot, Paris 7 :
 - APC : AstroParticule et Cosmologie
- 1 unité mixte de service, créée en 2002, qui regroupe les services communs et centraux.
- 1 unité de Formation et Enseignement

1.2 Matériels utilisés

J'ai utilisé divers ordinateurs, de l'ordinateur personnel (PC sous Linux) à des serveurs multi-cœurs avec des dizaines d'utilisateurs, tout au laboratoire qu'à la DIO.

La DIO (Division Informatique de l'Observatoire) est un service chargé de tous les aspects informatiques mutualisés de l'Observatoire de Paris : système et réseau, calcul, stockage ainsi que du système d'information, de la bureautique des services centraux et de l'infrastructure technique de l'Observatoire virtuel.

Pour accéder aux services numériques gérés par la DIO un compte LDAP m'a été créé. Via ce compte je pouvais aussi accéder au courriel électronique (@obsmp.fr) et au réseau sans-fils interne. Ce compte LDAP permet aussi d'atteindre le réseau interne via un rebond SSH sur un des deux "firewall" (physiquement 1 à Meudon et 1 à Paris).

Le travail a été effectué sur une machine personnelle, et des serveurs du LERMA et de la DIO :

- ARAMIS
un serveur sous CentOS Linux version 5.9,
avec 8 cœurs (deux Intel® Xeon® CPU X5450 @ 3.00GHz) 64 bits
et 16 Go de mémoire vive.
- M2PARIS (DIO UFE)
un serveur sous Debian Linux version 7,
avec 16 cœurs (deux Intel® Xeon® CPU L5520 @ 2.27GHz) 64 bits
et 12 Go de mémoire vive.
- MEDUSA
un serveur sous CentOS Linux version 5.9,
avec 48 cœurs (quatre AMD Opteron® 6176 SE @ 2.3GHz) 64 bits
et 128 Go de mémoire vive.
- HAKA
une machine personnelle sous Ubuntu Linux version 12.04.2,
avec Intel® Core®2 CPU 6400 @ 2.13GHz 32 bits
et 2 Go de mémoire vive.

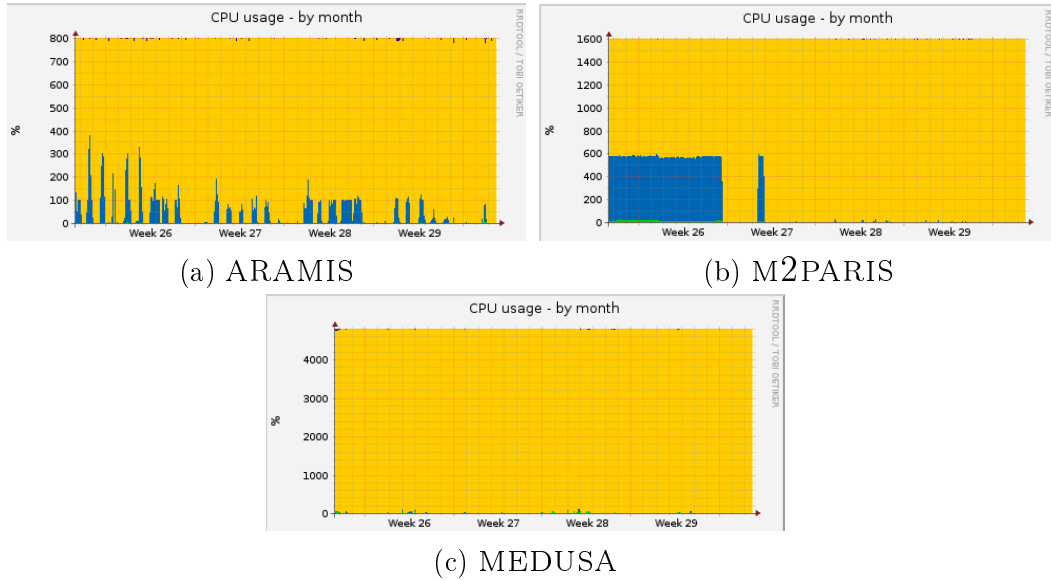


FIGURE 2 – La charge des serveurs pour le mois de juillet 2013 vu par l’outil Mumin

J’utilisais la machine personnelle (HAKA) pour connecter aux autres machines plus puissantes. Compiler un projet, tel que GDL qui contient plus de 120 000 lignes de codes sur un machine comme HAKA va prendre environ 10 minutes, quand sur un serveur comme ARAMIS la compilation va prendre au maximum 2 minutes quand la machine est partiellement chargée.

La charge des serveurs a été une des raisons d’avoir plusieurs machines à disposition, sur la Figure 2 on voit que la machine ARAMIS est partiellement chargé (Figure 2a), la machine M2PARIS est chargé de manière plus permanente (Figure 2b) et la charge de la machine MEDUSA est négligeable (Figure 2c). Le fait d’avoir à disposition des machines qui sont chargés a ses avantages, bien sur la majorité de travail est fait sur une machine avec la charge insignifiante pour diminuer le temps de compilation. Dès que la correction de code est terminé il faut tester le logiciel dans des environnements variés pour conforter son comportement. Notamment la dépendance aux versions de GCC et les problèmes liées aux architectures (32/64 bits). Dans un environnement chargé a été trouvé le problème avec les nombres des threads de OpenMP défini dans GDL (plus de détails dans le section 3.3, page 20).

L’autre a été l’exigence de tester son compatibilité avec des systèmes d’exploitations et des architectures variées.

1.3 Logiciels utilisés

Les systèmes d'exploitation que j'ai utilisés pendant le stage étaient les différentes distributions de Linux et parfois Mac OS X. Pendant le stage je travaillais uniquement en mode commande. Les logiciels les plus utilisés sont cités ci-dessous :

GNU Compiler Collection (GCC)

Crée en 1987 **GCC** est un ensemble de compilateurs portés sur les majorités de systèmes d'exploitation (Linux, Unix, Windows etc...) et de microprocesseurs (x86, AMD64, ARM, SPARC etc...). Il a été créé par le projet GNU comme un compilateur de langage C, mais après quelques années de développement il est devenu beaucoup plus fonctionnel que juste un compilateur.

Aujourd'hui **GCC** est très extensible et adaptable aux besoins des utilisateurs, il contient aussi les compilateurs pour nombreux langages de programmation, tels que C, C++, Objective-C, Java, Ada, Fortran et d'autres langages, plus ou moins connus. **GCC** est le compilateur fourni avec nombreux systèmes d'exploitation comme Linux, BSD, Mac OS X ou BeOS. La plupart des logiciels libres sont compilée avec **GCC**, car en plus d'être performant, il est libre au sens de la licence GNU GPL.

D'après tout ça on comprend pourquoi on a choisi **GCC** comme compilateur pour le projet GDL.

GNU Debugger (GDB)

GDB a été créé en 1988 par Richard Stallman comme le débogueur standard du projet GNU. Comme pour **GCC** il existe des ports de **GDB** sur différent types de microprocesseurs et systèmes d'exploitation, il supporte différents langages de programmation et il est distribué sous les termes de la licence GNU GPL.

GDB propose de nombreuses fonctionnalités pour tracer et contrôler l'exécution d'un programme informatique pas à pas. On peut surveiller ou modifier (si il y a besoin) une variable interne, on peut aussi faire un appel à fonction, indépendamment du comportement prévu dans le programme. **GDB** permet le débogage distant selon gdbserver, qui est lancé sur la même machine que le programme et le débogage est effectué d'une autre. Cette technique est souvent utilisée pour déboguer un programme sur un système embarqué.



(a) GCC



(b) GDB



(c) VALGRIND

FIGURE 3 – Les logos

Valgrind

L'une des plus grandes difficultés de C++ est de gérer correctement la mémoire, alors pour diminuer les problèmes liés à la fuite de mémoire (l'absence de désallocation de l'espace utilisé, double désallocation de l'espace utilisé etc...) j'ai utilisé **Valgrind**.

A l'origine il a été créé comme un débogueur de mémoire pour Linux sur architecture x86, mais depuis il est devenu une plate-forme d'analyse dynamique pour les contrôleurs de la mémoire et "profilers". Il est disponible pour les systèmes d'exploitation Linux, Mac OS X et Android sur les microprocesseurs x86, x86-64, PowerPC et ARM. Grâce à sa licence GNU GPL v2, il existe des ports non-officiels sur d'autres systèmes et architectures.

Essentiellement **Valgrind** est une machine virtuelle utilisant la technique de compilation à la volée, le programme d'origine ne s'exécute pas directement sur un processeur, avant d'exécution il est traduit dans une représentation intermédiaire simplifiée, les instructions sont ajoutées par un outil ("memcheck", "profiler" ou un outil externe) et finalement ces instructions sont traduites dans un code machine. A cause des procédures de traduction sous **Valgrind** l'exécution d'un programme est au moins 5 fois plus lente.

Il est possible de combiner **Valgrind** avec **GDB** pour profiter des avantages de ces deux excellents outils.

2 Détails du sujet de stage

2.1 Présentation d'IDL

IDL - Interactive Data Language est un langage interprété sous licence propriétaire. Créé en 1977 il a attiré l'attention des astronomes, médecins et d'autres scientifiques qui ont eu besoin d'un langage assez simple à apprendre et assez puissant pour traiter les données massives avec des fonctionnalités graphiques. La popularité d'**IDL** augmentait très vite à cause de ses fonctionnalités très pratiques pour les chercheurs, qui n'ont pas été programmeurs, mais qui ont besoin d'un outil pour faire des calculs scientifique. Entre ces fonctionnalités on trouve :

- typage dynamique, le type est attribué automatiquement pendant la définition de sa valeur (il y a une quinzaine de types) ;
- capacités d'un langage de programmation orientée objet ;
- interprétation par machine virtuelle ;
- affichage de tracés (1D, 2D, 3D) sans utilisation des bibliothèques graphiques externes ;
- **IDL** est un langage vectorisé, les opérations sur scalaires sont généralisées pour pouvoir les utiliser sur les vecteurs, matrices et les tableaux à une dimension supérieure ;
- multitâche, certaines fonctions et procédures sont parallélisés.

Grâce à ces avantages et une disponibilité sur les systèmes d'exploitations les plus utilisés y compris Windows, Mac OS X, Linux et Solaris, aujourd'hui il existe de nombreuses bibliothèques écrites complètement en syntaxe **IDL**. Ce sont les modules spécifiques à un domaine, qui ont été créés par les scientifiques pour simplifier l'utilisation de langage **IDL** dans ses domaines. Par exemple : *IDL Astronomy User's Library* ou plus court *AstroLib* créée par Wayne Landsman. Cette bibliothèque contient des procédures fréquemment utilisées par les astronomes, notamment pour relire certaines formats de données (FITS).

IDL est un langage de programmation dont l'interpréteur du langage (ou compilateur) est propriétaire, ce qui signifie qu'il n'est pas libre. Cela signifie aussi que les scientifiques n'ont pas accès au code source et qu'il ne peuvent ouvrir qu'une quantité limitée de sessions au travers d'un système de jetons.

Le système de jetons dans **IDL** est contraignant. Déjà, l'abonnement pour un jeton par an coûte plus de 1000 €. En plus, il y a un nombre limité d'ouverture d'**IDL**, car à chaque démarrage du logiciel, ce dernier va réserver des jetons disponibles sur un serveur. Cela implique aussi que l'utilisateur doit

rester connecté au serveur détenteur de jetons, ce qui peut très vite s'avérer compliqué en déplacement professionnel, durant des présentations ou lors de démonstration pendant des conférences.

Mais ce qui pose encore plus de problèmes aux scientifiques, c'est la pérennité du logiciel. En effet, ne pas avoir les codes source à disposition implique le fait que si la société en charge de garder le logiciel à jour décide de cesser ses activités, personne ne pourra continuer le développement. Dans le passé, la NASA et l'ESA ont contraint l'éditeur, qui voulait ne garder que Windows, à poursuivre le développement sous Mac OS X et Linux. Or, la technique évolue tous les jours, et les logiciels tels qu'est **IDL** se doivent de la suivre. De plus, s'il n'est plus possible d'utiliser **IDL** dans le futur, les nombreuses bibliothèques écrites en syntaxe **IDL** deviendront inutiles et il faudra recoder tous les algorithmes dans un autre langage de programmation, ce qui serait très long et extrêmement cher.

2.2 Présentation de GDL

GDL - GNU Data Language est un compilateur libre compatible avec la syntaxe **IDL**. Il a été créé par Marc Schellens en 2004 et il est développé par une communauté internationale de volontaires. **GDL** a été créé pour bénéficier des avantages d'**IDL** en ajoutant la liberté de compilateur pour assurer la pérennité de logiciel lui même et les bibliothèques tierces écrites en syntaxe **IDL**. L'interprétation de syntaxe est basé sur ANTLR.

GDL est déjà une alternative viable couvrant la plupart des fonctionnalités d'**IDL** et peut être utilisé pour compiler de larges "pipelines" de codes et bibliothèques en syntaxe **IDL**, lire les données de divers formats utilisées par les scientifiques (comme FITS, HDF, NetCDF, XDR), faire des calculs complexes et retourner un résultat exact. Grâce aux packagers les versions pré-compilées et pré-configurées sont mises à disposition pour la distribution de Linux les plus populaires (Debian, Ubuntu, Fedora, Gentoo, ArchLinux...), FreeBSD et Mac OS X. Pour les utilisateurs plus expérimentés la version la plus récente du code est disponible sur un CVS (accès en lecture seule à tous), qui peut être compilée sur presque tous les systèmes basés sur Unix.

La compatibilité avec Windows n'est pas assurée, mais toutes les nouvelles fonctionnalités prennent en compte la compatibilité avec ce système et les fonctionnalités déjà présentes dans **GDL** continuent d'évoluer pour être compatibles avec tous les systèmes possibles.

Le choix de bonnes bibliothèques externes et de bons algorithmes internes, l'usage de OpenMP, donnent une bonne performance globale à GDL. Sur les

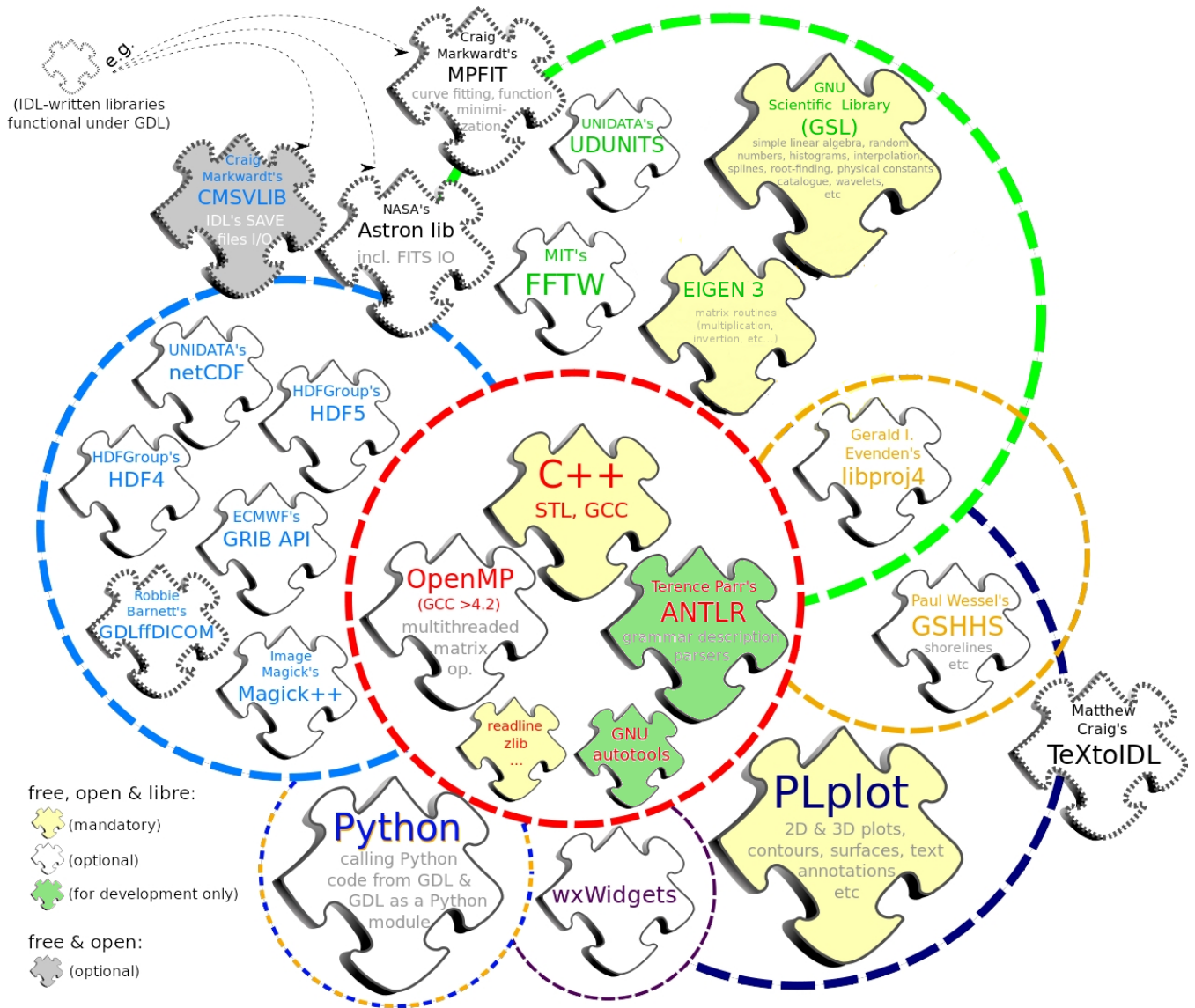


FIGURE 4 – Dépendances de GDL

multi-cœurs, les gains attendus sont généralement constatés. Les tests de performances actuels montrent que les calculs avec **GDL** sont aussi rapides (parfois plus vite), que les calculs avec son adversaire propriétaire et ce, en respectant le gain par nombre de cœurs. L'ultime déficience connue dans le benchmark de référence (IDL TIME_TEST_4) est SMOOTH, sujet du second stagiaire. Le plus grand inconvénient de **GDL** est l'absence de plusieurs fonctions et procédures présentes dans **IDL**, alors ma mission a été de di-

minuer le nombre de ces fonctions/procédures manquantes et de tester ces nouvelles fonctionnalités en comparant les performances de **GDL** et d'**IDL**.

Le projet **GDL** utilise la licence GNU GPL v2/v3. Les bibliothèques tierces utilisée dans ce projet sont uniquement des bibliothèques libres (la notion d'un logiciel libre/licence libre est définie dans la section 2.3, page 15) avec la licence GPL ou un autre, compatible avec GPL. Une exception unique concerne une fonctionnalité importante d'un format propriétaire de fichier interne à **IDL**. Il se trouve qu'une librairie externe open source à la licence incompatible avec la GNU GPL permet de gérer ce format. Il revient à l'utilisateur d'ajouter lui même ce code.

2.3 Logiciel libre

La notion de Logiciel Libre (Free Software) est inventée par Richard Stallman dans le début des années 1980. Une première définition est publiée en 1986 par FSF (Free Software Foundation), une organisation américaine à but non lucratif, dont la mission mondiale est la promotion du logiciel libre et la défense des utilisateurs. Cette définition se traduit par une tente d'une licence définissant les droits et les devoirs des utilisateurs en s'appuyant sur le droit d'auteur (Copyright). Rarement, la FSF modifie la notion de Logiciel libre afin de la clarifier ou pour résoudre des questions portant sur des points difficiles. Les deux principales licences GNU GPL sont les versions 2 et 3.

Selon la dernière définition, un programme est un logiciel libre si, en tant qu'utilisateur de ce programme, vous avez les quatre libertés essentielles :

liberté 0

la liberté d'exécuter le programme, pour tous les usages ;

liberté 1

la liberté d'étudier le fonctionnement du programme, et de le modifier pour qu'il effectue vos tâches informatiques comme vous le souhaitez ;
l'accès au code source est une condition nécessaire ;

liberté 2

la liberté de redistribuer des copies, donc d'aider votre voisin ;

liberté 3

la liberté de distribuer aux autres des copies de vos versions modifiées ;
en faisant cela, vous donnez à toute la communauté une possibilité de profiter de vos changements ; l'accès au code source est une condition nécessaire.

Pour protéger la [non]liberté d'un logiciel et les droit d'utilisateur il existe plusieurs type de licences. La FSF classe les licences en trois catégories :

- les licences libres compatibles avec GPL ;
- les licences libres non compatibles avec GPL ;
- les licences non libres/propriétaires.



FIGURE 5 – Logo de GPL version 3

GNU General Public License

La Licence Publique Générale GNU (GNU GPL ou GPL) est une licence des logiciels libres. Elle a été créé par Richard Stallman en 1989 pour permettre au grand nombre de projets de partager leur code source. C'est la licence le plus utilisée dans le monde de logiciels libres. La dernière version est "GNU GPL version 3" publiée en 2007. Ses termes autorisent l'utilisateur de logiciel sous licence GPL à : distribuer ce logiciel, le modifier, distribuer la version modifié. L'utilisateur n'est pas autorisé à : changer la licence, distribuer ce logiciel sans code source, commercialiser ce logiciel.

Ce ne signifie pas qu'on peut pas faire de l'argent avec les logiciels libres. L'entreprise multi-nationale Red Hat avec un chiffre d'affaire d'un milliard de dollars est dédiée aux logiciels libre et Open Source, qui a basé son modèle d'affaires sur la vente des services et supporte pour les logiciels libres, notamment le noyau Linux et le système d'exploitation RHEL (Red Hat Enterprise Linux).

Berkeley software distribution license

La licence BSD est une licence libre compatible avec la GPL. Ses termes sont moins restrictives que celles de GPL. Un logiciel sous licence BSD en différence de GPL peut être utilisé dans un autre projet commercial et l'accès au code source modifié n'est pas garanti.

End User License Agreement

Contrat de Licence Utilisateur Final est un licence non libre/propriétaire, ses termes sont définis par l'entreprise publiant le logiciel. Il est utilisé pour mettre les différent restrictions (restriction de distribution ; restriction sur le nombre d'utilisateurs, qui peuvent utiliser ce logiciel).

3 Travail réalisé

3.1 Les premiers pas

Mon stage a commencé par mon immersion dans le projet GDL. Ma première tâche a consisté à compiler le code source de GDL sur les différentes machines (HAKA, ARAMIS...) avec les différentes options (en activant ou non certaines bibliothèques externes facultatives). Pour chaque machine la procédure de compilation a été différente, car les versions de logiciels/bibliothèques installés ont été différentes ou sur certaines machines les logiciels/bibliothèques nécessaires pour la compilation n'étaient pas installés du tout. Les deux cas veulent dire que ça va prendre encore plus du temps de compiler un projet sur de telles machines. Trouver une version de logiciel qui va être compatible avec des logiciels/bibliothèques déjà installé ce n'est pas une tâche triviale d'autant que le projet GDL contient 5 bibliothèques obligatoires et plus d'une dizaine de bibliothèques optionnelles.

Le build du projet GDL est effectué par *Autotools*, mais dans la version suivante 0.9.4 *CMake* va remplacer *Autotools*, ce changement est déjà accompli dans la version CVS du projet. Pour mieux comprendre la raison pourquoi le projet va abandonner *Autotools* voici le bref description de ces deux outils :

Autotools

Autotools (ou GNU build system) est un ensemble des projets GNU utilisés pour le build de projet sur différents systèmes d'exploitations basés sur Unix ; ces outils peuvent être utilisés sur Windows, mais c'est une procédure très restrictive (compilation est possible seulement avec GCC) et prend plus de temps. Il est basé sur shell script, alors il ne nécessite pas l'installation des outils *Autotools* pour faire un build. Ces outils sont très populaires dans le monde Linux. Les outils *Autotools* ont de nombreux inconvénients à cause desquels les développeurs cherchent d'autres solutions. Les problèmes qui repoussent des programmeurs sont :

- la syntaxe très compliquée, qui rend difficile l'écriture d'un script de configuration ;
- large nombre des outils composant avec les syntaxes différents ;
- les messages d'erreurs difficiles à comprendre ;
- création de larges scripts de configuration même pour un projet basique ;
- difficilement extensible, difficulté à rajouter des fonctionnalités non standard.

CMake

Cross Platform Make ou simplement *CMake* est un moteur de production créé dans le début des années 2000 pour être compatible avec divers plate-formes (à peu près tous les systèmes basée sur Unix, Windows : Borland, Visual C++, cygwin...) et divers environnements de développement (KDevelop, XCode, Visual Studio...). Il devient un outils de plus en plus utilisé à cause de ces nombreux avantages :

- ne dépend que d'un compilateur C++ quelconque ;
- syntaxe très simple à apprendre ;
- il génère un seul Makefile pour tous les plate-formes supportées ;
- les messages d'erreurs aident à comprendre le problème ;
- supporte la plate-forme des tests, qui facilite la création des suites des tests ;
- la performance améliorée par rapport aux *Autotools*.

Au début de projet GDL l'utilisation de *CMake* a été impossible à cause de sa licence non compatible avec GPL (GDL utilise la licence GPL), mais la migration sur la licence BSD a aidé à surmonter l'obstacle de non compatibilité. Le point le plus faible de *CMake* représente sa documentation.

Par exemple, si on veut compiler sur ARAMIS la version CVS du projet, on peut utiliser la suite des commandes suivants (au début on suppose qu'on est dans la répertoire `~/sources/gnudatalanguage/gdl/`) :

```
mkdir compil build
cd compil
cmake .. -DHDF=OFF -DPSLIB=OFF
        -DPLPLOTDIR=~ /sources/plplot-5.9.6/build/
-DCMAKE_INSTALL_PREFIX=~ /sources/gnudatalanguage/gdl/build/
make -j 8
make check
```

3.2 Fonction d'inversion de matrice

3.2.1 Contexte

Ma première mission d'écriture de code a été le développement de la fonction d'inversion de matrice. En fait la fonction d'inversion de matrice (INVERT) était déjà présente dans le projet GDL, mais elle a été écrite en utilisant la librairie GSL (GSL représente des outils de calculs numériques en

mathématiques appliquées, elle fait partie du projet GNU et est distribuée selon les termes de la licence GNU GPL). Depuis début 2013, l'équipe GDL mène des essais très fructueux avec la librairie Eigen. INVERT est un code important, il fallait voir si un gain de temps notable existe avec Eigen par rapport à la GSL. Alors j'ai dû réécrire INVERT en utilisant la librairie Eigen (librairie de C++ contenant des templates qui implémentent l'algèbre linéaire et des opérations sur les matrices, sous la licence libre MPL2, qui représente l'hybride de BSD et GPL).

3.2.2 Réalisation

La philosophie de GDL est directement inspirée du monde libre, ainsi, si un code existe déjà et qu'il est possible de l'intégrer dans le projet grâce à sa licence, il est inutile de perdre du temps à le réécrire. On peut aussi choisir entre plusieurs librairies externes en fonction de critères : simplicité, taille, performance, évolution. C'est pourquoi, plutôt que de coder l'algorithme d'inversion de matrice j'ai préféré utiliser les routines existantes et largement testées.

Dans la fonction d'inversion, la matrice est donnée sous les types du langage C/C++. Pour mieux traiter les données Eigen à ses propres types de données et pour convertir les données existant de type non-Eigen à un type interne d'Eigen on a le mapping. La classe Map de bibliothèque Eigen contient des routines de mapping. Ces routines sont très simples à utiliser et elles sont très performantes.

Exemple de mapping :

```
float array[rows*cols];  
Map<MatrixXf> m(array,rows,cols);
```

3.2.3 Tests

Avant de publier dans le CVS la nouvelle fonction d'inversion d'une matrice il faut la bien tester. Les tests sont nombreux et divers :

- tests de justesse,
- tests de types,
- tests de correspondance avec la syntaxe IDL,
- tests de performance.

Ces tests sont nécessaires pour assurer la qualité de projet et d'éviter les problèmes dans le futur. Les résultats des tests de performance seront discutées dans la section suivante.

3.3 Nombre de threads

3.3.1 Contexte

Quand la fonction d'inversion de matrice utilisant la librairie Eigen (INVERT_EIGEN) a été prête, ses benchmarks sur une machine chargée ont montré un résultat complètement différent de celui de la machine non chargée. Par rapport à IDL la même fonction d'inversion de matrice sur une machine non chargée a été plus rapide pour tous les types, mais sur la machine chargée (même s'il y a eu un seul cœur occupé et 15 disponibles) pour certaines types la performance a baissé de 10 fois ou même plus. D'après l'analyse de la documentation de la librairie Eigen et la consultation (par courriel électronique) avec les autres membres de la communauté (Marc Schellens, Gilles Duvert...) on a trouvé que le problème est causé par le nombre de threads défini statiquement et indépendant de la charge de machine pour OpenMP (Eigen utilise OpenMP pour bénéficier des multi-cœurs).

Cette occasion a produit un devoir inattendu : j'ai dû écrire une fonction qui va définir le nombre de threads dynamiquement avant chaque exécution du GDL et cette fonction va prendre en considération la charge de la machine sur laquelle elle est exécutée.

3.3.2 Réalisation

Sous les systèmes d'exploitations basées sur Unix on calcule le nombre de threads optimal (*suggested_num_threads*) par la formule suivante :

$$suggested_num_threads = nbofproc - avload \quad (1)$$

Le nombre de processeurs sur la machine (*nbofproc*) est connu pour OpenMP et c'est un simple appel à sa fonction interne. Alors, il nous reste de trouver la charge moyenne de la machine (*avload*). La charge moyenne de la machine est définie comme un nombre réel de type X.XX (par exemple : *avload*=2.64 sur une machine avec 8 processeurs virtuels, où 2 signifie, que deux processeurs sont entièrement occupés et .64 signifie, que le troisième est chargé partiellement ; les cinq processeurs sont libres), de dernière(s) 1, 5

ou 15 minutes. Dans notre cas on a décidé d'utiliser le moyen de 5 derniers minutes, car une minute c'est un intervalle très court pour faire une conclusion sur la charge de la machine et 15 minutes c'est trop long (on veut pas prendre en compte la charge de la machine par un logiciel qui a déjà terminé son fonctionnement). Puisque le nombre de processeurs est un nombre entier et la charge moyenne est un nombre réel, la valeur de la charge moyenne est arrondi au plus proche.

Sous le système d'exploitation Windows le nombre de threads optimal est calculé par la formule un peu différente :

$$suggested_num_threads = nbofproc - avload * nbofproc / 100 \quad (2)$$

La différence vient du format de présentation de charge de la machine, qui est représenté sous Windows comme un pourcentage d'occupation de tous les processeurs virtuels (par exemple : *avload*=50 sur une machine avec 8 processeurs virtuels signifie que 4 processeurs sont occupés et les 4 restants sont libres).

Le code source de la fonction *set_num_threads* est disponible dans l'Annexe A sur page i.

3.3.3 Étude des performances

L'aspect très important dans la concurrence entre GDL et IDL c'est la performance, afin qu'ils sont utilisés pour les calculs scientifiques. Signifiant, que les données à traiter sont assez larges (par exemple : les matrices avec des dimensions supérieures à 200). La performance de la majorité des fonctions de GDL est très proche ou même meilleure que la performance d'IDL, mais pour attirer plus d'utilisateurs il faut avoir un meilleur indice de performance que IDL pour tous les fonctions sur tous les systèmes d'exploitations.

Pour trouver les indices de performance des différents fonctions d'inversion de matrice une suite de tests a été écrite, ces tests se composent de calcul du temps d'inversion des matrices aux très larges dimensions (255, 256, 300, 500). Les tests ont été faits sur ARAMIS, qui pour le moment de tests a eu 2 cœurs occupés en permanence par des autres processus. Cette condition a été idéale pour tester la performance de cette nouvelle caractéristique qui a été introduite par la fonction de définition de nombre de threads.

D'après les tests, la fonction d'inversion de matrice utilisant la bibliothèque Eigen donne une meilleur performance que celle utilisant la bibliothèque GSL. Mais sur une machine chargée les deux fonctions deviennent

Type	GDL INVERT ^a		IDL	GDL INVERT ^b	
	GSL	EIGEN		GSL	EIGEN
BYTE	7.383	6.249	0.543	0.663	0.397
INT	5.242	6.125	0.530	0.657	0.374
LONG	5.108	5.605	0.537	0.656	0.365
FLOAT	7.768	5.283	0.533	0.630	0.339
DOUBLE	5.645	6.580	1.161	0.690	0.466
COMPLEX	10.218	9.467	3.991	2.004	1.651
DCOMPLEX	13.543	2.675	4.365	2.160	2.089
UINT	9.204	7.416	0.520	0.671	0.399
ULONG	5.546	5.238	0.522	0.666	0.397

TABLE 1 – Résultats des tests de performance sur la fonction INVERT, utilisant la GSL ou Eigen

a. sans `set_num_threads` sur une machine chargé

b. avec `set_num_threads` ou sans `set_num_threads`, mais sur une machine non chargée

incomparable avec la fonction de IDL, qui est largement plus rapide. Alors, après la fonction `set_num_threads` la performance de la fonction d'inversion de matrice avec GSL (qui est utilisé dans la plupart de fonctions codés sous C++ pour les calculs scientifiques) est proche de la fonction présent dans IDL et parfois est meilleur.

Une fois les problèmes du nombre des threads optimale sont réglés, la fonction d'inversion de matrice avec Eigen a les meilleurs résultats pour tous les types (la Table 1 montre les temps de calcul pour les multiples fonctions). En plus d'une bonne performance, Eigen est simple à utiliser et il se charge de la gestion de la mémoire, qui dispense l'utilisateur de gérer les éventuelles fuites mémoire. C'est pourquoi Eigen a été choisi comme une bibliothèque pour les calculs matriciels pour les nouvelles fonctionnalités, qui vont faire partie de GDL version 0.9.4.

En conséquence, la définition dynamique du nombre de threads a permis d'obtenir un meilleur indice de performance pour l'inversion de matrice aussi bien que pour tous les fonctions utilisant la parallélisation avec OpenMP dès que les machines sont un peu chargées.

3.4 Factorisation de Cholesky

La tâche suivante a été de coder sous C++ la factorisation (décomposition) de Cholesky d'une matrice. Dans IDL on a quatre procédures/-fonctions concernant la factorisation de Cholesky (CHOLDC, CHOLSOL, LA_CHOLDC, LA_CHOLSOL). Tous ces procédures/fonctions étaient absentes dans GDL et des utilisateurs souhaitent en bénéficier. Certaines codes en syntaxe IDL/GDL (PlanckSkyModel, iCosmo) utilisent CHOLDC et ne pouvaient être totalement audités. Ma mission a été de les rajouter en utilisant la bibliothèque Eigen.

Cette factorisation est utilisée pour la résolution des systèmes d'équations linéaires, simulations avec méthode de Monte-Carlo, filtre de Kalman, l'inversion d'une matrice hermitienne.

3.4.1 Description de factorisation de Cholesky

La factorisation de Cholesky, nommée d'après André-Louis Cholesky, consiste, pour une matrice symétrique définie positive \mathbf{A} , à déterminer une matrice triangulaire inférieure \mathbf{L} telle que :

$$\mathbf{A} = \mathbf{L}\mathbf{L}^T \quad (3)$$

Où \mathbf{L}^T est la transposée de la matrice triangulaire inférieure \mathbf{L} . On peut imposer que les éléments diagonaux de la matrice \mathbf{L} soient tous positifs, et la factorisation correspondante est alors unique.

$$\begin{pmatrix} 4 & 12 & -16 \\ 12 & 37 & -43 \\ -16 & -43 & 98 \end{pmatrix} = \begin{pmatrix} 2 & & \\ 6 & 1 & \\ -8 & 5 & 3 \end{pmatrix} \begin{pmatrix} 2 & 6 & -8 \\ & 1 & 5 \\ & & 3 \end{pmatrix}$$

FIGURE 6 – Exemple de factorisation de Cholesky

Pour la résolution d'un système d'équations linéaires, la factorisation de Cholesky est approximativement deux fois plus efficace que la décomposition LU. Alors c'est très important d'avoir dans GDL les routines concernant la factorisation de Cholesky.

3.4.2 La réalisation de CHOLDC, CHOLSOL

Le première ensemble de factorisation de Cholesky que j'ai codé a été la procédure **CHOLDC**. Cette procédure est utilisée pour trouver une matrice triangulaire inférieure d'une matrice symétrique définie positive donnée. Coder cette procédure a été une tâche sans peine, car j'ai déjà eu l'expérience de codage avec la librairie Eigen, que j'ai obtenu pendant ma mission précédent (la fonction d'inversion de matrice). En plus Eigen contient tous les algorithmes de suite de Cholesky, alors j'ai dû à faire la répartition des appels des algorithmes pour les types appropriés du côté de GDL.

Ensuite j'ai implémenté la fonction **CHOLSOL**. Cette fonction calcule la solution pour un système d'équations linéaires $\mathbf{Ax}=\mathbf{B}$. Cette fonction prend comme premier argument la matrice calculée par **CHOLDC**, mais GDL en différence d'IDL utilise la matrice triangulaire supérieure qui contient les valeurs initiales. La différence est à cause de librairie Eigen, qui ne support pas la résolution d'un système d'équations linéaires d'après les valeurs intermédiaires (la matrice décomposée).

```
GDL> A = [[ 6.0, 15.0, 55.0], $
GDL>      [15.0, 55.0, 225.0], $
GDL>      [55.0, 225.0, 979.0]]
GDL>
GDL> B = [9.5, 50.0, 237.0]
GDL>
GDL> CHOLDC, A, P
GDL>
GDL> PRINT, CHOLSOL(A, P, B)
      -0.500000   -1.000000    0.500000
GDL> □
```

(a) GDL

```
IDL> A = [[ 6.0, 15.0, 55.0], $
IDL>      [15.0, 55.0, 225.0], $
IDL>      [55.0, 225.0, 979.0]]
IDL>
IDL> B = [9.5, 50.0, 237.0]
IDL>
IDL> CHOLDC, A, P
IDL>
IDL> PRINT, CHOLSOL(A, P, B)
      -0.499998   -1.000000    0.500000
IDL> □
```

(b) IDL

FIGURE 7 – Exemple de résolution d'un système d'équations linéaires dans les interpréteurs GDL et IDL

La Figure 7 montre un exemple de résolution d'un système d'équations linéaires $\mathbf{Ax}=\mathbf{B}$ dans différents interpréteurs en utilisant la procédure **CHOLDC** et ensuite la fonction **CHOLSOL**. Pour résoudre un système d'équations linéaires dans IDL on est obligé d'exécuter les deux procédures, mais dans GDL on peut le résoudre sans utilisation de **CHOLDC**. Puisque dans la fonction **CHOLSOL** on fait la décomposition (c'est obligatoire en raison d'utilisation de la bibliothèque Eigen). Pour assurer la compatibilité on peut résoudre $\mathbf{Ax}=\mathbf{B}$ à la manière d'IDL (suite de deux commandes).

La besoin de calculer deux fois la même factorisation baisse la performance de GDL dans résolution d'un système d'équations linéaires, mais GDL a un autre avantage par rapport à IDL : c'est une solution plus exacte. Dans un exemple qui est représenté sur la Figure 7 la solution exacte est un vecteur avec les valeurs : $\{ -0.500000, -1.00000, 0.500000 \}$ - ce que retourne GDL, mais le vecteur calculé par IDL est un peu différent, la première valeur contient un erreur de -0.000002.

3.4.3 La réalisation de **LA_CHOLDC**, **LA_CHOLSOL**

La procédure **LA_CHOLDC** est la procédure **CHOLDC** généralisée pour trouver la décomposition d'une matrice à coefficients complexes (matrice hermitienne). Pour la factorisation d'une matrice à coefficients réels, en plus de la formule (3) (page 23) qui est utilisée par la procédure **CHOLDC**, **LA_CHOLDC** peut utiliser la formule (4), où \mathbf{U} est une matrice triangulaire supérieure et \mathbf{U}^T est sa transposée :

(4)

Pour une matrice hermitienne la procédure **LA_CHOLDC** peut utiliser l'une des formules suivantes :

$$\mathbf{A} = \mathbf{U}^H \mathbf{U} \quad (5)$$

$$\mathbf{A} = \mathbf{L} \mathbf{L}^H \quad (6)$$

\mathbf{H} signifie une matrice adjointe (aussi appelée matrice transconjugée).

La réalisation de **LA_CHOLDC** a été similaire à la réalisation de **CHOLDC**, car la bibliothèque Eigen gère les types complexes aussi bien que les autres types.

Le code complet des routines Cholesky est disponible dans le CVS du projet <http://smarturl.it/Cholesky>

3.5 Convolution

3.5.1 Contexte

La procédure **CONVOL** convole un tableau (1D, 2D, 3D...) avec un noyau et retourne un tableau modifié. La convolution est largement utilisée dans divers domaines des sciences pour traitement d'image, traitement de signal, différenciation et beaucoup d'autres opérations. Le traitement d'images fait un grand partie d'astronomie, comme les image obtenu par les satellites parfois ont des défauts, on a besoin d'une procédure de correction. Alors c'est très important d'avoir la procédure de convolution dans GDL.

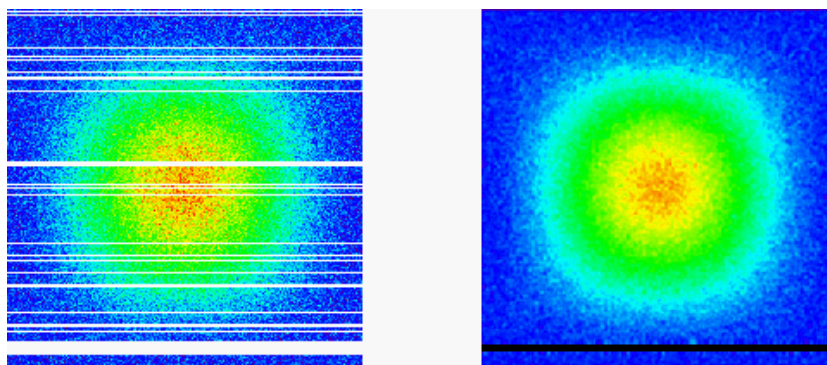


FIGURE 8 – Exemple d'utilisation de convolution en traitement d'image en IDL. Cette fonctionnalité n'est pas encore disponible dans GDL

La Figure 8 est un exemple de lissage d'image bruitée avec les données manquantes. Après la procédure de convolution la majorité des données manquant est restauré et le bruit est réduit. On peut appliquer les différents effets sur cette image, l'effet est défini par le noyau (aussi appelé filtre) avec le-quelle image est convolé.

3.5.2 Réalisation

La fonction de convolution est présente dans GDL depuis 2004, mais elle est resté inchangée malgré les changements dans cette fonction dans IDL (nouveaux mot-clés ajoutés). D'ajouter la prise en compte effective de certaines nouveaux mot-clés dans **CONVOL** de GDL m'a été confié. J'ai commencé la mise à jour de cette fonction par regroupement de plusieurs fonctions concernant la convolution dans un même fichier. Cette stratégie est efficace pour le gain du temps et compréhension de vieux code afin qu'on ne cherche pas les parties du code.

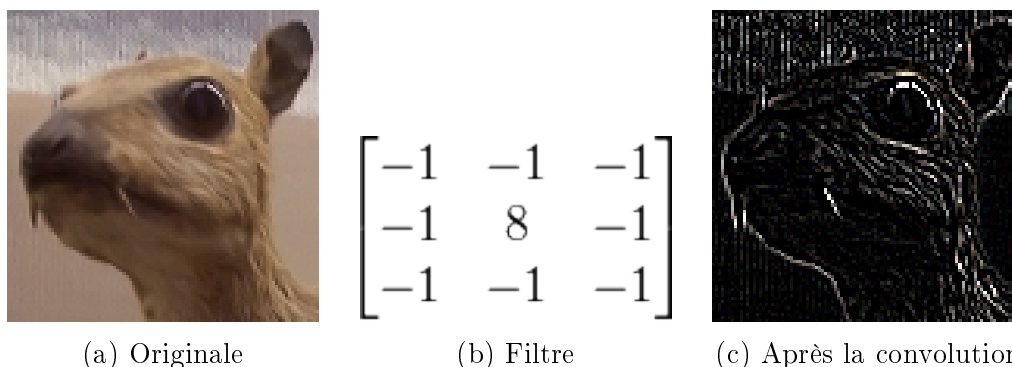


FIGURE 9 – Exemple d'utilisation de convolution pour détecter les contours

Le code complet de fonction **CONVOL** avec tous les mot-clés déjà ajoutés est très large et la compréhension du code a pris beaucoup de temps. L'étude de code m'a pris environs une semaine. En plus des difficultés pour comprendre le code, comme il n'est pas en maintenance depuis longtemps la partie de code a perdu son efficacité. Donc, trouver les parties faibles du code et les corriger a fait la partie de cette mission.

Pendant la réalisation de mis à jour de la fonction de convolution la tâche la plus difficile a été d'assurer son fonctionnement pour tous les types. La grande partie du code est généralisé pour tous les types (en utilisant les Templates). Mais, comme les différents types ont différents caractéristiques on a eu besoin de réécrire la partie de la logique plusieurs fois en fonctionnement du type utilisé et pendant la compilation inclure le code approprié à ce type. Par exemple le fait que les nombres complexes n'ont pas d'ordre a causé des problèmes quand on a eu besoin de les comparer.

Le code de la fonction de convolution est dans l'Annexe B, page iii. Ce code ne représente pas le code complet de **CONVOL** de GDL, le code complet est disponible dans le CVS du projet <http://smarturl.it/Convol>

3.5.3 Tests

Des tests pour la convolution étaient absent dans la suite de tests du projet. Alors j'ai écrit des tests en syntaxe IDL/GDL pour assurer que tous les types sont bien gérés. Après l'application de ce test aux différents environnements, on a trouvé que les type **ULONG**, **ULONG64** ont des problèmes avec GCC version jusqu'à 4.4 et le projet ne compile pas. Pour l'instant ces deux types sont interdits dans **CONVOL**.

3.6 Correction de bugs

Pendant mon stage dans l'intervalle entre les grandes missions j'ai eu des "petit" tâches. Ces tâches comprennent la correction de divers bugs dans GDL. L'un de ces bug a été un problème typique du langage C++ - fuite mémoire (appel de **delete** pour libérer la mémoire alloué par **new**). Un autre crash de GDL était causé par la copie de bloc de mémoire de **long** à **double**. Ces deux types ont les longueurs différentes pour les différentes architectures (32/64 bits). Sur une architecture de 32 bits les deux ont la longueur maximale de 32 bits, mais sur l'architecture de 64 bits ils ont les longueurs différentes (**long** - 32 bits, **double** - 64 bits). Ces fuites mémoire ont été trouvées grâce au logiciel Valgrind.

J'ai aussi corrigé un bug dans le code de build en syntaxe CMake. Le projet GDL utilise les bibliothèques ImageMagick ou GraphicsMagic, la deuxième est plus prioritaire. CMake n'arrive pas à trouver si ImageMagick est installé dans un répertoire différent de celui par défaut. En analysant le code j'ai trouvé que le problème venait des versions récentes de ImageMagick ou le nom du répertoire a été changé. C'est un problème classique, quand les deux logiciels ne sont pas en accord.

Un dernier bug significatif concernait un traitement d'une chaîne de caractères dans GDL. Le traitement marchait jusqu'au premier espace blanc, j'ai ajouté une boucle pour passer toute la chaîne de caractères. Cette modification a amélioré les appels des commandes d'interpréteur de commandes système (par exemple : **cd**, **cp**, **mkdir** ...).

Conclusion

Le travail que j'ai réalisé lors du stage sur le projet GDL sera utilisé dans la prochaine version du logiciel, qui va sortir à l'automne 2013. Ces corrections et améliorations vont aider GDL à devenir une alternative plus solide du logiciel propriétaire IDL (largement utilisé dans le domaine de l'astronomie). Tous les modifications sont mis dans le CVS et peuvent être téléchargées et compilées pour tester avant la sortie d'une version complète. Les utilisateurs ont déjà la possibilité de tester ces nouvelles fonctionnalités (les routines de Cholesky) et de ressentir la performance augmentée (grâce à la définition dynamique du nombre de threads pour OpenMP).

Ce stage au sein de l'Observatoire de Paris m'a été très profitable, et ce sur de nombreux plans. D'une part techniquement, j'ai pu découvrir un nouveau langage de programmation IDL/GDL, qui est un langage interprété. J'ai obtenu l'expérience de codage avec ce langage pendant que j'ai codé les différents tests de régression et de performance. Aussi j'ai approfondi mes connaissances en C++ durant plusieurs missions de correction de bugs et en ajoutant les nouvelles fonctionnalités dans GDL. Un autre apport de ce stage a été d'avoir utilisé plusieurs systèmes d'exploitation, qui m'a permis de me familiariser plus avec les différentes distributions de Linux (CentOS, Debian, Ubuntu). Ces systèmes étaient présents sur les serveurs aussi bien que sur les ordinateurs personnels.

Une des grandes difficultés techniques de ce stage fut d'avoir à modifier un code conçu et souvent modifié par d'autres personnes car il faut faire attention à bien comprendre pourquoi les développeurs précédents ont programmé de cette façon pour ensuite pouvoir apporter les modifications nécessaires. Mais le code qui n'est pas modifié pendant longtemps ne représente pas une logique triviale à comprendre non plus. En ce cas la difficulté est que l'auteur du code peut oublier la raison de choisir cette manière de codage. Dans tous les cas la difficulté du code augmente si le code n'est pas bien commenté et documenté. La documentation d'un logiciel est un aspect très important en informatique, car on a besoin d'utiliser les divers bibliothèques ou juste des petits morceaux du code de l'autre logiciel. Une bonne documentation facilite largement la vie des développeurs.

Les apports personnelles de ce stage sont tout aussi nombreux. Sur le plan professionnel, ce stage m'a apporté de nouvelles connaissances et a augmenté mes capacités de compréhension. Sur le plan humain, l'enrichissement est incontestable puisque j'ai pu développer mon indépendance et travailler en toute autonomie. J'ai par ailleurs pu juger de l'importance du relationnel entre les différents services, mais aussi entre les personnes internes à un service. Cette importance a dévoilé pendant les problèmes techniques (par exemple un problème avec la cable VGA), auxquelles la réaction a été très rapide et tout le monde a montré sa volonté d'aider.

Enfin, la communication fut un point important lors de ce stage. Le travail en équipe fait partie de la vie d'un ingénieur, ainsi, j'ai beaucoup appris de ce stage grâce au fait d'avoir travailler dans un laboratoire en collaboration avec un autre stagiaire. Aussi, le besoin de faire des réunions pour faire le point sur l'avancement du projet s'est très vite fait ressentir afin de vérifier que les idées convergeaient et pour décider des suites à donner aux activités.

Références

- [1] William H. Press, Brian P. Flannery. *Numerical Recipes in C : The Art of Scientific Computing*. Cambridge University Press. 2 edition. October 30, 1992.
- [2] Coulais, A. ; Schellens, M. ; Gales, J. ; Arabas, S. ; Boquien, M. ; Charnial, P. ; Messmer, P. ; Fillmore, D. ; Poplawski, O. ; Maret, S. ; Marchal, G. ; Galmiche, N. ; Mermet, T. *Status of GDL - GNU Data Language*. 434, Astronomical Data Analysis Software and Systems XIX. 2010. <http://www.aspbooks.org/publications/434/187.pdf>
- [3] Coulais, A. ; Schellens, M. ; Arabas, S. ; Lenoir, M. ; Noreksal, L. ; Erard, S. *Space Missions : Long Term Preservation of IDL-based Software using GDL*. 461, Astronomical Data Analysis Software and Systems XXI. 2012. <http://www.aspbooks.org/publications/461/615.pdf>
- [4] *Bilan Social*. L'Observatoire de Paris. Septembre 2010.
- [5] *Rapport d'Évaluation de l'Observatoire de Paris*. Agence d'Évaluation de la Recherche et de l'Enseignement Supérieur. Mars 2010. <http://www.aeres-evaluation.fr/content/download/13342/186416/ffile/aeres-s1-obsparis.pdf>
- [6] Sylwester Arabas, Alain Coulais. *Documentation de GDL*. Marc Schellens and The GDL team. Janvier 3, 2012. <http://gnudatalanguage.sourceforge.net/gdl.pdf>
- [7] *Documentation d'IDL*. <http://www.exelisvis.com/docs/>
- [8] *Documentation d'Eigen*. <http://eigen.tuxfamily.org/dox/>
- [9] *Description de matériels de l'Observatoire de Paris*. <http://dio.obspm.fr>
- [10] *Les divers articles*. <http://wikipedia.org>

Table des figures

- ★ Source de Figure 1 sur la page 6 :
Répartition des effectifs (chercheurs et EC) par composante et par corps
http://ca.obspm.fr/IMG/pdf/09_Bilan-social-2010.pdf
- ★ Source de Figure 2a sur la page 9 :
ARAMIS
<https://sionet.obspm.fr/munin/lerma-a111/aramis/cpu-month.png>
- ★ Source de Figure 2b sur la page 9 :
M2PARIS
<https://sionet.obspm.fr/munin/ufe/m2dsg-pro.obspm.fr/cpu-month.png>
- ★ Source de Figure 2c sur la page 9 :
MEDUSA
<https://sionet.obspm.fr/munin/lerma-b15/medusa/cpu-month.png>
- ★ Source de Figure 3a sur la page 11 :
GCC
<http://upload.wikimedia.org/wikipedia/commons/thumb/a/a9/Gccegg.svg/500px-Gccegg.svg.png>
- ★ Source de Figure 3b sur la page 11 :
GDB
<https://upload.wikimedia.org/wikipedia/commons/6/6a/Archer.jpg>
- ★ Source de Figure 3c sur la page 11 :
VALGRIND
https://upload.wikimedia.org/wikipedia/fr/f/f9/Valgrind_logo.png
- ★ Source de Figure 4 sur la page 14 :
Dépendances de GDL
Figure par Sylwester Arabas
- ★ Source de Figure 5 sur la page 16 :
Logo de GPL version 3
http://upload.wikimedia.org/wikipedia/commons/thumb/9/93/GPLv3_Logo.svg/500px-GPLv3_Logo.svg.png
- ★ Source de Figure 6 sur la page 23 :
Exemple de factorisation de Cholesky
<http://upload.wikimedia.org/math/9/7/3/9733349e9de16e972daca756204d97db.png>

- ★ Source de Figure 8 sur la page 26 :
Exemple d'utilisation de convolution en traitement d'image en IDL.
Cette fonctionnalité n'est pas encore disponible dans GDL
http://www.exelisvis.com/docs/html/images/convol_example2.gif
- ★ Source de Figure 9a sur la page 27 :
Originale
<http://upload.wikimedia.org/wikipedia/commons/5/50/Vd-Orig.png>
- ★ Source de Figure 9b sur la page 27 :
Filtre
<http://upload.wikimedia.org/math/4/e/1/4e13c64b515f5d652e0b18c71d279425.png>
- ★ Source de Figure 9c sur la page 27 :
Après la convolution
<http://upload.wikimedia.org/wikipedia/commons/6/6d/Vd-Edge3.png>

Abréviations

ALMA	Atacama Large Millimeter/sub-millimeter Array
ANTLR	ANother Tool for Language Recognition
BSD	Berkeley Software Distribution
CLI	Command Line Interface
CNAP	Conseil National des Astronomes et Physiciens
CNRS	Centre National de la Recherche Scientifique
CPU	Central Processing Unit
CVS	Concurrent Versions System
DIO	Division Informatique de l'Observatoire
EC	Enseignant-Chercheur
ESA	European Space Agency
ENS	École Normale Supérieure
EPSCP	Établissement Public à caractère Scientifique, Culturel et Professionnel
FITS	Flexible Image Transport System
FSF	Free Software Foundation
HDF	Hierarchical Data Format
IDL	Interactive Data Language
INSU	Institut National des Sciences de l'Univers
GDL	GNU Data Language
GPL	General Public License
GSL	GNU Scientific Library
LDAP	Lightweight Directory Access Protocol
MESR	Ministère de l'Enseignement Supérieur et de la Recherche
MPL	Mozilla Public License
NASA	National Aeronautics and Space Administration
NetCDF	Network Common Data Form
PRAG	PRofesseur AGrégé
SSH	Secure SHell

UMR Unité Mixte de Recherche

VGA Video Graphics Array

XDR External Data Representation

Annexes

A Code de *set_num_threads*

```
1 #ifndef _OPENMP
2 int get_suggested_omp_num_threads() {
3     return 1;
4 }
5 #endif
6
7 #if defined _OPENMP
8 int get_suggested_omp_num_threads() {
9
10     int default_num_threads=1, suggested_num_threads=1;
11
12     char* env_var_c;
13     env_var_c = getenv ("OMP_NUM_THREADS");
14     if(env_var_c)
15     {
16         return atoi(env_var_c);
17     }
18     // cout<<"OMP_NUM_THREADS is not defined"<<endl;
19
20     //set number of threads for appropriate OS
21     int avload, nbofproc=omp_get_num_procs();
22     FILE *iff;
23
24     #if defined(__APPLE__) || defined(__MACH__) || defined(__FreeBSD__)
25         || defined(__NetBSD__) || defined(__OpenBSD__) ||
26         defined(__bsdi__) || defined(__DragonFly__)
27     cout<<"is MAC/*BSD"<<endl;
28     iff= popen("echo $(sysctl -n vm.loadavg|cut -d\" \" -f2)", "r");
29     if (!iff)
30     {
31         return default_num_threads;
32     }
33
34     #elif defined(__linux__) || defined(__gnu_linux__) || defined(linux)
35     iff= popen("cat /proc/loadavg |cut -d\" \" -f2", "r");
36     if (!iff)
37     {
38         return default_num_threads;
39     }
40 }
```

```

37     }
38
39 #elif defined (__unix) || (__unix__)
40     iff=freopen("/proc/loadavg","r",stderr);
41     fclose(stderr);
42     if(!iff)
43     {
44         cout<<"your OS is not supported"<<endl;
45         return default_num_threads;
46     }
47     iff= popen("cat /proc/loadavg 2>/dev/null|cut -d\" \" -f2", "r");
48     if (!iff)
49     {
50         return default_num_threads;
51     }
52
53 #elif defined(_WIN32) || defined(__WIN32__) || defined(__WINDOWS__)
54     iff= popen("wmic cpu get loadpercentage|more +1", "r");
55     if (!iff)
56     {
57         return default_num_threads;
58     }
59     char buffer[4];
60     char* c;
61     c=fgets(buffer, sizeof(buffer), iff);
62     if(!c)
63     {
64         return default_num_threads;
65     }
66     pclose(iff);
67     int cout=0;
68     while(buffer[count]!='\0' && buffer[count]!=' ')count++;
69     for(int i=1,j=1;i<=count;i++,j*=10)
70         avload+=(buffer[count-i]-'0')*j;
71     suggested_num_threads=nbofproc-(int)(avload*((float)nbofproc/100)+0.5);
72     return suggested_num_threads;
73
74 #else
75     cout<<"Can't define your OS"<<endl;
76     return default_num_threads;
77 #endif
78
79     char buffer[4];

```

```

80  char* c;
81  c=fgets(buffer, sizeof(buffer), iff);
82  if(!c)
83      {
84          return default_num_threads;
85      }
86  pclose(iff);
87  avload=(buffer[0]-'0')+((buffer[2]-'0')>5?1:0);
88
89  suggested_num_threads=nbofproc-avload;
90  return suggested_num_threads;
91  }
92  #endif

```

B Fonction de convolution

```

1
2  #ifdef CONVOL_BYTE__
3
4  template<>
5  BaseGDL* Data_<SpDByte>::Convol( BaseGDL* kIn, BaseGDL* scaleIn,
6      BaseGDL* bias, bool center, bool normalize, int edgeMode)
7  {
8      Data_<SpDLong>* kernel = static_cast<Data_<SpDLong>*>( kIn);
9      DLong scale = (*static_cast<Data_<SpDInt>*>( scaleIn))[0];
10     // the result to be returned
11     Data_* res = New( dim, BaseGDL::ZERO);
12     DInt* ker = static_cast<DInt*>( kernel->DataAddr());
13     Data_<SpDLong>* biast=static_cast<Data_<SpDLong>*>( bias);
14     DLong* biasd = static_cast<DLong*>( biast->DataAddr());
15 #else
16 #ifdef CONVOL_UINT__
17
18  template<>
19  BaseGDL* Data_<SpDUInt>::Convol( BaseGDL* kIn, BaseGDL* scaleIn,
20      BaseGDL* bias, bool center, bool normalize, int edgeMode)
21  {
22      Data_* kernel = static_cast<Data_*>( kIn);
23      DLong scale = (*static_cast<Data_<SpDUInt>*>( scaleIn))[0];
24     // the result to be returned

```

```

24 Data_* res = New( dim, BaseGDL::ZERO);
25 Ty* ker = &(*kernel)[0];
26 Data_* biast=static_cast<Data_*>( bias);
27 Ty* biasd = &(*biast)[0];
28 #else
29
30
31 template<class Sp>
32 BaseGDL* Data_<Sp>::Conv( BaseGDL* kIn, BaseGDL* scaleIn,
    BaseGDL* bias, bool center, bool normalize, int edgeMode)
33 {
34 Data_* kernel = static_cast<Data_*>( kIn);
35 Ty scale = (*static_cast<Data_*>( scaleIn))[0];
36 // the result to be returned
37 Data_* res = New( this->dim, BaseGDL::ZERO);
38 Ty* ker = &(*kernel)[0];
39 Data_* biast=static_cast<Data_*>( bias);
40 Ty* biasd = &(*biast)[0];
41 #endif
42 #endif
43
44 if( scale == this->zero) scale = 1;
45
46 SizeT nA = N_Elements();
47 SizeT nK = kernel->N_Elements();
48
49 if(normalize)
50 {
51
52     DDouble tmp=0;
53     for ( SizeT ind=0; ind<nK; ind++ )
54         tmp+=abs(ker[ind]);
55     scale=tmp;
56
57 #ifdef CONVOL_BYTE__
58     tmp=0;
59     for ( SizeT ind=0; ind<nK; ind++ )
60         if(ker[ind]<0)
61             tmp+=abs(ker[ind]);
62     biasd[0]=tmp*255/scale;
63     if( biasd[0]<0)
64         biasd[0]=0;
65     else

```



```

66         if( biasd[0]>255)
67             biasd[0]=255;
68     #else
69     #ifdef CONVOL_UINT__
70         tmp=0;
71         for ( SizeT ind=0; ind<nK; ind++ )
72             if(ker[ind]<0)
73                 tmp+=abs(ker[ind]);
74         biasd[0]=tmp*65535/scale;
75         if( biasd[0]<0)
76             biasd[0]=0;
77         else
78             if( biasd[0]>65535)
79                 biasd[0]=65535;
80     #else
81     biasd[0]=this->null;
82     #endif
83 #endif
84     }
85
86     SizeT nDim = this->Rank(); // number of dimension to run over
87
88     SizeT kStride[MAXRANK+1];
89     kernel->Dim().Stride( kStride, nDim);
90
91     // setup kIxArr[ nDim * nK] the offset array
92     // this handles center
93     long* kIxArr = new long[ nDim * nK];
94     ArrayGuard<long> kIxArrGuard( kIxArr); // guard it
95     for( SizeT k=0; k<nK; ++k)
96     {
97         kIxArr[ k * nDim + 0] = -(k % kernel->Dim( 0));
98         if( center) kIxArr[ k * nDim + 0] = -(kIxArr[ k * nDim + 0] +
99                                     kernel->Dim( 0) / 2);
100     for( SizeT kSp=1; kSp<nDim; ++kSp)
101     {
102         SizeT kDim = kernel->Dim( kSp);
103         if( kDim == 0) kDim = 1;
104         kIxArr[ k * nDim + kSp] = -((k / kStride[kSp]) % kDim);
105         if( center) kIxArr[ k * nDim + kSp] = -(kIxArr[ k * nDim +
106                                     kSp] +
107                                     kDim / 2);
108     }
109 }

```

```

108     }
109
110     SizeT aStride[ MAXRANK + 1];
111     this->dim.Stride( aStride, nDim);
112
113     long aInitIx[ MAXRANK+1];
114     for( SizeT aSp=0; aSp<=nDim; ++aSp) aInitIx[ aSp] = 0;
115
116     bool regArr[ MAXRANK];
117
118     long aBeg[ MAXRANK];
119     long aEnd[ MAXRANK];
120     for( SizeT aSp=0; aSp<nDim; ++aSp)
121     {
122         SizeT kDim = kernel->Dim( aSp);
123         if( kDim == 0) kDim = 1;
124         aBeg[ aSp] = (center) ? kDim/2 : kDim-1; // >=
125         regArr[ aSp] = !aBeg[ aSp];
126         aEnd[ aSp] = (center) ? this->dim[aSp]-(kDim-1)/2 :
            this->dim[aSp]; // <
127     }
128
129     Ty* ddP = &(*this)[0];
130
131     // some loop constants
132     SizeT dim0 = this->dim[0];
133     SizeT aBeg0 = aBeg[0];
134     SizeT aEnd0 = aEnd[0];
135     SizeT dim0_1 = dim0 - 1;
136     SizeT dim0_aEnd0 = dim0 - aEnd[0];
137     SizeT kDim0 = kernel->Dim( 0);
138     SizeT kDim0_nDim = kDim0 * nDim;
139
140     #define INCLUDE_CONVOL_INC_CPP
141
142     if( edgeMode == 0)
143     {
144         #include "convol_inc0.cpp"
145     }
146     else if( edgeMode == 1)
147     {
148         #include "convol_inc1.cpp"
149     }

```

```
150     else if( edgeMode == 2)
151     {
152     #include "convol_inc2.cpp"
153     }
154
155
156 #undef INCLUDE_CONVOL_INC_CPP
157
158
159 if(biasd[0] != this->zero)
160 {
161     for(SizeT indi=0;indi<nA;indi++)
162         (*res)[indi] += biasd[0];
163 }
164
165 return res;
166 }
```
