

Développement de maquettes de solveurs d'écoulements compressibles en Volumes Finis non structurés pour des clusters de GPU TESLA

Jean-Marie LE GOUEZ (ONERA Châtillon, Département Simulation Numérique des Ecoulements et Aéroacoustique) jean-marie.le_gouez@onera.fr

Les maillages non structurés apportent une grande souplesse pour la représentation de domaines d'écoulement de géométrie complexe. Toutefois l'efficacité algorithmique des schémas numériques est handicapée par les accès indirects à la mémoire due au à la connectivité arbitraire entre les cellules de discrétisation voisines.

De manière classique, des techniques de partitionnement permettent de réduire la dimension des zones mémoire auxquelles un thread de calcul particulier doit accéder.

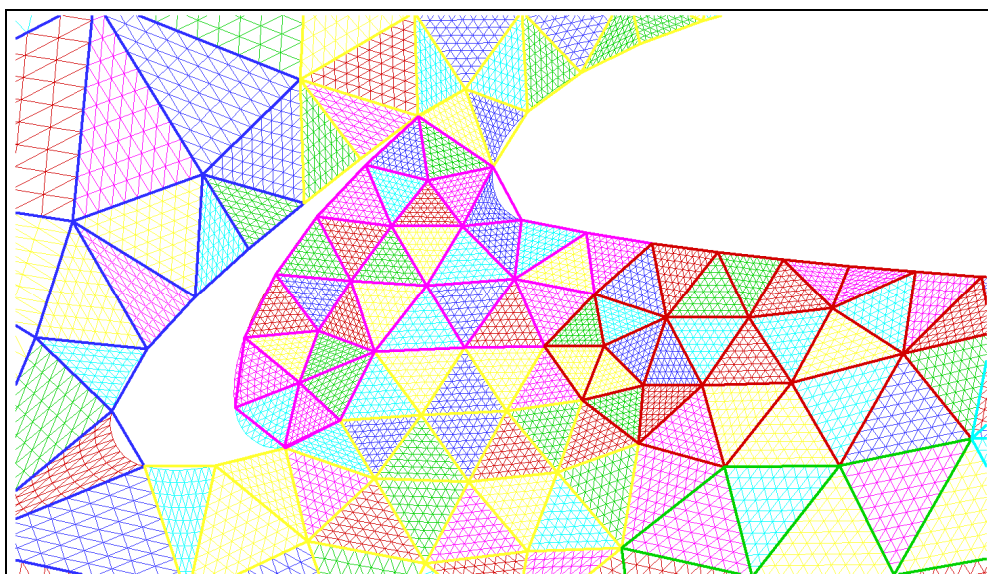
Les GPU TESLA présentent une hiérarchie de mémoires : mémoire globale, mémoire de travail partagée par les threads d'un multiprocesseur. Afin d'exploiter au mieux cette hiérarchie, on a mis en place des techniques de partitionnement de domaine à deux niveaux et un raffinement générique :

- blocs de grande taille destinés à une exécution sur un GPU, avec des échanges entre ces domaines partitionnés gérés par MPI sur le réseau,
- sous-partitionnement fin de blocs de cellules d'une taille typique de 32-64 éléments, destinés à une exécution sur un multiprocesseur, les éléments internes à ces partitions étant connectés par des échanges en mémoire globale de tableaux de variables d'adresses « proches »,
- raffinement générique des cellules de chaque bloc, destiné à imposer une structuration fine des données en paquets adjacents pour permettre des accès coalescents à la mémoire dans toutes les phases des algorithmes : par cellule, par face (calcul des flux), par nœud.

Ces techniques ont été mises en œuvre dans une maquette à base FORTRAN qui adopte le même type de partitionnement hiérarchique dans les structures de haut niveau qu'il manipule, le parallélisme interne au nœud étant géré par des boucles OpenMP sur les blocs issus du sous-partitionnement. Des noyaux CUDA correspondants aux différentes sous-routines de calcul ont été codés, leur exécution est activée depuis le FORTRAN à travers une interface en C qui réalise également le réordonnancement des données internes aux structures pour permettre l'accès coalescent à la mémoire.

Ces méthodes sont destinées aux maillages 2D en triangles et 3D en tétraèdres ou hexaèdres quelconques

Une deuxième maquette, destinée aux configurations 2.5D (typiquement une direction homogène en envergure pour les écoulements aérodynamiques autour de profils élancés, associée à un maillage 2D arbitraire des plans d'envergure) adopte une structuration des données différente, sur la base de vecteurs monodimensionnels périodiques. La programmation FORTRAN (qui reste basée sur MPI et OpenMP) se prête ainsi à une vectorisation poussée, tandis que les structures manipulées par les noyaux Cuda sont organisées selon la direction homogène pour garantir l'accès coalescent à la mémoire globale de chaque GPU.



Sous-partitionnement en blocs et raffinement générique
(inspire par les techniques graphiques de "tessellation" dans le rendu surfacique)