

ppOpen-AT: ポストペタスケール時代の数値シミュレーション基盤ソフトウェア ppOpen-HPC のための自動チューニング基盤

ppOpen-AT: An Auto-tuning Infrastructure for Numerical Simulation Infrastructure ppOpen-HPC Towards To Post Petascale Era

東京大学・情報基盤センター 片桐 孝洋(Takahiro Katagiri)
Information Technology Center,
The University of Tokyo

1. はじめに

先進計算機のハードウェア構成がますます複雑化している。非均質メモリ、多階層キャッシュ、チップ上に多数の計算機要素(コア)をもつマルチコア CPU, そして, GPU に代表される演算アクセラレータ, を搭載した計算機が普及している。ここ数年以内のごく近い将来に, CPU と GPU を混載した計算機環境(非均質(ヘテロジニアス)計算機環境と呼ぶ)が主流になると予想されている。

このように複雑化した計算機環境では, 工学や物理学で必要となる科学技術計算プログラムの最適化が, 高性能計算を専門としない科学者やエンジニアにとって困難になる。この状況は, 応用数理学者にとっても例外ではない。たとえば, 数理的な手法を使ってある手法の収束を 2 倍効率良くしたとしても, その提案手法の演算が非均質計算機環境に不向きであり従来実装の 2 倍演算効率が悪くなれば, 提案手法の実際の効果は無い。したがって, 数値計算法を提案するにも, 計算機環境の特徴を知っておく必要がある。

そこで我々は, 大規模並列計算機上での実用的なシミュレーションコード開発と演算最適化を支援するソフトウェア基盤 ppOpen-HPC を開発している。ここでの“pp”とは, post-peta の略語であり, ペタスケール計算機の後に来る計算機環境のことを指す。ポストペタ環境では, 非均質な計算ノードを持つことが予想されている。つまり, マルチコア CPU と, 演算アクセラレータである GPU の混合構成となることが予想されている。

2. ppOpen-HPC と AT 基盤 ppOpen-AT

2.1 ppOpen-HPC 概要

ppOpen-HPC が対象とするのは, 現在主流の数値計算法によるプログラムである。主に FEM (Finite Element Method), FDM (Finite Difference Method), BEM (Boundary Element Method) および DEM (Discrete Element Method) といった数値計算法を対象としている。

一方, 自動チューニング (Auto-tuning (AT)) 技術が, 先進的計算機環境での性能チューニングで必須となるといわれている。ppOpen-HPC は AT 技術を採用して開発される。ppOpen-HPC のための AT 基盤ソフトウェアを ppOpen-AT と呼ぶ。

図 1 に ppOpen-HPC における ppOpen-AT の位置づけとその機能をのせる。

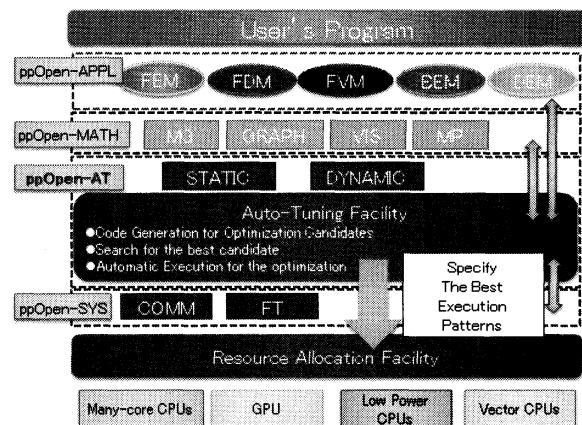


図 1 ppOpen-HPC における ppOpen-AT の位置づけとその機能

2.2 ppOpen-AT 概要

ppOpen-ATは、ppOpen-MATHとppOpen-SYSにAT機能を提供する。ppOpen-ATのAT機能を利用することで、計算機資源選択—ここではCPUとGPUの自動切り替えの機能—を実現する。

図1の計算機資源選択は、ppOpen-ATのAT機能と、ユーザ環境ごとに異なるコンパイラ最適化機能を合わせて使うことにより、ppOpen-AT専用プリプロセッサを通して、コード最適化機能の1つとして実現される。このことで幅広い計算機環境でのATを可能にする。ppOpen-ATの最適化対象は、ppOpen-HPCで実装されている、(1)数値計算法カーネル；(2)プログラム；(3)実装方式；に限定する。主にデータアクセスパターンの最適化を行う。ppOpen-APPLの原始コード(C, C++, Fortran90)が対象であり、コンパイラの「静的最適化」の範疇である。この機能をppOpen-AT/STATICと記載する。将来的には、実行時最適化機能の提供も検討している。実行時最適化機能をppOpen-AT/DYNAMICと記載する。

ppOpen-AT/STATICは以下の機能をもつ：

- (1) オリジナルのC, C++, Fortran90コードから自動生成されたCUDAコード、もしくはOpenCLコードに対する最適化機能；
- (2) DEM, FDM, およびBEMのコードに対し、特定のディレクティブを使うことによる直接のGPUコード生成機能；

2.3 計算機資源選択のための機能

ppOpen-ATが提供する計算機資源選択機能を実現するには、ABCLibScriptのオリジナルなディレクティブ機能[1]に対し拡張が必要である。この拡張は、非均質環境における計算機資源選択を実現するものである。以下の2通りの方法を前提とする：

- (1) ABCLibScript 専用の拡張ディレクティブ“allocate(auto)”を利用する方法[2]；
- (2) ABCLibScript オリジナルのディレクティブ“select region”を利用し、CPUコードとGPUコードの2種を開発者が用意したうえで、コード切り替え機能をAT化する方法；

本稿では、(2)の方法について説明し、予備評価を行う。

3. 予備評価

3.1 計算機環境

本稿で利用した、CPUとGPUの混載環境の概要を以下に述べる。

- CPU 環境：
Xeon W3520 (Nehalem Quad-core, 2.67GHz). 主メモリ：12GB. OS: CentOS5.5 x86_64.
- GPU 環境：
Tesla C2050. CUDA toolkit : version3.2. GPU Driver 260.19.26.
- コンパイラ：
PGI コンパイラ version 11.0-0. OpenMP によるスレッド並列化. PGI アクセレータ機能によりCUDAコードを自動生成する. コンパイラオプション：“-Msafepr=all -fastsse -lm -mp -ta=nvidia -Mlist -Minfo=accel,mp”

3.2 GPUコードについて

本予備評価では、ppOpen-AT/STATICにより選択されるGPUコード(CUDAコード)は、PGIアクセレータ機能により自動生成されるCUDAコードを利用した。なお[2]の方式では、ABCLibScriptがPGIアクセレータのディレクティブを自動生成するので、ユーザがPGIアクセレータのディレクティブを記述しなくてもCUDAコードが自動生成されることに注意する。

PGIアクセレータ機能を使うことなくCUDAコードを開発者が用意すれば、より高性能なGPUコードを利用できる。だがここでの目的は、高性能なGPUコードを自動生成することではなく、CPUコードとGPUコードを自動選択する計算機資源選択のAT機能をppOpen-AT/STATICにより提供し、その有効性を示すことにある。このことで、AT機能付きのソフトウェアを低コスト—ここでは少ない開発工数—で開発ができることを示す。

3.3 問題の詳細

扱うシミュレーションプログラムは、京都大学の岩下武史氏のグループがppOpen-HPCの提供機能の一つとして開発しているコードである。境界要素法(Boundary Element Method, BEM)を利用した静電場の問題である。

上記シミュレーションの問題は、地上から50cmのところ半径25cmの金属球を配置し、金属球に+1Vの電位を与えたときに金属球表面に誘起され

る電荷を求めるものである。空間内の電場、ポテンシャル（電位）の値が分かる。

このシミュレーションでは連立一次方程式を解く。そのため用いられている数値解法は、前処理なしの BiCGStab 法である。連立一次方程式の係数行列は「密行列」として扱われている。いくつかのシミュレーションでは係数行列は疎行列となるが、BEM では境界接触面の影響がほぼすべての要素に影響を及ぼすという特性から、係数行列がほぼ密行列になる。なので、係数行列を密行列として扱うことは BEM では特殊ではない。

このシミュレーションにおける主演算は、BiCGStab 法の中で用いる「密行列・ベクトル積」である。

本問題から作られるテスト用の行列サイズは、以下の 4 種である：

- Small : 600 x 600,
- Medium : 2400 x 2400,
- Medium-Large : 15,000 x 15,000,
- Large : 21,600 x 21,600.

3.4 AT 記述の方法

想定する *ppOpen-AT/STATIC* による AT の実現方法は、CPU コードと GPU コードで異なる。図 2 に、CPU コードにおける「行列・ベクトル積ルーチン」の *ppOpen-AT/STATIC* 記述を載せる。

```
#pragma oat install unroll (row, col) region start
#pragma oat name MyMatVec
#pragma oat varied (row, col) from 1 to 4
#pragma omp parallel for private(col)
for (row=0; row < dim; row++){
    tmp = 0.0;
    for (col=0; col < dim; col++){
        tmp = tmp + mat[row*dim+col] * p[col];
    }
    q[row] = tmp;
}
#pragma oat install unroll (row, col) region end
```

図 2 CPU コードの AT 記述

図 2 では、“#pragma oat” で始まる指示行（プラグマ）が、*ppOpen-AT/STATIC* への AT 指示である。AT のタイミングは「インストール時」(install) であり、ソフトウェアのインストール時に AT を一度行い、その後、その AT 結果を利用した実装選択を実行時に行うことを意味している。

図 2 では、密行列・ベクトル積を構成する 2 重ループのループ導入変数 row, col のそれぞれに 1 段から 4 段までのアンローリングを適用したコードを自動生成する。実行時間を実測の上、最高速となる実装方式を自動選択する。自動生成するコードの種類は、ここでは、4x4=16 種類 (varied (row, col) from 1 to 4) である。

図 2 では、さらに OpenMP のディレクティブが適用されている。OpenMP でスレッド並列化されたうえ、アンローリングを適用したコードになる。ここでは、2 次元配列 mat を 1 次元化して利用している。これは、この環境では 1 次元化したほうが高速であったため、CPU コードでは 1 次元実装を採用した。一方、PGI アクセラタにより CUDA 化した GPU コードは、2 次元配列を用いた実装のほうが、1 次元化した実装より高速であった。そのため、2 次元配列を用いた実装を利用した。図 2 の AT は、このシミュレーションでは、初めに行う AT である。

図 2 の次に行う AT は、BiCGStab 法のメインループに対する CPU コードと GPU コードの切り替えとなる。この AT 記述を、図 3 に示す。

```
#pragma oat select region start
#pragma oat name SelectGPU
#pragma oat select sub region start
    pbicgstab(dim, a, rhs, sol, tor, max_steps);
#pragma oat select sub region end
#pragma oat select sub region start
    pbicgstab_gpu(dim,a, rhs, sol, tor, max_steps);
#pragma oat select sub region end
#pragma oat select region end
```

図 3 CPU と GPU との切り替え AT 記述

図 3 により、CPU コードと GPU コードの時間を実測し、適する方法を選択する AT 機能を有するコードが自動生成される。

3.5 性能評価結果

図 2 と図 3 の AT 機能について、*ppOpen-AT/STATIC* 専用プリプロセッサにより自動生成されるコードについて性能評価した。CPU コードは、8 スレッド並列実行である点に注意する。

最初の AT である、CPU での行列・ベクトル積の AT 効果を図 4 にのせる。図 4 では、8 スレッド並列実行、N=600 (Small) での実行時間を示している。コンパイラ最適化のみで明示的にアンローリン

グをしていない実行時間は $iusw=1$ の時であり、それは 0.095 秒であった。AT により、明示的にアンローリングを施したコードが選ばれた。この場合は、 $iusw=16$ の時であり、これは、外側ループ 4 段、内側ループ 4 段、アンローリングしたコードである。 $iusw=16$ の時の実行時間は、 0.061 秒であった。したがって、AT を搭載することで、約 1.55 倍の速度向上が達成できた。

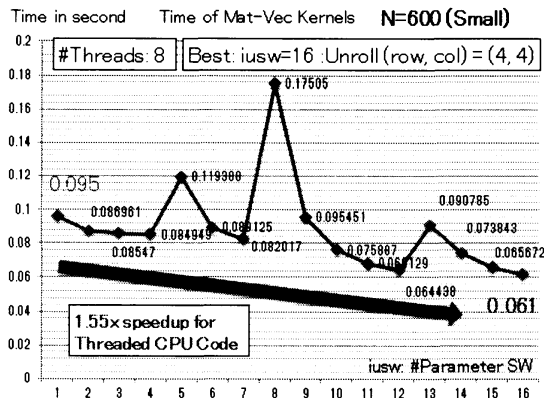


図 4 CPU の行列・ベクトル積の AT 効果 (BEM)

2 番目に行う AT である、CPU と GPU 選択の結果について表 1 に示す。なお、CPU での行列・ベクトル積は図 4 で示した最適化がされている。表 1 の結果は Small, Medium, Medium-Large, および Large とともに最適な実装が選択されている。

表 1 CPU-GPU 切り替え AT 効果 (BEM)

Size	CPU Code [Second]	GPU Code [Second]	Speedup
Small	0.002	0.022	0.09x
Medium	0.128	0.125	1.02x
Medium-Large	7.96	6.99	1.13x
Large	-	Out of Memory	-

表 1 から、問題サイズ Small では CPU 実行が GPU 実行より約 10 倍高速であることがわかる。しかし、問題サイズ Medium 以上になると、GPU コードが高速となることがわかる。したがって、*ppOpen-AT / STATIC* による CPU-GPU 資源切り替え AT は、この環境では有効といえる。

PGI アクセラータ機能による CUDA 化に際し、配列 *mat* の CPU から GPU への転送を最小化している。これは、PGI アクセラータのプラグマにより、

明示的にデータ転送のタイミングを指定することで実現される。したがって、ユーザが PGI アクセラータの機能を熟知していないと実現できない。しかし、*ppOpen-HPC* の利用方法を考えると、*ppOpen-HPC* 開発者であるライブラリ開発者が熟知していればよいわけで、*ppOpen-HPC* を利用するユーザ (エンドユーザ) は PGI コンパイラや AT 機能を一切知る必要はないといえる。すなわち、*ppOpen-HPC* は「ライブラリ」としてエンドユーザから利用されるので、エンドユーザは実装詳細について知る必要はない。

現在の実装は、技術的な問題がある。BiCGStab 法の反復ループ全体を GPU 化できない。この理由は、PGI アクセラータ機能の言語仕様上の制約から生じるものである。CUDA コードを手書きで開発すれば、BiCGStab 法の反復ループ全体を GPU 化できる。したがって、CUDA を用いるプログラミングをしたうえで、*ppOpen-AT* を用いる CPU-GPU の切り替え AT 化を行えば、さらに GPU コードが有効になり、GPU を用いることでさらなる速度向上が得られるはずである。再度の記載になるが、ここでの目的は、PGI アクセラータでの実装や CUDA コードでの手書き実装にかかわらず、CPU-GPU の切り替えを容易に行える *ppOpen-AT* の機能を提供することにある。この点において、本予備評価で有効性を確認できた。

4. おわりに

本原稿では、*ppOpen-AT / STATIC* の概略と、*ppOpen-HPC* に含まれる境界要素法 (BEM) プログラムに対し AT 機能を実装する 1 例を紹介した。

予備評価の結果、BEM で用いられる密行列・ベクトル積に対し、CPU 最適化の 1 例として OpenMP によるスレッド並列化時の AT 例を示した。一方、CPU-GPU の資源切り替えの 1 例として、PGI コンパイラを用いて自動的に CUDA コードを自動生成する方法による GPU コードを、AT 専用言語 *ppOpen-AT / STATIC* を用いて実装する例を示した。双方とも AT による速度向上が確認され、AT 方式の有効性が確認できた。

今後の予定として、*ppOpen-HPC* のプログラムに対して、*ppOpen-AT* の AT 適用事例と、効果的な新規 AT 機能を開発していく。特に BEM の例のような BLAS 基本カーネルの AT に限定するのではなく、FDM に現れる多様で多くの配列アクセスを含むループ

に適用できる AT 方式に適用することを考えている。また、FEM のプログラムにおいて、問題場の物理特性情報、メッシュ結合情報、を活用することで、行列情報だけでは実現できない効率の良いブロック化やアンローリング実装について、新たな AT 機能として開発していく予定である。従来から行っている BLAS レベルの AT ではなく、より上位階層（すなわち、アプリケーションに近い階層）での AT を実現していく。このことが、*ppOpen-HPC* プロジェクトでの AT 方式の新規性である。

関連研究について述べる。CPU-GPU の資源切り替えに関し、ストリーム処理に限定した AT 機能の方式が滝沢ら [3] により提案されている。滝沢らの方法は、ストリーム処理に限定し、少ない試行で GPU 性能を予測するモデルを提案し実装している。本アプローチの AT 方式は、*ppOpen-HPC* の処理に限定している。CPU-GPU 切り替え対象の処理拡大を目指しているため、AT 対象はストリーム処理に限定しない。反面として GPU や CPU の性能モデル化が難しくなる。現在実装されている AT の方法は、ユーザがある条件を与えたうえ、全探索するアプローチ（非モデル化ベース）が採用されている。この方法はコスト高と思われるが、モデル化できないような汎用的な処理に対しても、開発者情報を入れることで探索範囲の限定が実現できる。本プロジェクトでは *ppOpen-HPC* で現れる処理に限定し、現実的な時間で AT を行う AT 方式の実現を目指している。

謝辞: 本研究は、JST CREST 領域「ポストペタスケール高性能計算に資するシステムソフトウェア技術の創出」、H23 年度採択課題「自動チューニング機構を有するアプリケーション開発・実行環境」（代表：中島研吾 東大教授）の支援による。

日頃ご議論いただく *ppOpen-HPC* プロジェクトの諸氏に感謝します。また、BEM コードの提供およびサンプルデータの提供に関し、岩下武史氏（京都大学学術情報メディアセンター）、美船健氏（京都大学大学院工学研究科）、高橋康人氏（同志社大学理工学部電気工学科）の諸氏に感謝いたします。

参 考 文 献

- [1] Takahiro Katagiri, Kenji Kise, Hiroki Honda, and Toshitsugu Yuba, ABCLibScript: A Directive to Support Specification of An Auto-tuning Facility for Numerical Software, *Parallel Computing*, Vol. 32, No.1, pp.92-112

(2006).

- [2] 片桐孝洋, 大島聡史, 平澤将一, 本多弘樹: HxABCLibScript: 非均質計算機向け自動チューニング記述言語拡張, 情報処理学会研究報告-ハイパフォーマンスコンピューティング (HPC), Vol.2011-HPC-129, No.13, pp.1-8 (2011).
- [3] 滝沢寛之, 白取寛貴, 佐藤功人, 小林広明: SPRAT: 実行時自動チューニング機能を備えるストリーム処理記述用言語, 情報処理学会論文誌: ACS, Vol.1, No.2, pp.207-220 (2008).