

課題名 (タイトル) :

全電子計算に基づくタンパク質反応シミュレーションの研究

利用者氏名 :

○木寺 詔紀*
佐藤 文俊*
平野 敏行**
恒川 直樹**
上村 典子**
松田 潤一**

所属 :

*和光研究所 次世代計算科学研究開発プログラム
次世代生命体統合シミュレーション研究推進グループ 分子スケール研究開発チーム
**東京大学生産技術研究所

1. 本課題の研究の背景、目的、関係するプロジェクトとの関係

本課題は「次世代生命体統合シミュレーションソフトウェアの研究開発」プロジェクト、分子スケール研究開発チームに参画中の、全電子計算に基づくタンパク質反応シミュレーションプログラム ProteinDF のチューニングのために実施された研究である。次世代スーパーコンピュータは 8 万ノード・64 万コア以上の計算ノードから構成されており、高度に並列化されたプログラムが必要とされる。RICC は国内屈指の約 8000 並列のプログラムが実行可能な環境である。RICC を用いて ProteinDF が次世代スーパーコンピュータ上で効率的に並列計算できるように開発・テストすることを目的とした。

2. 具体的な利用内容、計算方法

ProteinDF はタンパク質の全電子計算を遂行するために、全ノードから参照できる共有ファイルシステムを有効利用するように構築されていた。大規模分子の量子化学計算を行う場合、行列要素の計算ならびに行列演算に大量のメモリ領域を必要とする。例えば、1000 残基規模のタンパク質の全電子量子化学計算を行うためには、一度に最低約 400 GB (80 GB×5) のメモリ領域を必要とする。このため、ProteinDF では同時に計算する行列を最低限に絞り、その時点で必要ではない行列はディスク領域に退避させる方法を採用している。これは非効率的であるが、計算ノードの搭載メモリ量が限られている現在において大規模分子の量子化学計算を行う次善の策といえる。ディスク領域に退避させるもう一つのメリットとして、リスタートが可能であ

る点も挙げられる。実際これにより、我々は世界最大の全電子量子化学計算を達成してきた。

一方 2009 年 12 月 25 日に公表された次世代スーパーコンピュータの仕様によると、利用者・プログラムは各ノードブロックに搭載されたローカルファイルシステムを利用する必要がある。ProteinDF は共有ファイルシステムを前提にコーディングされていたため、大規模行列をメモリ上に分散保持ならびにローカルファイルシステムへ分散保存、さらに分散された行列の参照・計算・格納を行うように書き換えが必要である。

また、次世代スーパーコンピュータでの 64 万並列を目標とした場合、MPI/OpenMP のハイブリッド並列は必須である。無論、効率的に並列計算を行うには、計算粒度に応じた MPI/OpenMP 比を求める必要がある。本年度は実行時にハイブリッド並列のバランスを変えられるように工夫を施した。

3. 結果

巨大行列を分散保持し、それを各ノードが参照・計算・格納するアプローチを現在 ProteinDF に組み込み作業中である。10000 並列(超)規模で効果的な演算を行うための、行列演算に適した行列分散保持方法と、分子積分に適した行列分散保持方法は異なっているので、それぞれの計算に適した保持方法に組み替えることが望ましいと考えている。現在の ProteinDF は SCF 繰り返し計算において、direct SCF 法と呼ばれる、差電子密度行列を用いて計算すべき分子積分の数を絞り込む手法を利用している。密度行列の性質を利用して上手に分散保持し、SCF 計算を加速する手法を現在コーディング中である。

ハイブリッド並列化はほぼ完了しており、細部のチューニングを行っている段階である。for ループ内の計算強度がほぼ一定な箇所の並列化は、static に、すなわち領域分割によって並列計算を行うように OpenMP 指示文を挿入している。一方、for ループ内部の計算強度の異なる計算箇所については、環境変数 OMP_SCHEDULE により実行時にユーザーが並列計算手法を選択できるようにコーディングしている。これにより、少ないスレッド数の場合は領域分割による並列計算、多いスレッド数の場合は動的にタスクを割り当てる並列計算を行うように指示することができる。次世代スーパーコンピュータでは、各ノードで OpenMP による 8 スレッド並列を行うことができるので、動的にタスク分割をする並列計算方法の方が高速に動作すると予想される。

このようにプログラムで可能な限りの工夫を施したものの、ハイブリッド並列の効率は、OpenMP に関わる C++コンパイラの出来によるところが大きい。同一コード、同一環境においてテストした場合、GNU C++コンパイラでビルドされた実行バイナリと、intel 社製コンパイラのそれを比較すると、40%近く計算時間に差が生まれることがある。ProteinDF は標準 C++で記述されているため、GNU C++コンパイラ、intel 社製 C++コンパイラ、富士通社製 C++コンパイラで正常にビルドし、動作可能である。次世代スーパーコンピュータの場合、各ノードで OpenMP による 8 スレッド並列が行われることになるため、次世代スーパーコンピュータにおける ProteinDF の並列性能は C++コンパイラの性能に寄与するところが多い。

4. まとめ

次世代スーパーコンピュータ上で効率良く ProteinDF を実行できるようにするため、巨大行列分散保持・計算方式の改良とハイブリッド並列の性能向上を行った。巨大行列を分散保持した場合の分子積分ルーチンは現在改良作業中である。

5. 今後の計画・展望

現在、作業中の巨大行列の分散保持による分子積分計算ルーチンを組み込み、パフォーマンステストを行う。問題の大きさや並列数を変化させることによりスケラビリティを測定し、改善を行う。これにより、次世代スーパーコンピュータに適した ProteinDF が完成し、前人未踏の大規模タンパク質量子化学計算が達成出来るであろう。

6. RICC の継続利用を希望の場合は、これまで利用した状況（どの程度研究が進んだか、研究においてどこまで計算出来て、何が出来ていないか）や、継続して利用する際に行う具体的な内容

RICC では次世代スーパーコンピュータを想定した ProteinDF の実行・利用が可能であるため、ローカルファイルシステムの利用方法も含め、効率的な利用方法を探る。8000 コアを使用した並列計算を実施し、スケラビリティの結果を外挿することにより、次世代スーパーコンピュータの 64 万コアでもパフォーマンスが低下しないようにコーディングを行う。

7. 一般利用で演算時間を使い切れなかった理由
該当なし。

8. 利用研究成果が無かった場合の理由

本課題は 2010 年 1 月から出発しており、これを利用した研究報告は行っていない。