

SR11000 を使い始めた人のための コンパイラのちょっとしたコツ

佐藤周行

スーパーコンピューティング研究部門

概要

この原稿の（実質的な点は）次にまとめられます。

- いじるのはプログラムではなくてコンパイラオプションのほう。
- コンパイラの頭はあまりよくないことは事実だが、人間ほどではない。
- プログラムは素直に書くほうが結局は速い（下手の考え休むに似たり）。

1 はじめに

この原稿は、SR11000 を使い始めた利用者が、東大センターでの標準コンパイラである日立のコンパイラをどう利用するかについて、少しノウハウを伝えられたらなあというものです。今回は Fortran を取り上げますが、C でも状況は（基本的に）変わりません。

ん？わかってますって。常日頃から、どういじってもキャッシュの再利用なんか土台無理なプログラムの改良に苦労しているのに、行列積なんていうやさしい問題を例にとって「ほら、こんなにうまくいくでしょう」なんていうのは詐欺だって。そういう人たちにとって、この稿はあまり参考にならないかもしれません。でも、コンパイラオプションが何を制御して、どのように性能に影響するかをながめるには少しは参考になるかもしれません。また、性能を上げようと思ったらプログラムをいじらないといけない GNU のコンパイラと違って、「コンパイラオプションをいじる」ことで性能を上げるのが日立コンパイラのコツだということを理解する良い機会かもしれません。

また、この稿は性能向上の王道であるプロファイリングについてふれていません。大きな問題ですから、稿をあらためてふれることにします。

2 日立 Fortran コンパイラ

ここで相手に選ぶのは日立最適化 Fortran90 コンパイラです。名前のとおりいわゆる Fortran90 とよばれる規格に準拠していますが、Fortran95 にも対応しています¹。このコンパイラは、命令スケジューリングやキャッシュ、プリフェッチなどの CPU に特化した最適化を行なうのはもちろん、1980 年代から連綿と続いてきたループの再構成による最適化（当時はベクトル化を対象としていましたが）の並列アーキテクチャやキャッシュアーキテクチャ

¹ ちなみに最新の規格は Fortran2003 です。これを実装したコンパイラは世界中探してもまだありません。規格が決まって（2004 年）から 3 年目だというのに...

への適用が大々的に行なわれています。また、高速ライブラリを用意し、イディオム認識によってあるコード部分をそのままごっそり置換するような最適化も平気で行ないます。

並列化に関する能力を見るのは別の機会に譲り、今回は、基本的な最適化の性能を見てみることにしましょう。

2.1 最適化オプション

最適化オプションは基本的に以下の3つです。マニュアルに書いてある説明をつけておきます。

-03	プログラム全体で大域的な最適化をする
-04	-03に加えて制御構造の変換，演算順序の変更などを含め，プログラム全体を最適化する。
-0s	実行速度が速くなるように，次に示すコンパイルオプションを自動的に設定する。 -nolcheck -nodochk -approx -noconvcheck -disbracket -divmove -expmove -fma -invariant_if -ischedule=3 -loopdistribute -loopexpand -loopfuse -loopinterchange -loopreroll -04 -parallel=2 -prefetch -prod -rapidcall -scope -swpl -noargchk -nobreak -noerstmt -noagochk

そのほかに最強の最適化オプション-0ss というのもありますが，今回は使いません。最適化オプションを指定しなければ -03 が適用されます。最適化を一切抑止する-00 というオプションもありますが，特別な理由がないかぎり使うことはないでしょう。

ある特定の最適化オプションを抑止したい場合はnoを先頭につけます。たとえば-loopfuseを無効にしたい場合は-noloopfuseとします。

筆者の感触では，-03 はコンパイラの教科書² に書いてあるような最適化がなされます。-04 でもいわゆる「普通」にがんばっているコンパイラなら実装されているようなものです（それでも出てくるコードはだいぶ上質です）。

日立コンパイラの日立コンパイラたる所以は-0sの最適化の凄さにあるといっても過言ではないでしょう。日立（に限らず，いわゆる国産3社）は昔からコンパイラを作ってきて，ベクトルコンピュータの時代は製品としては世界一の品質を誇ってきたわけですが，その技術の歴史が全部詰め込まれている感じがします。

ここでは，ループの再構成（交換，分散，融合，タイリングなど）³ をはじめとした，いわゆる「ソースプログラムの書き換え」レベルの最適化が適用されます。どれをどのくらい適用するかについてはコンパイラが判断します。そういうわけで

1. いじるのはプログラムではなくてコンパイラオプションのほう。

一般論ですが，この「判断」はコンパイラに任せるのが賢明です。コンパイラも時々判断ミスをすることがありますが，

2. コンパイラの頭はあまりよくないことは事実だが、人間ほどではない。

² これにはまるとコンピュータサイエンティストになっちゃいますから，計算科学者は読んではいけません（？）

³ ちなみに，ループ最適化として皆がすぐ思いつくループアンローリングは再構成には一般に含めません。適用される理論が違うからです。

2.2 最適化のリストの表示

ある特定の最適化オプションを指定したときに、どの最適化がどう適用されるか実感がわからないかもしれません。そういう時は`-loglist`を指定しましょう。たとえば次のようにします。

```
% f90 -O0s -nodisbracket -noparallel -loglist m.f
...
```

上の例では、リストが`m.log`に出力されます。出力の形式はたとえば次のようになります。
**に続いて、コメントが挿入されます。

```
subroutine mmjki(a,b,c,n)
  implicit none
  integer n,i,j,k
  double precision a(n,n),b(n,n),c(n,n)
**
** Outer loop unrolled (2 times).
**
  do j = 1,n
**
** IF test is invariant in loop so moved to outside.
** Outer loop unrolled (4 times).
**
    do k = 1,n
**
** SWPL applied.
** 6 streams (A[]x4, C[]x2) pre-fetch applied.
**
      do i = 1,n
        c(i,j) = c(i,j) + a(i,k) * b(k,j)
      end do
    end do
  end do
end
```

もっと詳細に見たいという方にはアセンブリコードを直接見るのが一つの方法です。オプションは`-S`。

```
% f90 -O0s -S -loglist m.f
...
```

出力は上の場合、`m.s`に出てきます。アセンブリコードを読む必要のある場合が存在することは否定しませんが、これにはまると、生産性その他大きな問題が生じます。最後の手段にしておきましょう。

3 実験

与太話をいくらしても実感がわきませんから、具体的に最適化オプションをいろいろ変えてそれが性能にどのくらい関係するかを実感してみましょう。

実験に使うのは行列積です。行列積はLU分解とともに、最適化の余地が大きく、かつプログラムの構造が単純で、数値計算の素養がない人にも理解しやすいので、最適化の効果を見せるための格好のターゲットになってきました⁴。

3.1 実験環境

実験は2007年7月17日に行ないました。日立 Fortran コンパイラのバージョンは V01-05-/C で、Interactive 環境で実行しました。

interactive 環境では、他のユーザと CPU をシェアしていますから、いろいろ邪魔が入ります。センターのバッチ環境では、OS の割り込みが入る可能性はありますが、ほぼ安定した環境でプログラムが実行できると思ってよいでしょう。その意味で、interactive 環境で行なった今回の実験は少し割り引いて考えてください。

3.2 実験プログラム

実験プログラムとしていずれも行列積を計算する以下のものを用意しました。表1にあげておきます。

解説

IJK は、定義に素直にプログラムを書いたものです。JKI は、データアクセスの足跡を考えれば、一番効率的だろうと思われているものです⁵。IJKT は、行列積の改善のときに必ず問題になる A のデータの再利用を、ループの入れ替えではなくて、必要なデータをあらかじめ集めることで改善しました。IJKTT は、それを一歩進めて、行列を転置してしまったものです。作業用に余分な配列が必要になります。最後に INTRINSIC は Fortran90 の組み込み関数 `matmul` を使いました。

3.3 実験結果

行列のサイズを 200 と 2000 にし、最適化オプションを 4 種類用意しました。結果を MFLOPS を単位として表 2 に書き出します。INTRINSIC については、一度 $2 \cdot n^3$ の掛け算をしたあとに n^2 の足し算をしているようにみえますが、やっていることは他のルーチンと同じですから、これも計算量は $2 \cdot n^3$ とすることにします。

⁴ もちろん、現実のプログラムの中に頻繁に現れることも理由のひとつです。

⁵ なぜ？この説明を講義その他で今までされたことがない人は、これを確認しましょう。ついでに、データの格納方法についての C と Fortran の違いを一度確認しておきましょう。ちなみに C では最適なパターンが IKJ になります。column major と row major というのがキーワードです。

表 1: 実験に使ったプログラム

IJK:

```
do i = 1,n
  do j = 1,n
    do k = 1,n
      c(i,j)=c(i,j)+a(i,k)*b(k,j)
    end do
  end do
end do
```

JKI:

```
do j = 1,n
  do k = 1,n
    do i = 1,n
      c(i,j)=c(i,j)+a(i,k)*b(k,j)
    end do
  end do
end do
```

IJKT:

```
do i = 1,n
  d = a(i,:)
  do j = 1,n
    do k = 1,n
      c(i,j)=c(i,j)+d(k)*b(k,j)
    end do
  end do
end do
```

IJKTT:

```
d = transpose(a)
do i = 1,n
  do j = 1,n
    do k = 1,n
      c(i,j)=c(i,j)+d(k,i)*b(k,j)
    end do
  end do
end do
```

INTRINSIC:

```
c = c+matmul(a,b)
```

表 2: 実験結果

n=200: (MFLOPS)

オプション	IJK	JKI	IJKT	IJKTT	Intrinsic
-03 -noparallel	758.8	1250.5	757.0	748.1	216.7
-04 -noparallel	790.1	2344.6	788.8	778.7	755.1
-0s -nodisbracket -noparallel	5110.5	5123.4	2818.2	5062.7	1117.9
-0s -noparallel	6940.7	6949.6	4305.7	5065.9	6147.3

n=2000: (MFLOPS)

オプション	IJK	JKI	IJKT	IJKTT	Intrinsic
-03 -noparallel	141.1	902.1	607.5	622.9	134.4
-04 -noparallel	199.4	1243.0	628.5	641.6	198.6
-0s -nodisbracket -noparallel	3333.1	3333.3	1480.1	3834.5	295.5
-0s -noparallel	6116.5	6116.5	1660.0	3819.4	5998.1

3.4 観察

まず観察すべきは、最適化オプションで性能が劇的に異なるということです。キャッシュが効く（サイズ 200）と、一桁ですみますが⁶、二桁になるとそれは大問題でしょう。

さてと、行列のサイズが 200 ではキャッシュの効果が如実に現れますからデータアクセスのコストをあまり考えない最適化で大丈夫です。ただ、2000 になるとキャッシュがあふれてしまいますから、いわゆる「キャッシュ最適化」をきちんとしないと性能が出ません⁷。

さて、見るべきポイントは 200 と 2000 の差だけではありません。つっこみどころは実は満載なんですけど、ここではポイントを 3 つにしぼりましょう。

1. データアクセスのコストを考えた最適化が必須で、それは最適化レベルを上げないときには有効。
2. IJK と JKI の差は、最適化のレベルを上げると消えてしまう。
3. -0s にすると、結局同じコードになってしまう。

-03 では、m.log には何も出てきません。-04 になるとソフトウェアパイプラインが適用されるというメッセージが出てきます。この効果は数割にとどまるというのは表を見ればよくわかります。-0s -nodisbracket にすると、日立コンパイラの頭の中がだんだん透けて見えるようになってきます。リストで確認してみましょう。

```
subroutine mmijk(a,b,c,n)
  implicit none
  integer n,i,j,k
  double precision a(n,n),b(n,n),c(n,n)
**
** SWPL applied.
** 6 streams (A[]x4, C[]x2) pre-fetch applied.
**
  do i = 1,n
**
** Loop interchanging (I,J,K)->(J,K,I) applied.
** Outer loop unrolled (2 times).
**
    do j = 1,n
**
** IF test is invariant in loop so moved to outside.
** Outer loop unrolled (4 times).
**
      do k = 1,n
```

⁶ 原稿の都合上、こう書きましたが、性能が 7 倍くらいちがうと、1 日で終わるジョブが 1 週間流れるということになります。この感覚は無視すべきではないと考えます。

⁷ だからわかってますって。みなさんのプログラムがキャッシュ最適化の恩恵を必ずしも受けない、受けようと思ったらアルゴリズムからの再検討が必要になるということは。キャッシュ最適化がすべてを解決するような幻想を与えるつもりはありません。この意味で、今もってベクトル計算機を作り続けている会社、またその利用者からの恒常的な批判は十分意識する必要があります。

```

        c(i,j) = c(i,j) + a(i,k) * b(k,j)
    end do
end do
end do
end do

```

このリストと、`-loglist` の例として出したリストの違いがループの交換だけというのについて確認しておきましょう。

一般に-03 や-04 では、メモリ、キャッシュ、レジスタの階層を持つデータアクセスのコストを考慮した最適化がなされないようです。性能を出そうと思ったら-0s が必須です。これに背を向けてたとえば IJKT とか IJKTТ などのように自分でデータアクセスを最適化する手もあります。これもある程度効果がありますが、結局、-0s -nodisbracket 以上になると投資の効果が消えてくるということでしょうか。

投資の効果が消えるということでは、IJK と JKI の差の解消も見逃せません。ちょっと前に「JKI がなぜ速いか理解せよ」と書きましたが、日立コンパイラはこころへんも解析して最適のデータアクセスのパターンを選択してくれます。このご利益が実感できるのも-0s 以上になります。

ちなみに-nodisbracket の有無で出てくるコードに差があるのは、日立の親切によるもので演算順序の違いを（規格上無視するのが許されているのに）利用者のコードに書かれていた意思を尊重することによるものです。-nodisbracket をつけずに-0s にすると、問答無用で日立の開発した高速行列積ライブラリが呼び出されます。これをリストで観察しましょう。

-0s でのリスト:

```

**
**   行列積ライブラリ化 (_hf_mat_matmul).
**
    do j = 1,n
        do k = 1,n
            do i = 1,n
                c(i,j) = c(i,j) + a(i,k) * b(k,j)
            end do
        end do
    end do
end do

```

本当は 5 種類のコードが全部置き換えられなければならないのですが⁸、それでも IJK, JKI, Intrinsic が同じコードに落ち着きました⁹。問題は、とても速いコードに置き換えるためには、それなりにプログラムをシンプルに保たなくてははいけないということです。IJKT や IJKTТ のように、「変な」工夫をすると、コンパイラがその意図を見つけきれないことがおきます。そこで

3. プログラムは素直に書くほうが結局は速い(下手の考え休むに似たり)。

ここではふれませんでした。並列化でも同じことが言えます。プログラムはできるだけシンプルにして保守するのがお勧めです。

⁸ こころへんができるという論文を昔書いたことがあるのですが、まあいいや。

⁹ ただし、Intrinsic については、イディオム認識にもう少し改良の余地がありそうです。

4 おわりに

ここでは、行列積を例にとって日立のコンパイラがどのように最適化を行なっているかについて解説をしました。まとめは概要に書いた3点に集約されます。繰り返します。

- いじるのはプログラムではなくてコンパイラオプションのほう。
- コンパイラの頭はあまりよくないことは事実だが、人間ほどではない。
- プログラムは素直に書くほうが結局は速い(下手の考え休むに似たり)。

もっともこれはコンパイラの研究者にとってとても都合の良い結論ですから、計算科学者はあまり信用してはいけません。ただし、コンピュータサイエンティストから何かアドバイスがあったときにそれを受容するゆとりを持っていたいただければ、チューニングがさらに進むと思います。