富士通詳細 プロファイラの利用事例

2022-12-14 17th HPC-Phys勉強会 R-CCSソフトウェア開発技術ユニット 似鳥啓吾

内容

- 詳細プロファイラfappの使い方
 - 雰囲気で使っているので非公式です、あくまでも1ユーザーの体験と 感想として
- Nek5000/RS axhelmカーネルのチューニング
 - R-CCS辻さんの仕事の再現実験+文書化
 - プロファイラを取ったのはチューニングの後
- Bridge++ QXS Wilsonカーネルのチューニング
 - 金森さんからの宿題
 - チューニング前にプロファイラを取った例

参考文献

- 「富岳」もしくはFX1000ユーザー限定となります
- ポータル->マニュアル->言語マニュアル
 - プロファイラ使用手引書
- ポータル->プログラミングガイド
 - プログラミングガイド(チューニング編)
 - プログラミングガイド(プロセッサ編)

基本プロファイラと詳細プロファイラ

- 基本プロファイラ: fipp: Fujitsu Instant Performance Profiler
 - 割り込みによるサンプリングベース
 - 関数ごとの消費時間など
 - 一部CPU動作状況(Gflopsやメモリスループット)
- 詳細プロファイラ: fapp: Fujitsu Advanced Performance Profiler
 - ハードウェアカウンターベース
 - 最大17回同じプログラムを繰り返し実行
 - CSV形式に変換後Excelマクロで可視化できる

• 共通事項

- 実行時のDLL切り替えにより通常実行とプロファイリングを切り替え
- ソースコードにstart/stopの埋め込みが必要だが、プロファイリング専用のbranchを作らなくても済む
- ただし富士通以外の環境ともソースを共通化するにはマクロ関数でAPIを 消せるようにする必要がある

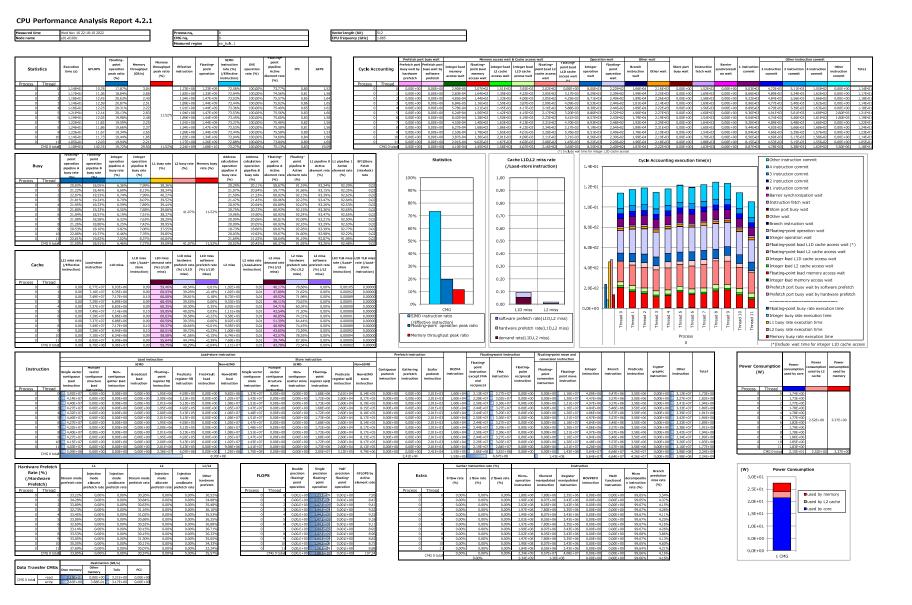
ソースコードの編集 (fapp)

- name, number, levelの順で指定
 - name+numberで区間を識別
 - levelは0なら常に計測、実行時にオプションで対象を絞れる
- fapp_start/fapp_stopのネストも可能

実行スクリプトの例

```
for i in {1..17}
do
fapp -C -d ./rep${i} -Hevent=pa${i},method=fast mpiexec ${BIN} ${Px} ${Py} ${Pz} ${Nx} ${Ny} ${Nz} ${Nt} ${Nthreads} ${Iter}
done
for i in {1..17}
do
fapp -A -d ./rep${i} -Icpupa,nompi -tcsv -o pa${i}.csv
done
```

- ホストでcsvファイルを作成するときはfapppxコマンド
- -Hmethod=fastだとOSを経由しないのでオーバーヘッドが小 さい
- 後々結果をscpする用に/optからxlsmファイルへのリンクを 作っておくとよい



• 出力例:50インチぐらいの4K TVで見るとよい

Nek5000/RS axhelmカーネル

- スペクトル要素方
 - 六面体要素の中に8³=512個の点がある
 - 各軸の8点は選点(直行多項式の根)になっている
 - 各軸に8×8行列を掛けることで勾配が求まる (7次の多項式になっていると思えばわかる)
- やっている計算

$$A = egin{pmatrix} D_1 \ D_2 \ D_3 \end{pmatrix}^T egin{pmatrix} G_{11} & G_{12} & G_{13} \ G_{12} & G_{22} & G_{23} \ G_{13} & G_{23} & G_{33} \end{pmatrix} egin{pmatrix} D_1 \ D_2 \ D_3 \end{pmatrix} \ (
abla ec{v},
abla ec{u}) = ec{v}^T A ec{u} \ egin{pmatrix} V ec{v},
abla ec{v} ec{v} \end{pmatrix}$$

- 勾配、幾何係数、勾配の転地の順に掛ける
- 幾何係数: 六面体の中の座標での微分を外の座標での微分に直す (+ ヤコビアンと密度で重み付け)
- 勾配は演算ネック、幾何係数はメモリネック(6語読んで9乗算)

As-is & tune-1

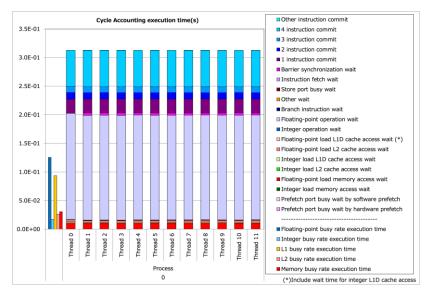
```
79
           for(int k = 0; k ; <math>++k) for(int j = 0; j ; <math>++j) // Loop-2
81
               for(int i = 0; i 
82
                   // const dlong id = i + j * p_Nq + k * p_Nq * p_Nq + element * p_Np;
83
                   const dlong gbase = element * p_Nggeo * p_Np + k * p_Nq * p_Nq + j * p_Nq + i;
                   const dfloat r_G00 = ggeo[gbase + p_G00ID * p_Np]; // #define p_G00ID 1
                   const dfloat r G01 = ggeo[gbase + p G01ID * p Np]; // #define p G01ID 2
85
                   const dfloat r_G11 = ggeo[gbase + p_G11ID * p_Np]; // #define p_G11ID 3
                   const dfloat r_G12 = ggeo[gbase + p_G12ID * p_Np]; // #define p_G12ID 4
                   const dfloat r_G02 = ggeo[gbase + p_G02ID * p_Np]; // #define p_G02ID 5
                   const dfloat r_G22 = ggeo[gbase + p_G22ID * p_Np]; // #define p_G22ID 6
90
91
                   dfloat gr = 0.f;
                   dfloat qs = 0.f;
93
                   dfloat qt = 0.f;
94
                   // Loop-2a (gradient)
95
                   for(int m = 0; m < p_Nq; ++m) {</pre>
96
                       qr += s_D[i][m] * s_q[k][j][m];
97
                       qs += s_D[j][m] * s_q[k][m][i];
98
                       qt += s_D[k][m] * s_q[m][j][i];
99
                   // Loop-2b (geometry factor)
                   s Gqr[k][j][i] = r G00 * qr + r G01 * qs + r G02 * qt;
.01
02
                   s_{gs[k][j][i]} = r_{go1} * qr + r_{go1} * qs + r_{go1} * qt;
.03
                   s_{qt}[k][j][i] = r_{q02} * qr + r_{q12} * qs + r_{q22} * qt;
```

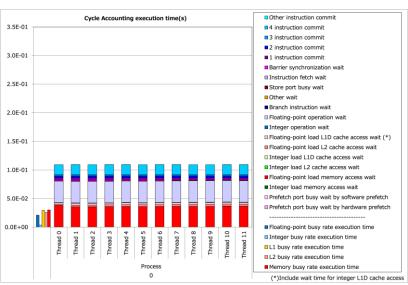
- 最内ループ変数mに対するSIMD化では不適切である
 - Gatherと水平加算の発生
 - 幾何係数の部分が一切 SIMD化されない

```
jwd82220-i "axhelm-0.cpp", line 78: このループで必要なプリフェッチの数が、ハードウェアプリフェッチの許容ま
jwd86620-i "axhelm-0.cpp", line 78: スケジューリング結果を得られなかったため、ソフトウェアパイプライニング
jwd82030-i "axhelm-0.cpp", line 80: このループをフルアンローリングしました。
jwd6004s-i "axhelm-0.cpp", line 94: リダクション演算を含むループ制御変数'm'のループをSIMD化しました。
jwd82090-i "axhelm-0.cpp", line 100: 多項式の演算順序を変更しました。
jwd86680-i "axhelm-0.cpp", line 106: 命令数が多いため、ソフトウェアパイプライニングを抑止しました。
jwd82030-i "axhelm-0.cpp", line 106: このループをフルアンローリングしました。
jwd82030-i "axhelm-0.cpp", line 108: このループをフルアンローリングしました。
jwd6004s-i "axhelm-0.cpp", line 112: リダクション演算を含むループ制御変数'm'のループをSIMD化しました。
```

ディレクティブでm-loopをアン ロールするとひとつ外側のループ変 数iに対してSIMD化される

Before/after





- 効果は抜群だが、、、
- 左側のPA情報からこのチューニングに到れるかどうかが問題
- どのループ変数でSIMD化したかはコンパイラのメッセージからわかる
 - mでは不適切でiならよいと気付くことができるか?
- ディレクティブは知らないかもしれない
 - mとiのループ順序を入れ替えることでも対処は可能

Tune-2 (ストライドアクセスの回避)

```
for(int m = 0; m < p_Nq; ++m) {
    qr += s_D[i][m] * s_q[k][j][m];
    qs += s_D[j][m] * s_q[k][m][i];
    qt += s_D[k][m] * s_q[m][j][i];
}</pre>
```

8×8の行列s_Dは要素 に非依存

		⊢				000-000	285000000000000000000000000000000000000		Load-store							
			Load instruction													
e		ᆫ	SIMD													
Instruction Process Thread			gle vector ontiguous load ostruction	Multiple vector contiguous structure load instruction	Non- contiguous gather load instruction	Broadcast load instruction	Floating- point register fill instruction	Predicate register fill instruction	First-fault load instruction							
0	0		2.10E+07	0.00E+0	1.02E+06	3.08E+06	2.34E+07	1.70E+01	0.00E+00							
0	1		2.10E+07	0.00E+00	1.02E+06	3.08E+06	2.34E+07	1.70E+01	0.00E+00							
0	2		2.10E+07	0.00E+00	1.02E+06	3.08E+06	2.34E+07	1.70E+01	0.00E+00							
0	3		2.10E+07	0.00E+0	1.02E+06	3.08E+06	2.34E+07	1.70E+01	0.00E+00							
0	4		2.10E+07	0.00E+0	1.02E+06	3.08E+06	2.34E+07	1.70E+01	0.00E+00							
0	5		2.10E+07	0.00E+0	1.02E+06	3.08E+06	2.34E+07	1.70E+01	0.00E+00							
0	6		2.10E+07	0.00E+0	1.02E+06	3.08E+06	2.34E+07	1.70E+01	0.00E+00							
0	7		2.10E+07	0.00E+00	1.02E+06	3.08E+06	2.34E+07	1.70E+01	0.00E+00							
0	8		2.10E+07	0.00E+0	1.02E+06	3.08E+06	2.34E+07	1.70E+01	0.00E+00							
0	9		2.10E+07	0.00E+0	1.02E+06	3.08E+06	2.34E+07	1.70E+01	0.00E+00							
0	10		2.10E+07	0.00E+00	1.02E+06	3.08E+06	2.34E+07	1.70E+01	0.00E+00							
0	11		2.10E+07	0.00E+0	1.02E+06	3.08E+06	2.34E+07	1.70E+01	0.00E+00							
	CMG 0 total		2.52E+08	0.00E+0	1.23E+07	3.69E+07	2.81E+08	2.04E+02								
	Ci io o total								8.76							

Gather/scatterの発生 はプロファイラで検出 可能

関数の入り口で転置行列s_Tを作っておく

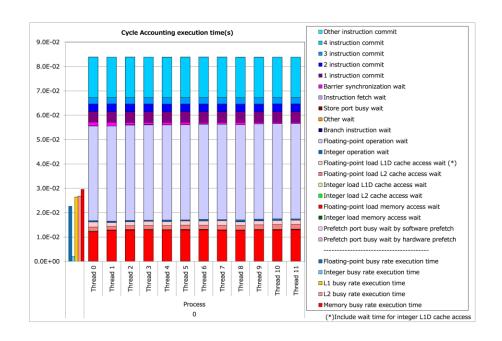
Tune-3 (ソフトウェアパイプラン)

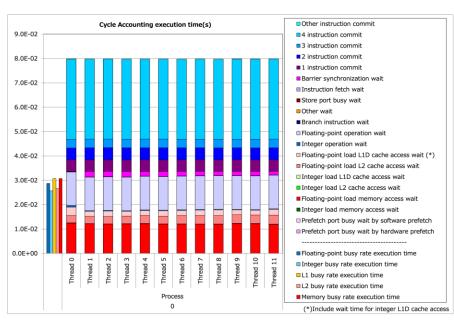
- mループではフルアンロール、iループはSIMD化
- その外側にj、kとループがあるが長さが8でソフトウェアパイプライニングを適用するには短い
 - 現状ではjについてフルアンロールされている
- k-loopとj-loopを手動でcollapseさせると長さが64になりswp が有効になる
 - OpenMPのcollapseは同時にスレッド並列化もするので使えない

```
for(int kj = 0; kj < p_Nq*p_Nq; ++kj) {
    const int j = kj % p_Nq;
    const int k = kj / p_Nq;</pre>
```

```
jwd8204o-i "axhelm-3.cpp", line 73: ループにソフトウェアパイプライニングを適用しました。
jwd8205o-i "axhelm-3.cpp", line 73: ループの繰返し数が64回以上の時、ソフトウェアパイプライニングを適用した
jwd8204o-i "axhelm-3.cpp", line 77: ループ制御変数'i'のループをSIMD化しました。
jwd8204o-i "axhelm-3.cpp", line 84: ループにソフトウェアパイプライニングを適用しました。
jwd8205o-i "axhelm-3.cpp", line 84: ループの繰返し数が48回以上の時、ソフトウェアパイプライニングを適用した
jwd6001s-i "axhelm-3.cpp", line 88: ループ制御変数'i'のループをSIMD化しました。
jwd8203o-i "axhelm-3.cpp", line 103: このループをフルアンローリングしました。
jwd8204o-i "axhelm-3.cpp", line 115: ループにソフトウェアパイプライニングを適用しました。
jwd8205o-i "axhelm-3.cpp", line 115: ループの繰返し数が48回以上の時、ソフトウェアパイプライニングを適用し
jwd6001s-i "axhelm-3.cpp", line 119: ループ制御変数'i'のループをSIMD化しました。
jwd8203o-i "axhelm-3.cpp", line 124: このループをフルアンローリングしました。
```

Before/after





- 待ち時間は減った?
- IPC $51.01 \rightarrow 2.01$ \(\text{a} \text{ } \text{p} \text{GIPS} \(\text{figure} 2.03 \rightarrow 4.02 \)
- 答え:整数命令が増えた、実行時間の減少は少しだけ

幾何係数の部分のデータ構造 ggeo[Nelement][7][512]

プリフェッチ

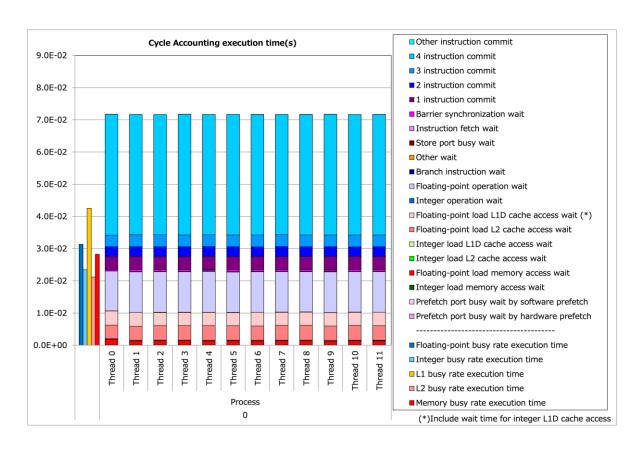
- 冗長なアドレス計算を手で取り除く
 - kとjを求めてからアドレス計算→kjから直接アドレス計算

```
for(int kj = 0; kj < p_Nq*p_Nq; ++kj) {
    const int j = kj % p_Nq;
    const int k = kj / p_Nq;

for(int i = 0; i < p_Nq; ++i) {
    // const dlong id = i + j * p_Nq + k * p_Nq * p_Nq + element * p_Np;
    const dlong gbase = element * p_Nggeo * p_Np + k * p_Nq * p_Nq + j * p_Nq + i;</pre>
```

- 急に速くなった
- アセンブリを見るとプリフェッチ命令が出るようになっていた
 - コンパイラからのメッセージはなし
 - PA情報からキャッシュミス情報を見ることでもわかったかも

最終版 (swp+prefetch)



- •L1ビジーが大き い(59.33%)
- 演算:ピークの 19.1%
- メモリスルー プット:ピーク の39.4%

axhelmまとめ

- プロファイラでわかること
 - メモリ待ち時間
 - 演算待ち時間
 - Gather/Scatter命令
 - ソースコード上の場所までは示してくれない
 - (スレッド間インバランス)
- チューニングの指針となったか?
 - 消したい待ち時間はある
 - 「このスイッチをオンにしたら待ち時間が消えます」という便利なものはない(-Kfastオプションから始めるなら)
 - 結局は試行錯誤による発見的手法になる
 - Before/afterの比較が情報量に富む

Bridge++ QXS (格子QCD)

- SIMDを2次元方向に使おうという実装
 - float16のSIMDなら例えばxy方向に4×4で
- Even-odd法との組み合わせ
- C++でtemplatを多用しておりコンパイルはClangモード
 - ACLE for SVEも多用
- 「正しく計算できるけれどflops値がいまいちなので」と金森 さんから調査と改善の依頼

格子QCDの計算

$$D_{x,y} = [1 + F(x)]\delta_{x,y} - \kappa \sum_{\mu=1}^{4} \left[(1 - \gamma_{\mu})U_{\mu}(x)\delta_{x+\hat{\mu},y} + (1 + \gamma_{\mu})U_{\mu}^{\dagger}(x - \hat{\mu})\delta_{x-\hat{\mu},y} \right],$$

- 格子点上にクォークのスピン (Dirac spinが3色で複素12語)
- リンク変数にゲージ (グルーオン) 4方向に3×3ユニタリ行列
- 9点(対角+8点)のステンシル計算

Naïveなデータ構造
real spin [NT][NZ][NY][NX][3][4][2];
real gauge[NT][NZ][NY][NX][4][3][3][2];

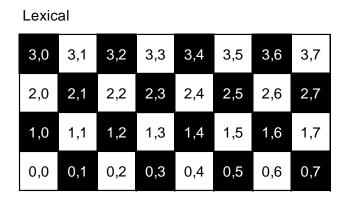
格子点あたり24語と72語、floatなら96 byteと144 byte

γ行列

• $(1 \pm \gamma_{\mu})$ を掛けるとスピンの自由度が複素4語から複素2語に減るのに注意

$$\gamma_{1} = \begin{pmatrix}
0 & 0 & 0 & i \\
0 & 0 & i & 0 \\
0 & -i & 0 & 0 \\
-i & 0 & 0 & 0
\end{pmatrix} \qquad
\gamma_{2} = \begin{pmatrix}
0 & 0 & 0 & 1 \\
0 & 0 & -1 & 0 \\
0 & -1 & 0 & 0 \\
1 & 0 & 0 & 0
\end{pmatrix} \qquad
\gamma_{3} = \begin{pmatrix}
0 & 0 & i & 0 \\
0 & 0 & 0 & -i \\
-i & 0 & 0 & 0 \\
0 & i & 0 & 0
\end{pmatrix} \qquad
\gamma_{4} = \begin{pmatrix}
1 & 0 & 0 & 0 \\
0 & 1 & 0 & 0 \\
0 & 0 & -1 & 0 \\
0 & 0 & 0 & -1
\end{pmatrix}$$

Even-odd法



Even	EVEII											
3,1	3,3	3,5	3,7									
2,0	2,2	2,4	2,6									
1,1	1,3	1,5	1,7									
0,0	0,2	0,4	0,6									

Odd			
3,0	3,2	3,4	3,6
2,1	2,3	2,5	2,7
1,0	1,2	1,4	1,6
0,1	0,3	0,5	0,7

$$\begin{pmatrix} D_{ee} & D_{eo} \\ D_{oe} & D_{oo} \end{pmatrix} \begin{pmatrix} x_e \\ x_o \end{pmatrix} = \begin{pmatrix} b_e \\ b_o \end{pmatrix}$$

$$\underbrace{\left(1 - D_{ee}^{-1} D_{eo} D_{oo}^{-1} D_{oe}\right)}_{\equiv D} x_e = D_{ee}^{-1} \left(b_e - D_{eo} D_{oo}^{-1} b_o\right)$$

$$x_o = D_{oo}^{-1} \left(b_o - D_{oe} x_e\right)$$

Odd

- D_{oe} : Even成分を読んでOdd成分に書き込み \Rightarrow D_{eo}
- *D_{ee}, D_{oo}*: Wilsonカーネルでは単なる対角行列
- •1回のmultの演算数は減らないが2 hopするので収束が早い
- CG法等で複数のベクトルを保持するメモリフットプリントが 半分

QWS & QXS

- QWS: QCD Wide SIMD
 - 重点課題9のターゲットアプリ
 - 単精度16語のSIMDとeven-odd法の組み合わせにより、MPIプロセス当たりのx軸の長さが**32の倍数**である必要
 - 目標の1924格子には使えるが、、、
- Bridge++ QXS
 - SIMDの軸を2次元化
 - x軸とy軸の4×4タイルをSIMDのレジスタに載せる
 - その分x方向とy方向のステンシルアクセスが非自明に

```
例えばこんなレイアウトに
float spin eo [NT][NZ][NV/4][NV/8] [3][4]
```

```
float spin_eo [NT][NZ][NY/4][NX/8] [3][4] [2] [16]; float gauge_eo[NT][NZ][NY/4][NX/8] [4][3][3][2] [16];
```

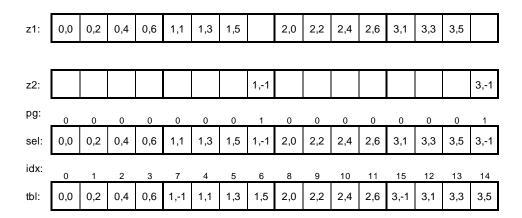
SIMD命令の複素数サポートについて

- 一部のSIMD命令では偶奇のレーンを束ねての複素数演算をサポートしている
 - BlueGeneや京、FX100(倍精度限定)
 - x86はshuffl系の命令を多用する
 - 第1オペランドの偶数レーン複製
 - 第1オペランドの奇数レーン複製
 - 第2オペランドの偶奇スワップと符号変更
- これがあると格子QCDの計算で嬉しい
 - 実効的にベクトル長が半分になる
 - 実効的にレジスタの本数が倍になる
 - スピンの12本をレジスタに置きっぱなしにできる
 - ゲージの9本をレジスタに置きっぱなしにできる (DWF)
- SVEの命令セットには入ったがA64FXではマイクロコード実装

ステンシルアクセス: X-

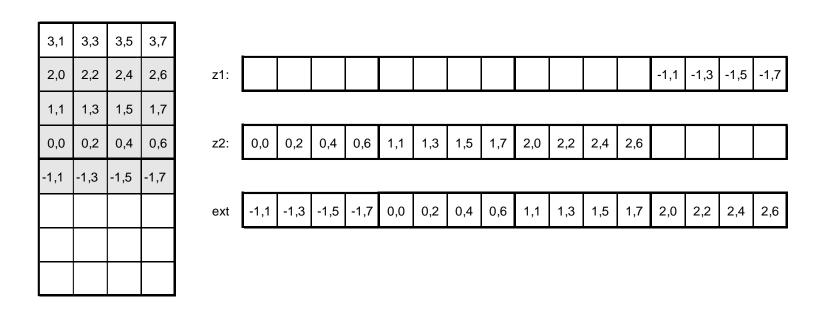
	3,-1	3,1	3,3	3,5	3,7
		2,0	2,2	2,4	2,6
	1,-1	1,1	1,3	1,5	1,7
		0,0	0,2	0,4	0,6

3,0	3,2	3,4	3,6
2,1	2,3	2,5	2,7
1,0	1,2	1,4	1,6
0,1	0,3	0,5	0,7



- Oddのひとつ左(X⁻)のevenからの読み込み(右にシフトして 足す)
- •yが奇数の行だけ左にずれる
- load, load, sel, tblで実現可能

ステンシルアクセス: Y-



- load, load, extでできる
- ext: 2つのベクトルを連結して即値でシフト

通信バッファのgather/scatte回避

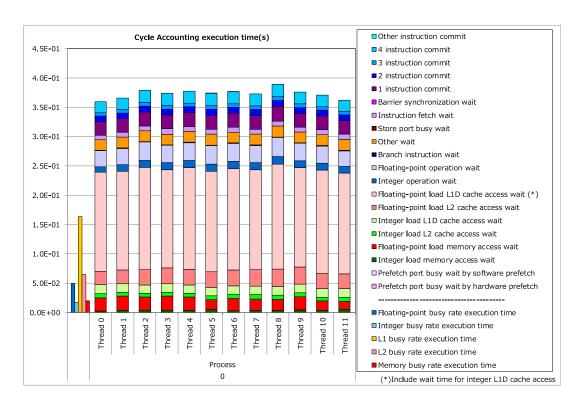
3,1	3,3	3,5	3,7	z1:								1,7								3,7
2,0	2,2	2,4	2,6	pg:	0	0	0	0	0	0	0	1	0	0	0	0	0	0	0	1
1,1	1,3	1,5	1,7	compact:	1,7	3,7														
0,0	0,2	0,4	0,6	idx:								0								1
tbl:												1,7								3,7

- 16語中2語だけ送信バッファに詰める (compact)
 - レジスタ内で連続にしscatterを回避
- 受信バッファから元の場所に戻す(tbl)
 - Gatherを回避

計測条件、チューニング

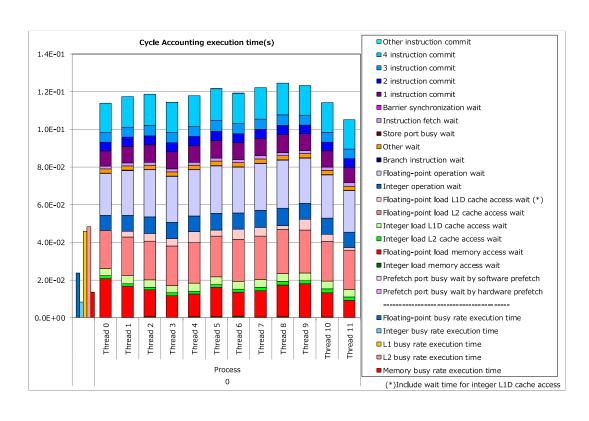
- A64FX 1ノード4プロセス、16⁴=64K格子
 - CMGあたり8×8×16×16
 - SIMDとEOも含めると 8×8×4×2
- カーネルとなる関数
 - eo1: 送信バッファの準備
 - 途中でスレッドの同期とMPI送信のキック
 - bulk: ローカルでできる計算
 - スレッドを同期してMPIの受信待ち
 - eo2: 受信バッファからの計算
- コンパイラオプション
 - -Rpass-missed=inline -mllvm -inline-threshold=1000 を追加
 - 8方向からのhopのカーネルで関数callがなくなるように

As-is状態



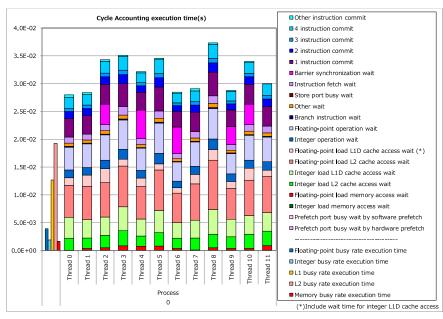
- Bulk部分
- L1アクセス待ち、 L1ビジー率が高い
- Gather/scatterが沢 山出ていた
 - 実は対角の部分、最 後にκを掛けてスト アするだけの部分
 - Vsimd_tという SIMDエミュレーション型
 - といっても単なる外側24、内側16の二重ループなのだが

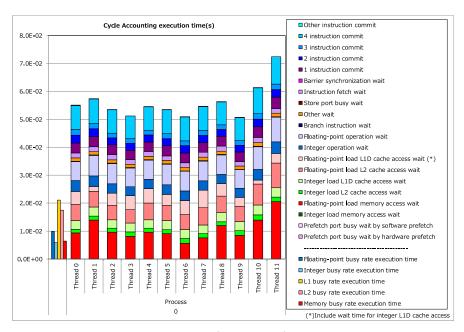
チューン(性能デバッグ)後



- Bulk
- 普通の見た目になっている
- 180 Gflops→360 Gflops
- 若干主記憶アクセス が発生しているのが 気になる
 - 6MiB/CMGぐらいの はず
- 若干のロードインバランス

ロードバランス問題





eo1 (export)

eo2 (import)

- 実行時間/演算数含め可視化してくれる
- 格子点数均等割にしているので、袖に触る/触らないで負荷バランスが崩れる

まとめ (プロファイラで何が見えたか)

- 手を動かせば必ず絵が出てくるのが有難い
- (意図せぬ) Gather/Scatter
- ロードインバランス
- 見ておくべきもの
 - ビジー率
 - ある程度のチューンができたコードは大抵複数のリソースが忙しい
 - 「実アプリだから主記憶帯域ネック」は本当?
- 見えにくいもの
 - ソースコードの何処でそれが起こっているか? (区間の細分化には限度がある)
 - アセンブリ出力を確認することの代替にはならない