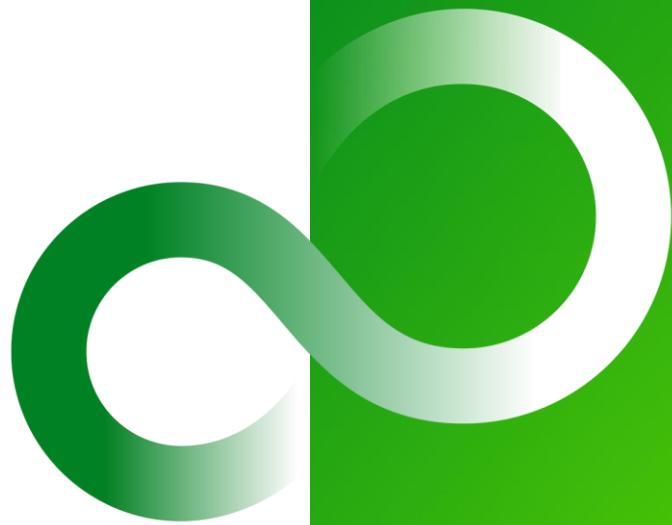


プログラミングガイド チューニング編

2023年3月

v2.2

富士通株式会社



富士通株式会社の許諾を得て一般公開しています。内容は国立研究開発法人理化学研究所にご確認ください。

■ 当資料は、A64FX プロセッサを対象としたアプリケーション開発者やチューニング実施者のためのものです。

■ 注意

- FortranとC/C++やtradモードとclangモードではコンパイラが異なるため、チューニングの方法が異なる場合や対応するチューニング方法がない場合があります
- tradモードとclangモードで類似のチューニング手法の場合、clangモードのチューニングを割愛します

■ 当資料と併せて以下も参考してください。

- Fortran 使用手引書
- C/C++ 使用手引書
- プログラミングガイド プロセッサ編
- プログラミングガイド プログラミング共通編
- プログラミングガイド Fortran編

■ 当資料の記載においては、以下の文章を参考にしています。

- A64FX 論理仕様書
- A64FX®Microarchitecture Manual
- ARM® Architecture Reference Manual
(ARMv8 , ARMv8.1 , ARMv8.2 , ARMv8.3)
- ARM® Architecture Reference Manual Supplement
The Scalable Vector Extension

■ 商標について

- Linux® は、Linus Torvalds 氏の米国およびそのほかの国における登録商標または商標です。
- Red Hat は、Red Hat, Inc. の米国およびそのほかの国における登録商標または商標です。
- ARMは、ARM Ltd.の英国およびその他の国における登録商標です。
- そのほかの会社名、製品名等の固有名詞は、各社の登録商標または商標です。
- 本資料に記載されているシステム名、製品名等には、必ずしも商標表示(®、™)を付記していません

■ 謝辞

- C/C++のチューニングについては、国立研究開発法人理化学研究所 計算科学研究センターから委託を受けて実施した事業の成果を基にしています。

アプリケーションのチューニングは、観点によって次のようなポイントがあります。当資料はCPUチューニングとスレッド並列チューニングについて説明しています。

当資料の説明範囲				
	1コア チューニング	スレッド並列 チューニング	プロセス並列 チューニング	超高並列 チューニング
チューニングポイント	<ul style="list-style-type: none">✓ I/Oを減らす✓ 演算量を減らす✓ SIMD化を促進する✓ メモリアクセスを減らす✓ キャッシュ利用率を高める	<ul style="list-style-type: none">✓ 並列化率を上げる✓ 並列化粒度を上げる✓ スレッド間の同期コストを削減する✓ ロードバランスを均等化する	<ul style="list-style-type: none">✓ 並列化率を上げる✓ 並列化粒度を上げる✓ プロセス間の通信コストを削減する✓ ロードバランスを均等化する	<ul style="list-style-type: none">✓ 並列化率を上げる✓ 並列化粒度を上げる✓ プロセス間の通信コストを削減する✓ ロードバランスを均等化する✓ グローバルな通信コストを削減する

目次 (1/2)

□ ボトルネックの調査

- [CPU性能解析レポート：全体](#)
- [CPU性能解析レポート：サイクルアカウンティングとは](#)
- [サイクルアカウンティングを活用したボトルネック抽出](#)
- [各種ビギー時間を活用した（バンド幅）ボトルネック抽出](#)
- [サイクルアカウンティングを活用したチューニング方法](#)
- [チューニングマップ](#)

□ 1コアチューニング

- [データアクセス待ち（データ局所性の向上）](#)
 - [ストリップマイニング](#)
 - [ループブロッキング](#)
 - [セクタキヤッショウ](#)
 - [ループ交換](#)
 - [ループ融合](#)
 - [配列マージ（インダイレクトアクセスの効率改善）](#)
 - [配列次元移動](#)
 - [アンロールアンドジャム](#)

□ データアクセス待ち（レイテンシの隠蔽）

- [インダイレクトアクセスプリフェッチ](#)
- [連続アクセスでないデータへのソフトウェアプリフェッチの利用](#)

□ データアクセス待ち（アクセス量の軽減）

- [高速ストア\(ZFILL\)](#)

□ データアクセス待ち（スラッシングの改善）

- [1次元目の配列要素を増やすパディング](#)
- [2次元目の配列要素を増やすパディング](#)
- [ダミー配列によるパディング](#)
- [ダミー配列によるパディング（サイズが異なる配列）](#)
- [配列マージ（スラッシングの改善）](#)
- [ループ分割（スラッシングの改善）](#)
- [ラージページ環境変数によるパディング](#)

目次 (2/2)

□ 演算待ち (SIMD化の促進)

- [ループピーリング](#)
- [定義関係が明らかでないループ](#)
- [ポインタ変数が含まれるループ](#)

□ 演算待ち (レイテンシの隠蔽)

- [ループ分割 \(ソフトウェアパイプラインング促進\)](#)
- [適切なアンローリング展開数の指定およびソフトウェアパイプラインングの抑止](#)
- [ストライピング \(インターリーブ\) 展開数の指定およびソフトウェアパイプラインングの抑止](#)
- [外側ループでのソフトウェアパイプラインング](#)
- [リローリング](#)
- [ループアンスイッチング](#)

□ マイクロアーキ依存の改善

- [Scatter Store命令の回避](#)
- [Gather Load命令の纏め上げの促進](#)
- [過剰SFIの回避](#)
- [Multiple Structures命令の活用](#)
- [ハードウェアプリフェッчの距離調整](#)
- [SVEのベクトルレジスタサイズ \(SIMD幅\)](#)
- [半精度実数型の活用](#)

□ スレッド並列チューニング

- [スレッド並列化率の改善](#)
 - [定義引用関係が明らかでないループの改善](#)
 - [ポインタ変数が含まれるループの改善](#)
 - [データ依存関係があるループの改善](#)
- [スレッド並列化実行効率の改善](#)
 - [False Sharingの改善](#)
 - [処理量が不規則なループの改善](#)
 - [適切な並列化次元での並列化](#)
- [ラージページの設定による実行効率の改善](#)
 - [ラージページのページングポリシー指定](#)
 - [ロックタイプの変更](#)
- [メモリ使用量の削減](#)
 - [OMP SINGLE から OMP MASTER への書き換え](#)

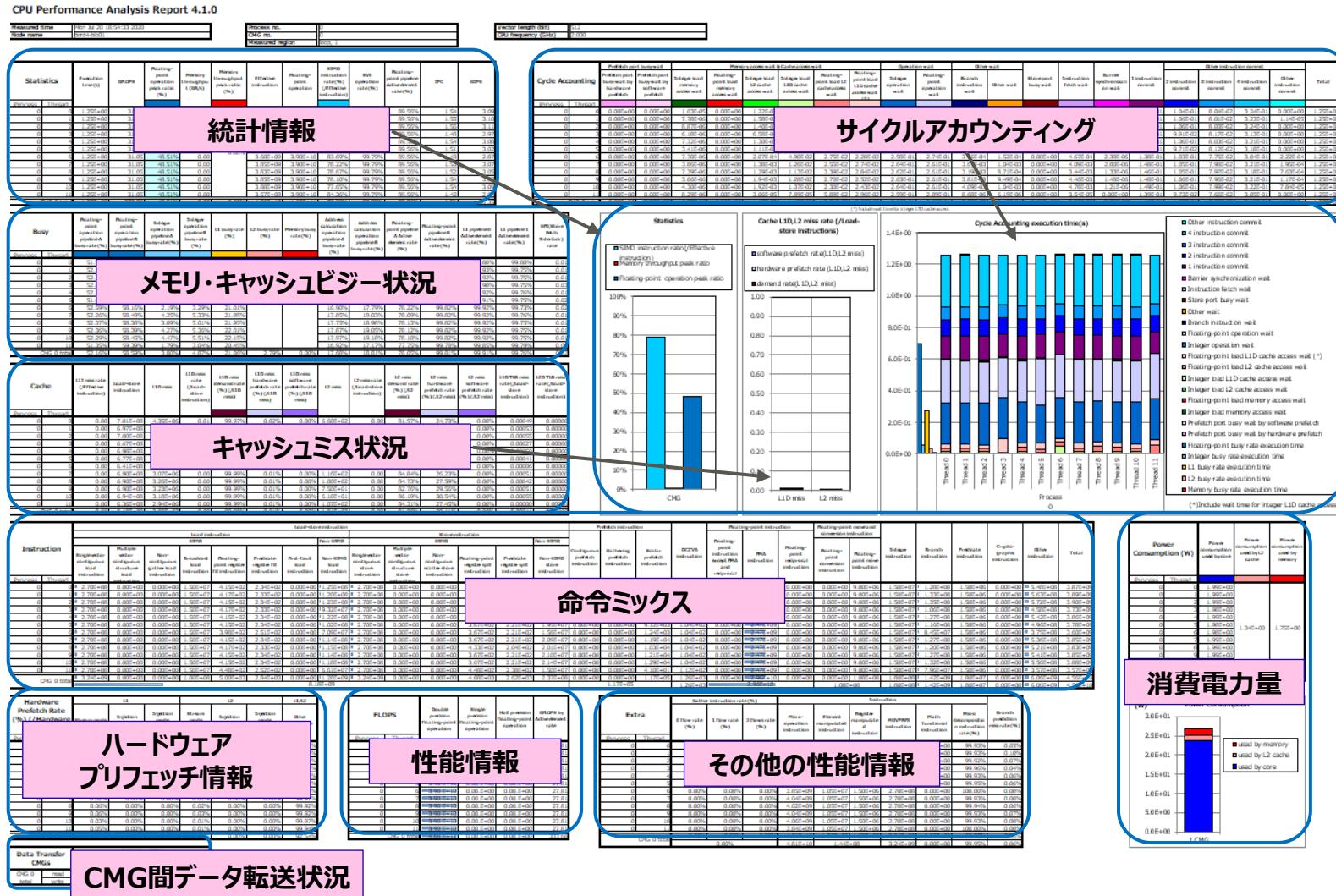
ボトルネックの調査

CPU性能解析レポート：全体

FUJITSU

性能ボトルネックの抽出手段には、CPU性能解析レポートの利用を薦めます。

CPU性能解析レポートでは、以下のような豊富な種類のPA(Performance Analysis)イベントを計測でき、アプリケーションプログラム実行時のCPU動作状態を調べることができます。



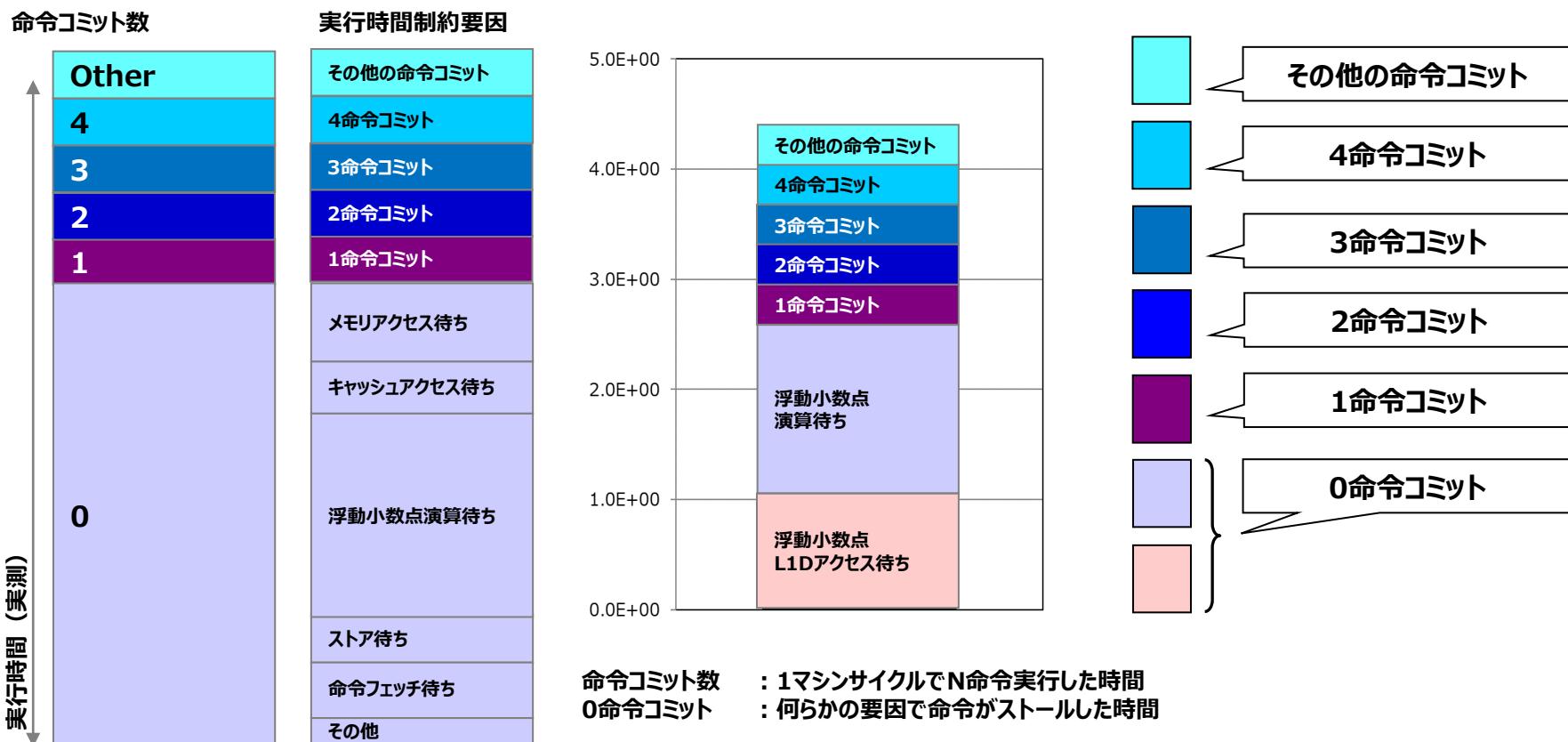
CPU性能解析レポート：サイクルアカウンティングとは

FUJITSU

サイクルアカウンティングとは、性能ボトルネックの要因分析手法です。

サイクルアカウンティング情報は、CPU性能解析レポートの右上部に掲載されています。

サイクルアカウンティングでは、あるアプリケーションプログラムを実行するためにかかった総時間（CPUサイクル数）をCPUの動作状態で分類してグラフ化します。そのグラフからCPU内のボトルネックが把握できるので、詳細な性能分析やチューニングを行うことができます。



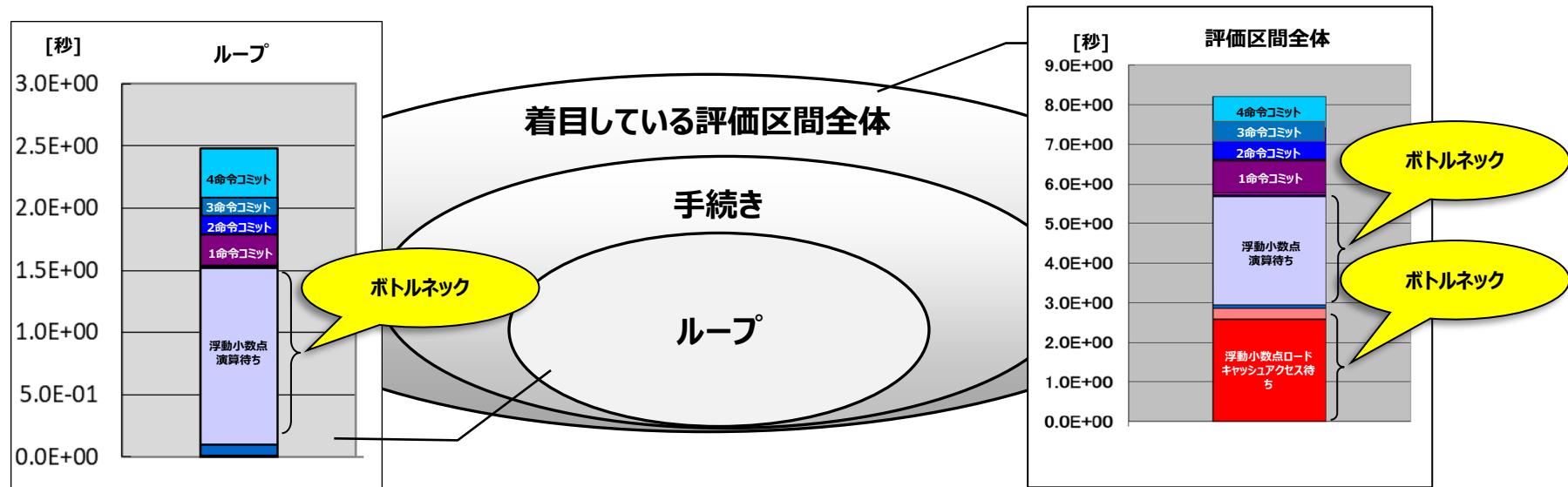
サイクルアカウンティングを活用したボトルネック抽出

FUJITSU

● ボトルネックの把握

チューニングの基本はボトルネックを把握することです。

評価区間のボトルネックはサイクルアカウンティングから判断することができます。



● チューニングへの活用

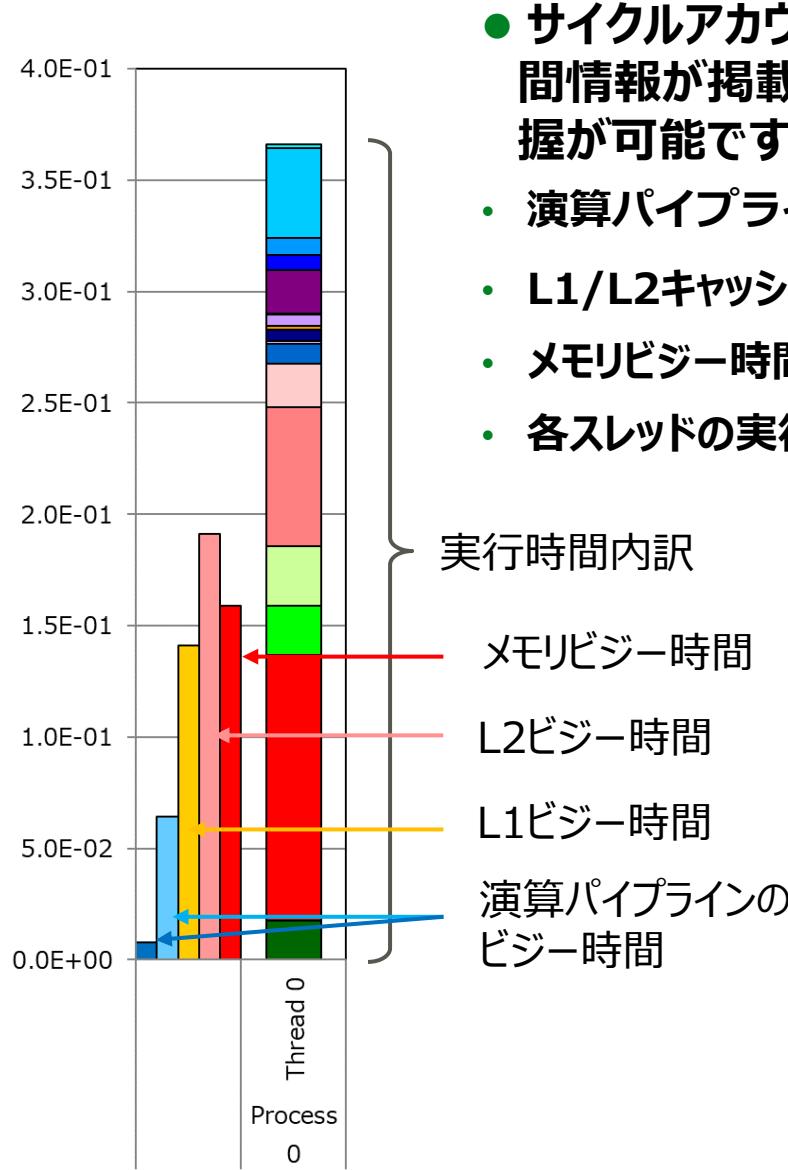
- ・ ボトルネックの改善にはどのような対策が必要か
- ・ どの程度改善することができるのか

を判断するためには

ループ単位までブレークダウンして分析した方が良い。

各種ビギー時間を活用した（バンド幅）ボトルネック抽出

FUJITSU



- サイクルアカウンティングの左側のグラフには以下のビギー時間情報が掲載されています。各種バンド幅のボトルネック把握が可能です。

- 演算パイプラインのビギー時間
- L1/L2キャッシュビギー時間
- メモリビギー時間
- 各スレッドの実行時間内訳

- メモリビギー時間
メモリへのデータ転送量が多い場合に発生します。
- L1/L2キャッシュビギー時間
L1/L2キャッシュへのデータ転送量が多い場合に発生します。
- 演算パイプラインのビギー時間
演算量が多い場合に発生します。

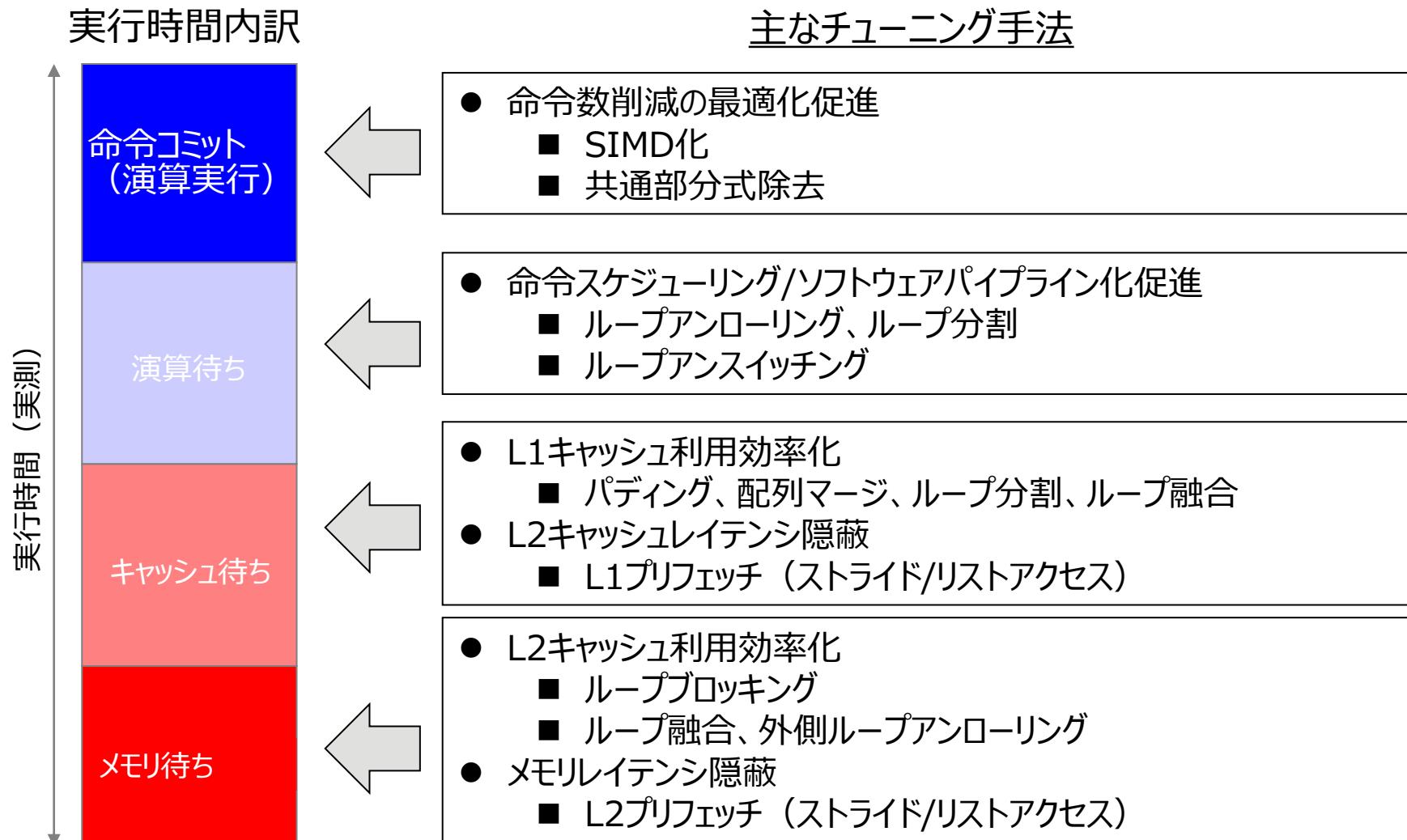
サイクルアカウンティングを活用したチューニング方法

FUJITSU

● サイクルアカウンティングの結果から、チューニング手法を選択します。

次に主なチューニング手法を示します。

より詳細な内容はチューニングマップを参照してください。



チューニングマップ：分類と状態

ボトルネックの分類	PAグラフから見える高コスト	PA情報から見える高コスト	状態
メモリネック	浮動小数点ロードメモリアクセス待ち	-	メモリレイテンシがボトルネックになっています。
	整数ロードメモリアクセス待ち	-	
	ストア待ち	-	ストア命令のコストがボトルネックになっています。
	メモリ・キャッシュビギー待ち	-	メモリスループットがボトルネックになっています。
	-	メモリビギー率が高い	
	-	L2ミス率が高い L2ミスdm率が高い	メモリレイテンシがボトルネックになっています。
L2キャッシュネック	浮動小数点ロードL2アクセス待ち	-	L2キャッシュレイテンシがボトルネックになっています。
	整数ロードL2アクセス待ち	-	
	-	L1ビギー率が高い	L2キャッシュスループットがボトルネックになっています。
	-	L1Dミス率が高い L1Dミスdm率が高い	L2キャッシュレイテンシがボトルネックになっています。
L1キャッシュネック	浮動小数点ロードL1Dアクセス待ち	-	L1キャッシュレイテンシがボトルネックになっています。
	整数ロードL1Dアクセス待ち	-	
	-	L1ビギー率が高い	L1キャッシュスループットがボトルネックになっています。
スケジューリングネック	浮動小数点演算待ち	-	演算命令のレイテンシがボトルネックになっています。
	整数演算待ち	-	
	分岐命令待ち	-	分岐命令がボトルネックになっています。
並列化ネック	パリア同期待ち	-	スレッド並列化されていない部分がボトルネックになっています。
ロードインバランスネック	パリア同期待ち	命令バランスのmax-min差が大きい	スレッド間のロードインバランスがボトルネックになっています。
TLBネック	-	mDTLBミス率が高い	TLBミスやTLBスラッシングがボトルネックになっています。
	-	uDTLBミス率が高い	TLBミスがボトルネックになっています。
命令フェッチ	命令フェッチ待ち	-	命令キャッシュミスやスラッシングがボトルネックになっています。
命令数ネック	命令コミット	その他の命令コミット	命令数がボトルネックになっています。
		4命令コミット 3命令コミット 2命令コミット 1命令コミット	
その他	その他の待ち	-	PAが正しく採取できていない可能性があります。

チューニングマップ (1/4)

● スループットネック

PAグラフから見える高コスト	PA情報から見える高コスト	状態	チューニング案
メモリ・キャッシュビジー待ち	-	メモリスループットがボトルネックになっています。	データアクセス待ちの改善 - 配列次元移動 - ループブロッキング - ストリップマイニング - 高速ストア (ZFILL)
-	メモリビジー率が高い	メモリスループットがボトルネックになっています。	データアクセス待ちの改善 - 配列次元移動 - ループブロッキング - ストリップマイニング - 高速ストア (ZFILL)
-	L2ミス率が高い L2ミスdm率が高い	メモリレイテンシがボトルネックになっています。	データアクセス待ちの改善 - 配列次元移動 - ループブロッキング - ストリップマイニング - セクタキャッシュ - ブリフェッチに関する改善 - スラッシングの改善
	L2ビジー率が高い	L2キャッシュスループットがボトルネックになっています。	データアクセス待ちの改善 - 配列次元移動 - ループブロッキング - ストリップマイニング
	L1ビジー率が高い	L1キャッシュスループットがボトルネックになっています。	データアクセス待ちの改善 - アルゴリズムの見直し

チューニングマップ (2/4)

● レイテンシック

PAグラフから見える高コスト	PA情報から見える高コスト	状態	チューニング案
浮動小数点ロードメモリアクセス待ち		メモリレイテンシがボトルネックになっています。	データアクセス待ちの改善 - 配列次元移動 - プリフェッчに関する改善 - ループブロッキング
整数ロードメモリアクセス待ち		L2キャッシュレイテンシがボトルネックになっています。	データアクセス待ちの改善 - 配列次元移動 - アンロールアンドジャム - プリフェッчに関する改善
浮動小数点ロードL2アクセス待ち		L2キャッシュレイテンシがボトルネックになっています。	データアクセス待ちの改善 - 配列次元移動 - スラッシングの改善
整数ロードL2アクセス待ち			
-	1Dミス率が高い L1Dミスdm率が高い		
浮動小数点ロードL1Dアクセス待ち		L1キャッシュレイテンシがボトルネックになっています。	命令スケジューリングの改善 マイクロアーキ依存の改善
整数ロードL1Dアクセス待ち			
浮動小数点演算待ち		演算命令のレイテンシがボトルネックになっています。	命令スケジューリングの改善
整数演算待ち			

チューニングマップ (3/4)

● 命令数ネック

PAグラフから見える高コスト	PA情報から見える高コスト	状態	チューニング案
命令コミット	その他の命令コミット	命令数がボトルネックになっています。	命令数ネックの改善 <ul style="list-style-type: none">- SIMD化促進- ソフトウェアパイプライン促進- プリフェッчに関する改善- インライン展開
	4命令コミット		
	3命令コミット		
	2命令コミット		
	1命令コミット		

● TLBネック

PAグラフから見える高コスト	PA情報から見える高コスト	状態	チューニング案
-	mDTLBミス率が高い	TLBミスやTLBスラッシングがボトルネックになっています。	TLBネックの改善 <ul style="list-style-type: none">- スラッシングの解消- 使用領域の変更- ラージページのオプションによる最適化
-	uDTLBミス率が高い	TLBミスがボトルネックになっています。	TLBネックの改善 <ul style="list-style-type: none">- ページサイズの拡張

チューニングマップ (4/4)

● その他

PAグラフから見える高コスト	PA情報から見える高コスト	状態	チューニング案
ストア待ち		ストア命令のコストがボトルネックになっています。	データアクセス待ちの改善 - 配列次元移動 - プリフェッчに関する改善 - 高速ストア (ZFILL)
分岐命令待ち		分岐命令がボトルネックになっています。	命令スケジューリングの改善 - IF文の除去 - マスク付きSIMD
バリア同期待ち		スレッド並列化されていない部分がボトルネックになっています。	スレッド並列化率の改善
	命令バランスのmax-min 差が大きい	スレッド間のロードインバランスがボトルネックになっています。	スレッド並列実行効率の改善
命令フェッチ待ち		命令キャッシュミスやスラッシングがボトルネックになっています。	命令フェッチの改善 - ループボディの縮小 - アルゴリズムの見直し - スラッシングの解消

1コアチューニング データアクセス待ち (データ局所性の向上)

- データ局所性とは
- ストリップマイニング
- ループブロッキング
- セクタキャッシュ
- ループ交換
- ループ融合
- 配列マージ（インダイレクトアクセスの効率改善）
- 配列次元移動
- アンロールアンドジャム

データ局所性とは

データ局所性とは、データの参照やアクセスが狭い範囲に集中している度合いです。データ局所性が低いと、キャッシュメモリ上のデータが利活用されず、メモリアクセス負荷が大きくなります。ソースチューニングによりデータ局所性を高めることで、メモリアクセスの負荷が軽減され、データアクセス待ちが改善することができます。

データ局所性の向上には、以下のような手法が有効です。

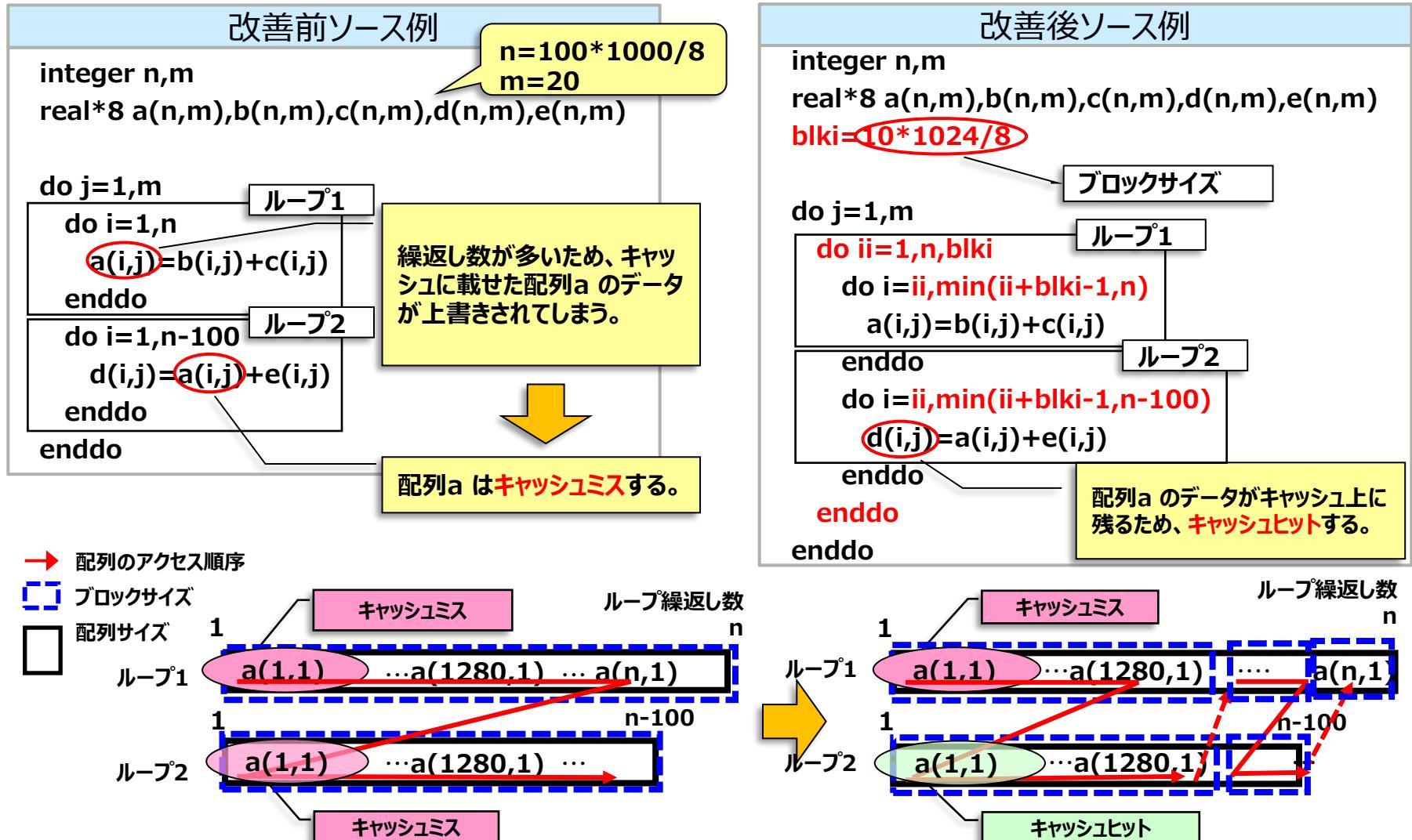
- ストリップマイニング
- ループブロッキング
- セクタキャッシュ
- ループ交換
- ループ融合
- 配列マージ(インダイレクトアクセスの効率改善)
- 配列次元移動
- アンロールアンドジャム

ストリップマイニング

- ストリップマイニングとは
- ストリップマイニング（改善前）
- ストリップマイニングの効果（ソースチューニング）

ストリップマイニングとは

ストリップマイニングとは、ループをより小さなセグメントやストリップに断片化することで、キャッシュ効率を向上させる手法です。



ループ¹の繰返し数が多く、配列データはキャッシュに載りきらないため、ループ²で再利用できません。そのため、メモリ・キャッシュビジー待ちが多くなっています。

改善前ソース

```

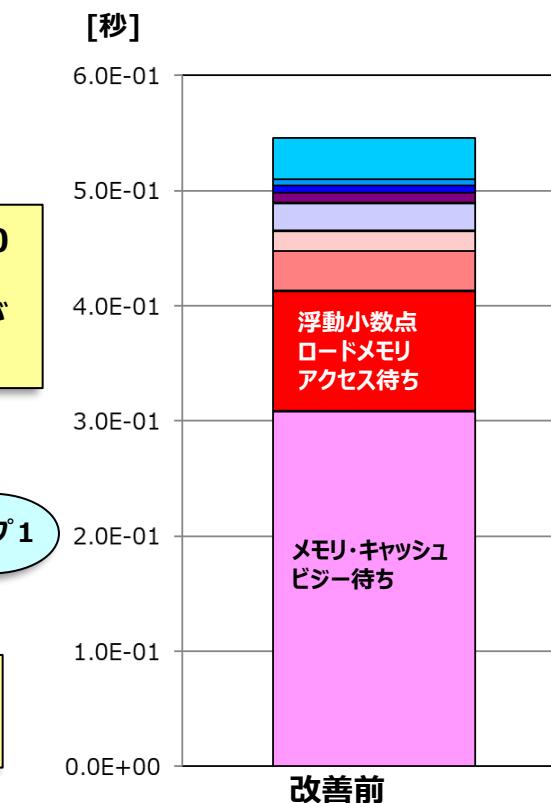
32      !$omp parallel do reduction(+:s1)
<<< Loop-information Start >>>
<<< [OPTIMIZATION]
<<< PREFETCH(HARD) Expected by compiler :
<<<   d, c, b, a, e
<<< Loop-information End >>>
33  1 p   do j=1,m
<<< Loop-information Start >>>
<<< [OPTIMIZATION]
<<< SIMD(VL: 8)
<<< PREFETCH(HARD) Expected by compiler :
<<<   b, c, a, d
<<< Loop-information End >>>
34  2 p  8v   do i=1,n
35  2 p  8v     s1 = s1 + a(i,j) * (s3 * b(i,j) + c(i,j) * (s2 + s3 * d(i,j)))
36  2 p  8v   enddo
<<< Loop-information Start >>>
<<< [OPTIMIZATION]
<<< SIMD(VL: 8)
<<< SOFTWARE PIPELINING
<<< PREFETCH(HARD) Expected by compiler :
<<<   b, a, d, c, e
<<< Loop-information End >>>
37  2 p  2v   do i=1,n-100
38  2 p  2v     e(i,j) = s2 * (a(i,j) + b(i,j) * (s3 + c(i,j) * d(i,j)))
39  2 p  2v   enddo
40  1 p   enddo

```

ループ¹繰返し数 : 375000
配列サイズの合計 : 12MB
繰返し数が多いため、配列データがキャッシュに載りきらない

配列アクセスがキャッシュミスする

ループ²



データの再利用ができないため、L1D miss rateとL2 miss rateがストリームアクセスの理論値である0.20前後になっている。

Cache

	L1I miss rate (/Effective instruction)	Load-store instruction	L1D miss	L1D miss rate (/Load-store instruction)	L1D miss demand rate (%) (/L1D miss)	L1D miss hardware prefetch rate (%) (/L1D miss)	L1D miss software prefetch rate (%) (/L1D miss)	L2 miss	L2 miss rate (/Load-store instruction)	L2 miss demand rate (%) (/L2 miss)	L2 miss hardware prefetch rate (%) (/L2 miss)	L2 miss software prefetch rate (%) (/L2 miss)
改善前	0.00	2.08E+09	4.23E+08	0.20	7.61%	92.40%	-0.01%	4.23E+08	0.20	5.17%	95.41%	0.00%

ストリップマイニングによって、キャッシュ効率が向上し、メモリ・キャッシュビジー待ちが改善されました。

改善後ソース

```

32      blki=4*1024/8
33
34      !$omp parallel do reduction(+:s1)
35 1 p       do j=1,m
36 <<< Loop-information Start >>>
37 <<< [OPTIMIZATION]
38 <<< PREFETCH(HARD) Expected by compiler :
39 <<< d, c, b, a, e
40 <<< Loop-information End >>>
41 2 p       do ii=1,n,blki
42 <<< Loop-information Start >>>
43 <<< [OPTIMIZATION]
44 <<< SIMD(VL: 8)
45 <<< PREFETCH(HARD) Expected by compiler :
46 <<< b, c, a, d
47 <<< Loop-information End >>>
48 3 p 8v     do i=ii,min(ii+blki-1,n)
49 3 p 8v       s1 = s1 + a(i,j) * (s3 * b(i,j) + &
50 3           c(i,j) * (s2 + s3 * d(i,j) ))
51 40 3 p 8v     enddo
52 <<< Loop-information Start >>>
53 <<< [OPTIMIZATION]
54 <<< SIMD(VL: 8)
55 <<< SOFTWARE PIPELINING
56 <<< PREFETCH(HARD) Expected by compiler :
57 <<< b, a, d, c, e
58 <<< Loop-information End >>>
59 41 3 p 2v    do i=ii,min(ii+blki-1,n-100)
60 42 3 p 2v      e(i,j) = s2 * (a(i,j) + &
61 43 3           b(i,j) * (s3 + c(i,j) * d(i,j)))
62 44 3 p 2v    enddo
63 45 2 p      enddo
64 46 1 p      enddo

```

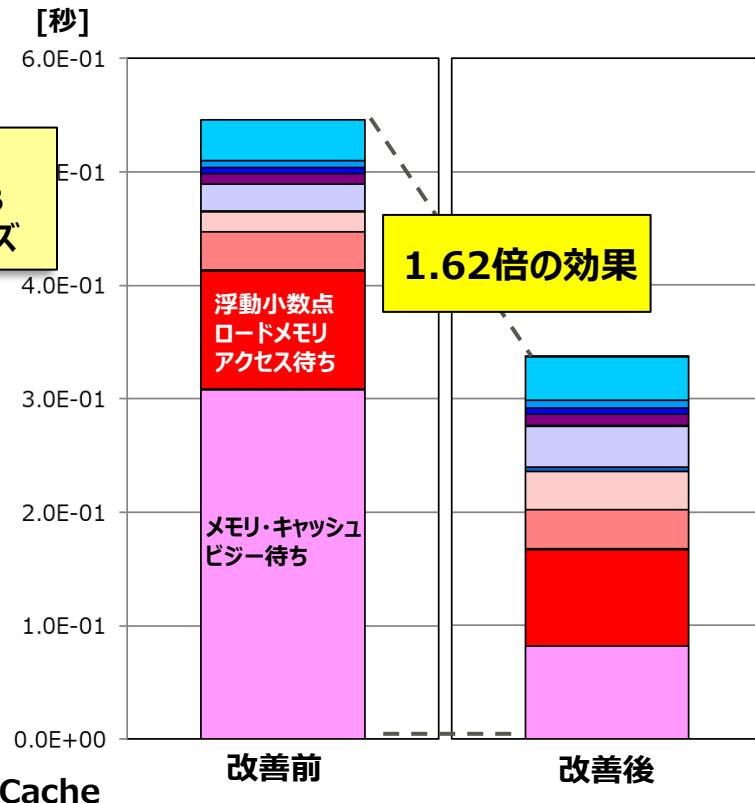
ブロックサイズ : 4KB
4KB×4ストリーム=16KB
→ L1キャッシュに載るサイズ

ループ 1

配列アクセスがキャッシュヒットする

ループ 2

L1DミスとL2ミスが減少した



	L1I miss rate (/Effective instruction)	Load-store instruction	L1D miss	L1D miss rate (/Load-store instruction)	L1D miss demand rate (%) (/L1D miss)
改善前	0.00	2.08E+09	4.23E+08	0.20	7.61%
改善後	0.00	1.95E+09	2.35E+08	0.12	15.24%

	L2 miss	L2 miss rate (/Load-store instruction)	L2 miss demand rate (%) (/L2 miss)
改善前	4.23E+08	0.20	5.17%
改善後	2.35E+08	0.12	7.84%

ループ¹の繰返し数が多く、配列データはキャッシュに載りきらないため、ループ²で再利用できません。そのため、メモリ・キャッシュビジー待ちが多くなっています。

改善前ソース

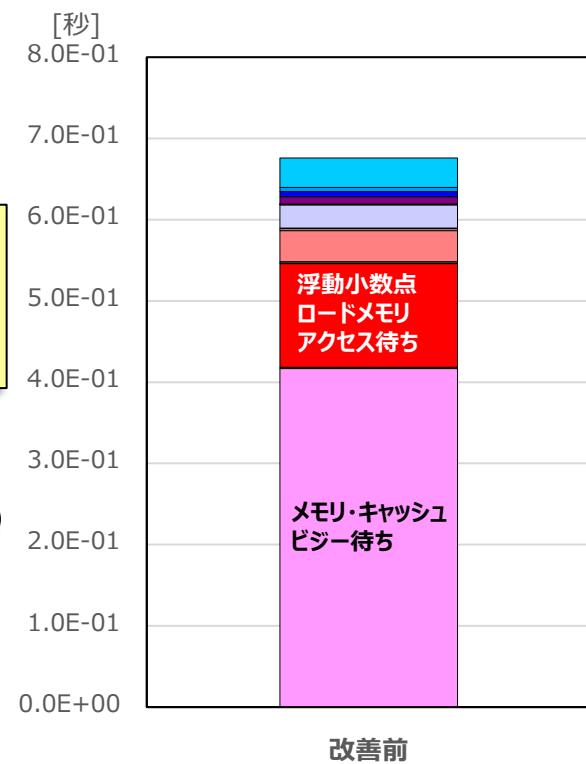
```

35 #pragma omp parallel for reduction(+:s1)
<<< Loop-information Start >>>
<<< [OPTIMIZATION]
<<< PREFETCH(HARD) Expected by compiler :
<<< (unknown)
<<< Loop-information End >>>
36 p   for(j=0;j<m;j++) {
<<< Loop-information Start >>>
<<< [OPTIMIZATION]
<<< SIMD(VL: 8)
<<< PREFETCH(HARD) Expected by compiler :
<<< (unknown)
<<< Loop-information End >>>
37 p   8v   for(i=0;i<n;i++) {
38 p   8v     s1 = s1 + a[j][i] * (s3 * b[j][i] + c[j][i] * (s2 + s3 * d[j][i]));
39 p   8v   }
<<< Loop-information Start >>>
<<< [OPTIMIZATION]
<<< SIMD(VL: 8)
<<< SOFTWARE PIPELINING(1)
<<< PREFETCH(HARD) Expected by compiler :
<<< (unknown)
<<< Loop-information End >>>
40 p   2v   for(i=0;i<n-100;i++) {
41 p   2v     e[j][i] = s2 * (a[j][i] + b[j][i] * (s3 + c[j][i] * d[j][i]));
42 p   2v   }
43 p   2v   }
44 }

```

ループ^{繰返し数} : 375000
配列サイズの合計 : 12MB
繰返し数が多いため、配列データがキャッシュに載りきらない

配列アクセスがキャッシュミスする



データの再利用ができないため、L1D miss rateとL2 miss rateがストリームアクセスの理論値である0.20前後になっている。

Cache	L1I miss rate (/Effective instruction)	Load-store instruction	L1D miss	L1D miss rate (/Load-store instruction)	L1D miss demand rate (%) (/L1D miss)	L1D miss hardware prefetch rate (%) (/L1D miss)	L1D miss software prefetch rate (%) (/L1D miss)	L2 miss	L2 miss rate (/Load-store instruction)	L2 miss demand rate (%) (/L2 miss)	L2 miss hardware prefetch rate (%) (/L2 miss)	L2 miss software prefetch rate (%) (/L2 miss)
改善前	0.00	2.17E+09	4.23E+08	0.19	9.14%	90.86%	0.00%	4.23E+08	0.19	5.42%	95.53%	0.00%

ストリップマイニングによって、キャッシュ効率が向上し、メモリ・キャッシュビジー待ちが改善されました。

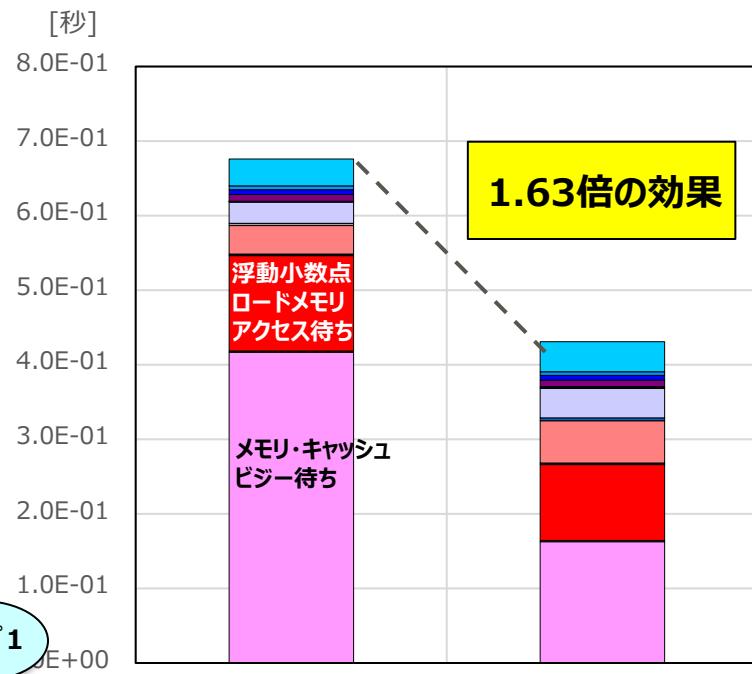
改善後ソース

```

37     blki=4*1024/8;
38
39 #pragma omp parallel for reduction(+:s
40 p     for(j=0;j<m;j++) {
41 p         <<< Loop-information Start >>>
42 p         <<< [OPTIMIZATION]
43 p         <<< PREFETCH(HARD) Expected by compiler :
44 p             <<< (unknown)
45 p             <<< Loop-information End >>>
46 p             for(ii=0;ii<n;ii+=blki) {
47 p                 int min1=MIN(ii+blki,n);
48 p                 <<< Loop-information Start >>>
49 p                 <<< [OPTIMIZATION]
50 p                 <<< SIMD(VL: 8)
51 p                 <<< PREFETCH(HARD) Expected by compiler :
52 p                     <<< (unknown)
53 p                     <<< Loop-information End >>>
54 p                     for(i=ii;i<min1;i++) {
55 p                         s1 = s1 + a[j][i] * (s3 * b[j][i] + c[j][i] * (s2 + s3 * d[j][i]));
56 p                         int min2=MIN(ii+blki,n-100);
57 p                         <<< Loop-information Start >>>
58 p                         <<< [OPTIMIZATION]
59 p                         <<< SIMD(VL: 8)
60 p                         <<< SOFTWARE PIPELINING(IPC: 2.3)
61 p                         <<< PREFETCH(HARD) Expected by compiler :
62 p                             <<< (unknown)
63 p                             <<< Loop-information End >>>
64 p                             for(i=ii;i<min2;i++) {
65 p                                 e[j][i] = s2 * (a[j][i] + b[j][i] * (s3 + c[j][i] * d[j][i]));
66 p                                 }
67 p                         }
68 p                     }
69 p                 }
70 p             }
71 p         }
72 p     }
73 p }
```

ブロックサイズ : 4KB
4KB×4ストリーム=16KB
→ L1キャッシュに載るサイズ

配列アクセスがキャッシュヒットする



ループ 1

ループ 2

L1DミスとL2ミスが減少した

Cache	L1I miss rate (/Effective instruction)	Load-store instruction	L1D miss	改善前		改善後	
				L1D miss rate (/Load-store instruction)	L1D miss demand rate (%) (/L1D miss)	L1D miss rate (/Load-store instruction)	L1D miss demand rate (%) (/L1D miss)
改善前	0.00	2.17E+09	4.23E+08	0.19	9.14%	改善後	0.00
改善後	0.00	1.99E+09	2.35E+08	0.12	19.90%		

L2 miss	L2 miss rate (/Load-store instruction)	L2 miss demand rate (%) (/L2 miss)
4.23E+08	0.19	5.42%
2.35E+08	0.12	7.77%

ループブロッキング

- ループブロッキングとは
- ループブロッキング（改善前）
- ループブロッキングの効果（ソースチューニング）
- ハードウェアプリフェッチの副作用

ループブロッキングとは (1/3)

ループブロッキングとは、キャッシュ効率を向上させるために、ソースを指定したブロッキングサイズに分割して実行することです。

2次元またはそれ以上の次元におけるストリップマイニングと見なすことができます。

改善前ソース例

```
subroutine sub(a,b,m,n)
    integer n,m
    real*8 a(m,n),b(n,m)
    do j=1,m
        do i=1,n
            b(i,j)=a(j,i)
        enddo
    enddo
end subroutine
```

配列a : ストライドアクセス
配列b : 連続アクセス



改善後ソース例

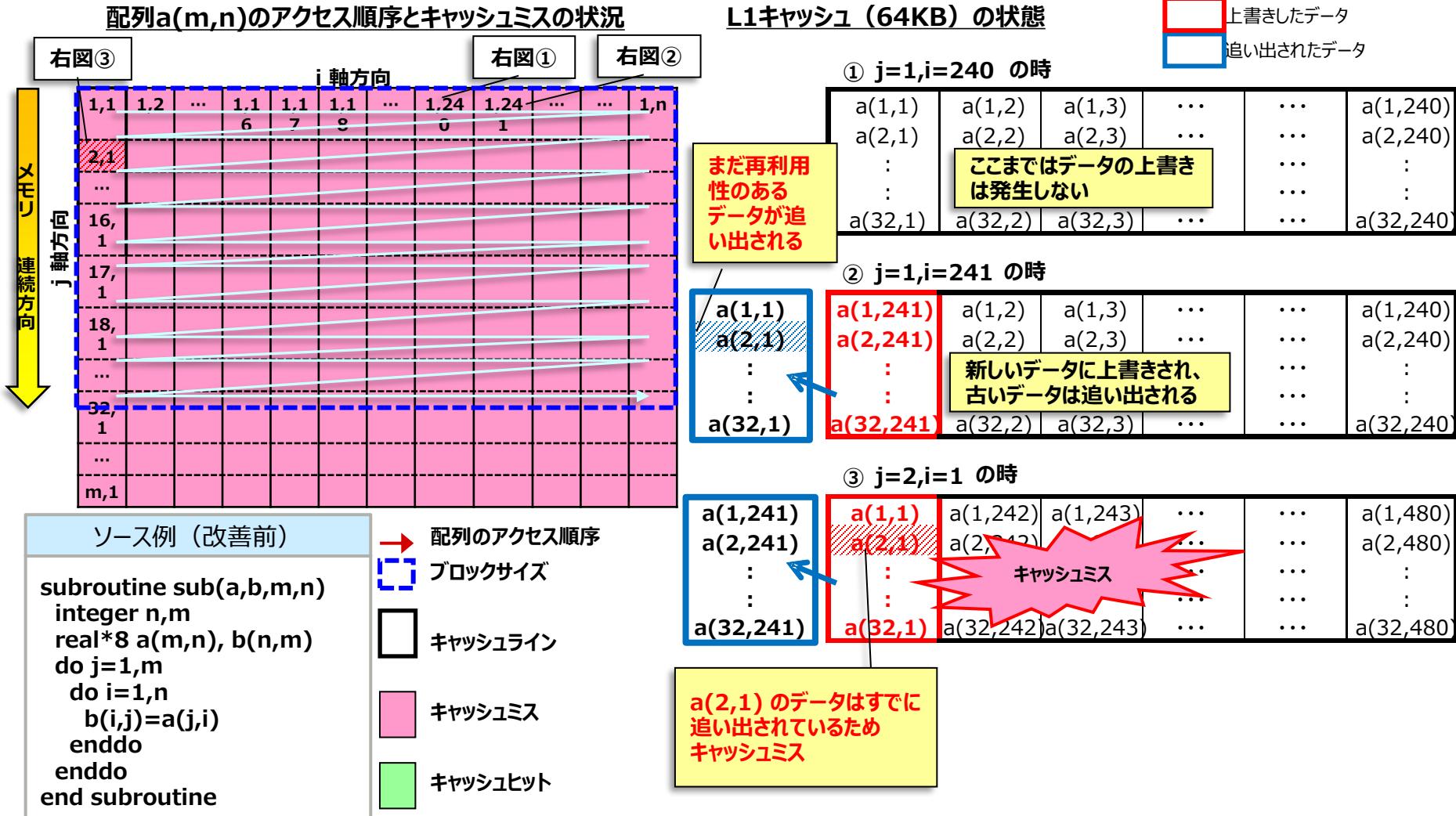
```
subroutine sub(a,b,m,n)
    parameter(blki=96,blkj=16)
    integer n,m
    real*8 a(m,n),b(n,m)
    do jj=1,m,blkj
        do ii=1,n,blki
            do j=jj,min(jj+blkj-1,m)
                do i=ii,min(ii+blki-1,n)
                    b(i,j)=a(j,i)
                enddo
            enddo
        enddo
    enddo
end subroutine
```

ブロックサイズ
1配列12Kbyte
(=96×16×8byte)

ループブロッキングとは (2/3)

● 配列アクセス (改善前)

配列aはストライドアクセスのため、iが更新されるたびにメモリアクセスが発生する。
その結果、a(1,1)でキャッシュに載せたデータは、a(2,1)でアクセスされる前に追い出されてしまう。



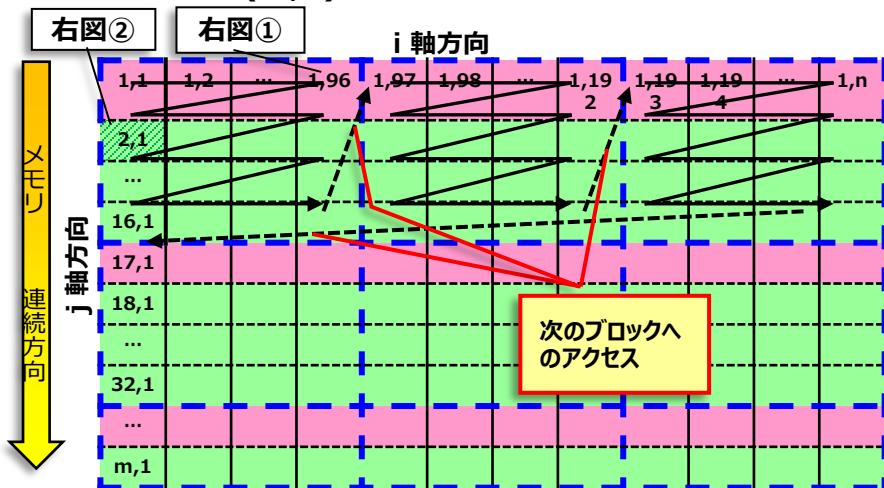
ループブロッキングとは (3/3)

- 配列アクセス（改善後） ブロックサイズ96×16

ループブロッキングにより1ブロックずつのアクセスを行う。

その結果、 $a(2,1)$ のアクセスがキャッシュにヒットするようになり、キャッシュ効率が向上する。

配列a(m,n)のアクセス順序とキャッシュミスの状況



ソース例（改善後）

```

subroutine sub(a,b,m,n)
parameter(blki=96,blkj=16)
integer n,m
real*8 a(m,n),b(n,m)
do jj=1,m, blkj
  do ii=1,n, blki
    do j=jj,min(jj+blkj-1,m)
      do i=ii,min(ii+blki-1,n)
        b(i,j)=a(j,i)
      enddo
    enddo
  enddo
end subroutine

```

- 配列のアクセス順序
- ブロックサイズ
- キャッシュライン
- キャッシュミス
- キャッシュヒット

L1キャッシュ（64KB）の状態

① $j=1, i=96$ の時

a(1,1)	a(1,2)	a(1,3)	...	a(1,96)	
a(2,1)	a(2,2)	a(2,3)	...	a(2,96)	
:	配列データが キャッシュに載る		:
:			:
a(32,1)	a(32,2)	a(32,3)	...	a(32,96)	

② $j=2, i=1$ の時

データがキャッシュに
残っているためキャッシュヒット

ポイント：
ループブロッキングを行うことで、データの連續性が損なわれ
ハードウェアプリフェッチが効かなくなることがあります。
その場合はソフトウェアプリフェッチを利用します。
詳しくは[ハードウェアプリフェッチの副作用](#)を参照してください。

配列aがストライドアクセスであるため、キャッシュの利用効率が悪くなり、浮動小数点ロードメモリアクセス待ちが多くなっています。

改善前ソース

```

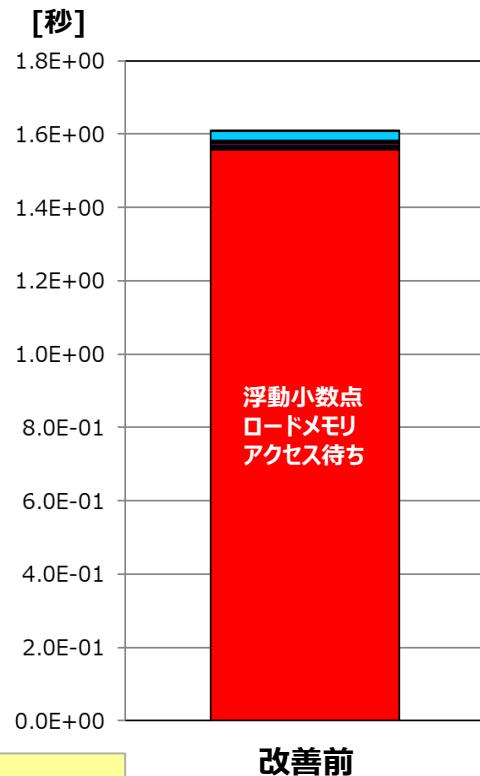
48   1      !$omp do
        <<< Loop-information Start >>>
        <<< [OPTIMIZATION]
        <<< PREFETCH(HARD) Expected by compiler :
        <<<     b
        <<< Loop-information End >>>
49   2 p      do j=1,n2
        <<< Loop-information Start >>>
        <<< [OPTIMIZATION]
        <<< SIMD(VL: 8)
        <<< SOFTWARE PIPELINING(IPC: 1.72, ITR: 176, MVE: 6, POL: S)
        <<< PREFETCH(HARD) Expected by compiler :
        <<<     b
        <<< Loop-information End >>>
50   3 p 2v    do i=1,n1
51   3 p 2v      b(i,j) = c0 + a(j,i)*(c1 + a(j,i)*(c2 + a(j,i)*(c3 + a(j,i)*
52   3           & (c4 + a(j,i)*(c5 + a(j,i)*(c6 + a(j,i)*(c7 + a(j,i)*
53   3           & (c8 + a(j,i)*c9)))))))
54   3 p 2v    enddo
55   2 p    enddo
56   1      !$omp enddo

```

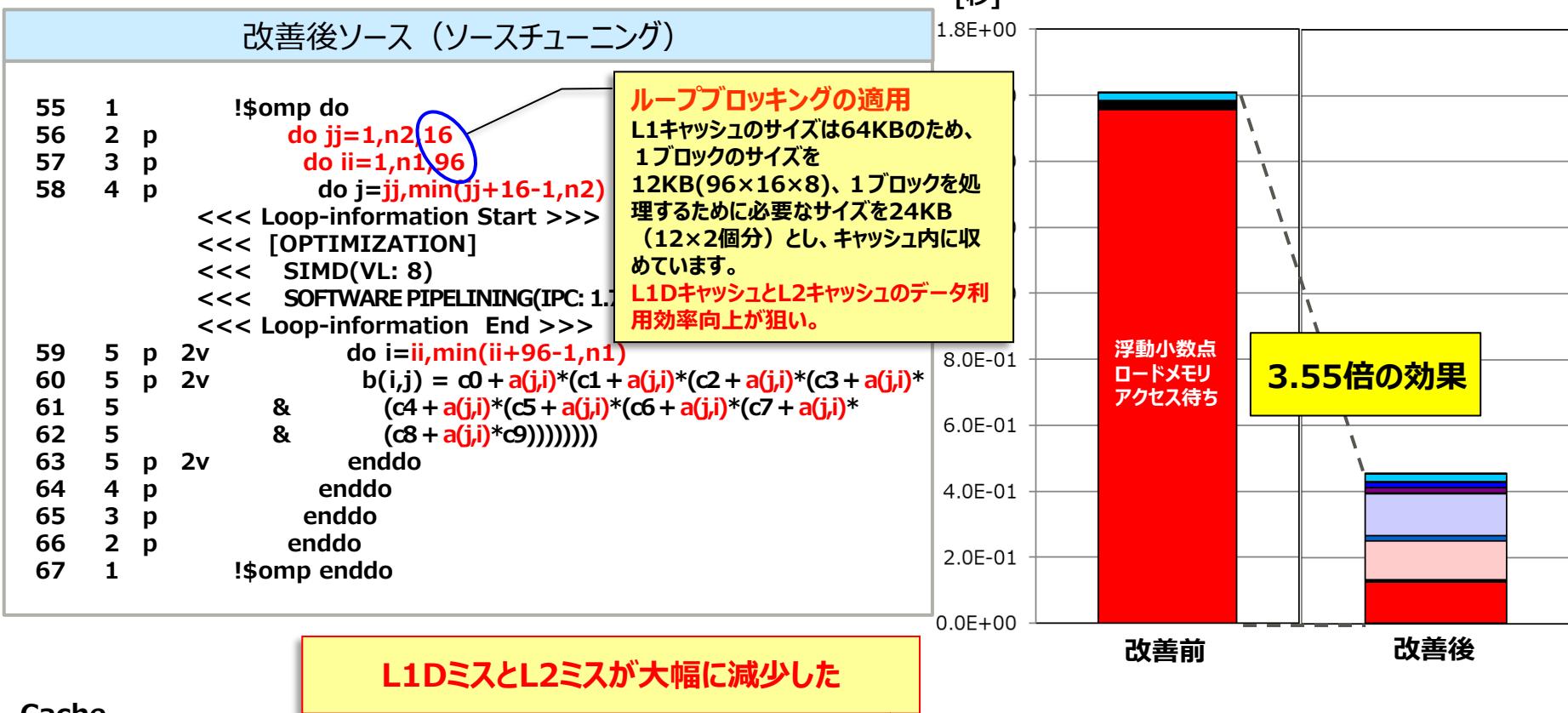
配列a のデータは、i のイタレーション時に一度
L1Dキャッシュに載るが、次の j のイタレーション
時にはデータが追い出されてしまっているた
め、キャッシュミスとなる

Cache

	L1I miss rate (/Effective instruction)	Load-store instruction	L1D miss	L1D miss rate (/Load-store instruction)	L1D miss demand rate (%) (/L1D miss)	L1D miss hardware prefetch rate (%) (/L1D miss)	L1D miss software prefetch rate (%) (/L1D miss)	L2 miss	L2 miss rate (/Load-store instruction)	L2 miss demand rate (%) (/L2 miss)	L2 miss hardware prefetch rate (%) (/L2 miss)	L2 miss software prefetch rate (%) (/L2 miss)
改善前	0.00	4.06E+08	1.28E+09	3.16	98.85%	1.15%	0.00%	1.31E+09	3.24	29.75%	71.85%	0.00%



ループブロッキングにより配列aのデータが再利用されることで、キャッシュ効率が向上し、浮動小数点ロードメモリアクセス待ちが改善されました。



Cache

	L1 miss rate (/Effective instruction)	Load-store instruction	L1D miss	L1D miss rate (/Load-store instruction)	L1D miss demand rate (%) (/L1D miss)	L1D miss hardware prefetch rate (%) (/L1D miss)	L1D miss software prefetch rate (%) (/L1D miss)	L2 miss	L2 miss rate (/Load-store instruction)	L2 miss demand rate (%) (/L2 miss)	L2 miss hardware prefetch rate (%) (/L2 miss)	L2 miss software prefetch rate (%) (/L2 miss)
改善前	0.00	4.06E+08	1.28E+09	3.16	98.85%	1.15%	0.00%	1.31E+09	3.24	29.75%	71.85%	0.00%
改善後	0.00	8.26E+08	1.69E+08	0.20	95.18%	4.82%	0.00%	1.56E+08	0.19	45.06%	61.56%	0.00%

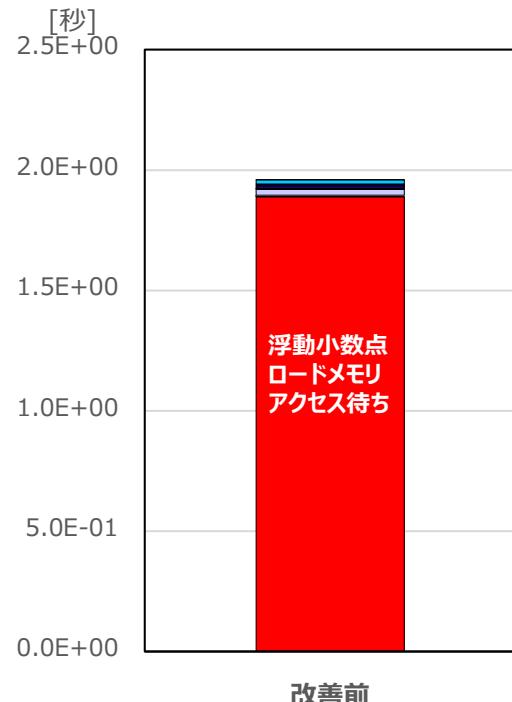
配列aがストライドアクセスであるため、キャッシュの利用効率が悪くなり、浮動小数点ロードメモリアクセス待ちが多くなっています。

改善前ソース

```

43     #pragma omp parallel
44     {
45         for(k=0;k<iter;k++) {
46             #pragma omp for
47             p           <<< Loop-information Start >>>
48             p           <<< [OPTIMIZATION]
49             p           <<< PREFETCH(HARD) Expected by compiler :
50             p           <<< (unknown)
51             p           <<< Loop-information End >>>
52             p           <<< Loop-information Start >>>
53             p           <<< [OPTIMIZATION]
54             p           <<< SIMD(VL: 8)
55             p           <<< SOFTWARE PIPELINING(IPC: 1.08, ITR: 128, MVE: 4, POL: S)
56             p           <<< PREFETCH(HARD) Expected by compiler :
57             p           <<< (unknown)
58             p           <<< Loop-information End >>>
59             p           2v           for(i=0;i<n1;i++) {
60             p           2v               b[j][i] = c0 + a[i][j]*(c1 + a[i][j]*(c2 + a[i][j]*(c3 + a[i][j]*
61             p           2v                   (c4 + a[i][j]*(c5 + a[i][j]*(c6 + a[i][j]*(c7 + a[i][j]*
62             p           2v                   (c8 + a[i][j]*c9))))));
63             p           }
64             p           }
65             p           }
66             p           }

```



配列a のデータは、i のイタレーション時に一度 L1Dキャッシュに載るが、次の j のイタレーション時にはデータが追い出されてしまっているため、キャッシュミスとなる

Cache	L1I miss rate (/Effective instruction)	Load-store instruction	L1D miss	L1D miss rate (/Load-store instruction)	L1D miss demand rate (%) (/L1D miss)	L1D miss hardware prefetch rate (%) (/L1D miss)	L1D miss software prefetch rate (%) (/L1D miss)	L2 miss	L2 miss rate (/Load-store instruction)	L2 miss demand rate (%) (/L2 miss)	L2 miss hardware prefetch rate (%) (/L2 miss)	L2 miss software prefetch rate (%) (/L2 miss)
改善前	0.00	4.00E+08	1.28E+09	3.21	98.53%	1.47%	0.00%	1.33E+09	3.32	28.24%	73.08%	0.00%

ループブロッキングにより配列aのデータが再利用されることで、キャッシュ効率が向上し、浮動小数点ロードメモリアクセス待ちが改善されました。

改善後ソース（ソースチューニング）

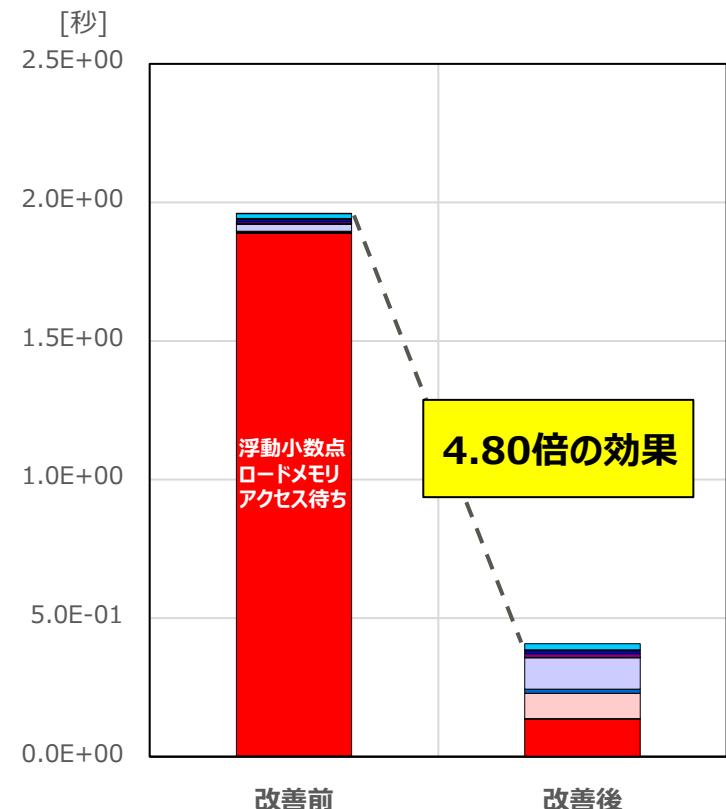
```

45      #pragma omp parallel
46      {
47          for(k=0;k<iter;k++) {
48              #pragma omp for
49              p
50              p
51          p
52          p
53          p
54          p
55          p
56          p
57          p
58          p
59          p
60          p
61          p
62      }

```

ループブロッキングの適用
L1キャッシュのサイズは64KBのため、1ブロックのサイズを12KB(96×16×8)、1ブロックを処理するために必要なサイズを24KB(12×2個分)とし、キャッシュ内に収めています。
L1DキャッシュとL2キャッシュのデータ利用効率向上が狙い。

L1DミスとL2ミスが大幅に減少した



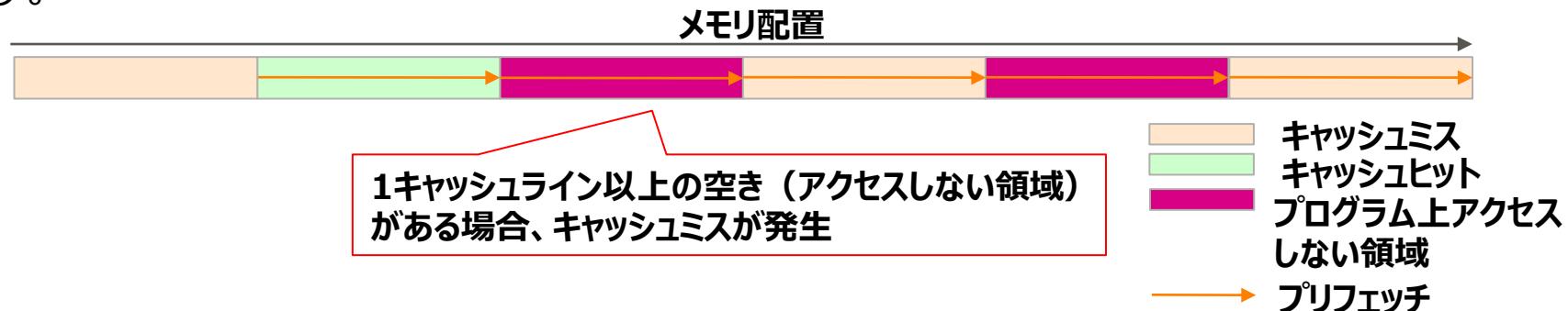
Cache	L1 miss rate (/Effective instruction)	Load-store instruction	L1D miss	L1D miss rate (/Load-store instruction)	L1D miss demand rate (%) (/L1D miss)	L1D miss hardware prefetch rate (%) (/L1D miss)	L1D miss software prefetch rate (%) (/L1D miss)	L2 miss	L2 miss rate (/Load-store instruction)	L2 miss demand rate (%) (/L2 miss)	L2 miss hardware prefetch rate (%) (/L2 miss)	L2 miss software prefetch rate (%) (/L2 miss)
改善前	0.00	4.00E+08	1.28E+09	3.21	98.53%	1.47%	0.00%	1.33E+09	3.32	28.24%	73.08%	0.00%
改善後	0.00	4.93E+08	1.70E+08	0.35	93.66%	6.34%	0.00%	1.40E+08	0.28	49.47%	53.12%	0.00%

- ループブロッキング、外側ループプリフェッチなどでは、ブロックサイズが小さくなることでデータが連続的でなくなり、ハードウェアプリフェッチが有効にならない場合があります。その場合はソフトウェアプリフェッチを利用する必要があります。
- ソフトウェアプリフェッチの指定方法は[ソフトウェアプリフェッチの利用](#)を参照してください。

● ハードウェアプリフェッチ

プログラムのメモリアクセスの規則性からハードウェアがデータアクセスを予測しプリフェッチします。

データ間に1キャッシュライン以上の空きがある場合はプリフェッチが失敗する可能性があります。



● ソフトウェアプリフェッチ

ソフトウェア（コンパイラ）がプログラムを解析し、prefetch命令を生成することでプリフェッチします。または、指示行指定で該当箇所にプリフェッチ命令を生成します。

セクタキヤッシュ

- ・ セクタキヤッシュとは
- ・ セクタキヤッシュの使用方法
- ・ セクタキヤッシュ：事例 1 （改善前）
- ・ セクタキヤッシュ：事例 1 （ソースチューニング）
- ・ セクタキヤッシュ：事例 2 （改善前）
- ・ セクタキヤッシュ：事例 2 （ソースチューニング）

セクタキャッシュとは

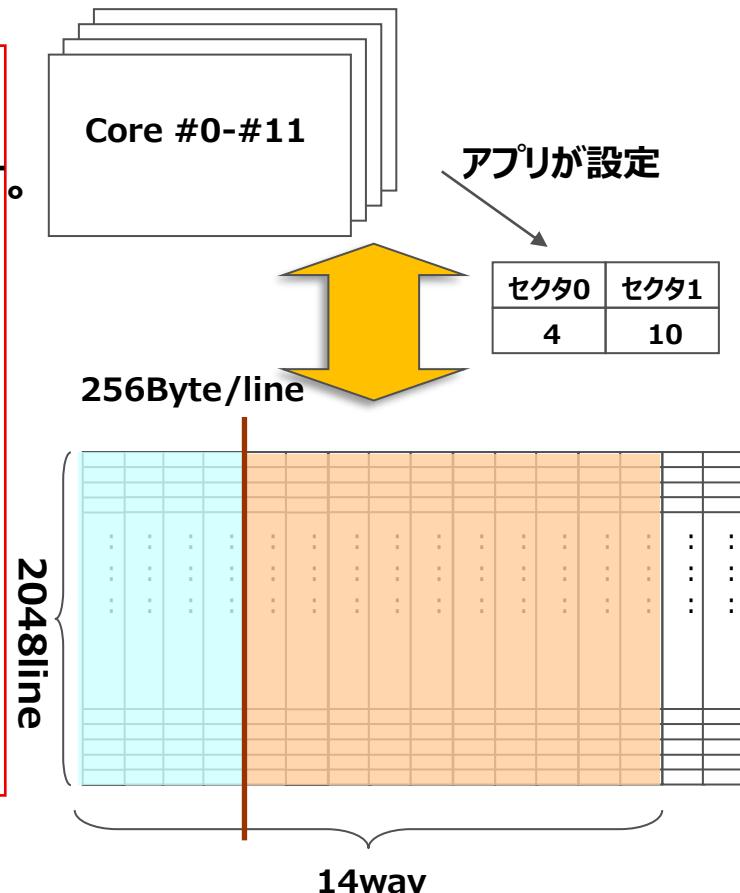
セクタキャッシュとは、再利用性のあるデータが、再利用性の無いデータによってキャッシュから追い出されることを防ぐことができるキャッシュ機構です。アプリが再利用性のあるデータと再利用性の無いデータをセクタ毎に分けて配置することができます。（再利用する配列はセクタ1を使用し、その他はセクタ0を使用）

L2キャッシュでの利用イメージ

■ セクタキャッシュの詳細

- L1Dキャッシュ・L2キャッシュいずれにも複数個のセクタを設定できます。セクタの最大数はL1Dが4(※1), L2は2です。
- 各セクタの容量はWAY数で指定します。
- 容量は目標値として働きます。
ハードは、ラインリプレース時に各セクタが指定された容量に近づくように制御します。
→容量オーバーしていても強制的に無効化しない
- セクタ内の追い出しはLRU（最近最も使われていないデータを最初に捨てる）アルゴリズムで制御します。
- セクタ0、1の用途はアプリケーションで決めることが可能です。
ただし、命令列はセクタ0に格納されます。
- 2次キャッシュはアシスタントコアが常時2wayを使用します。

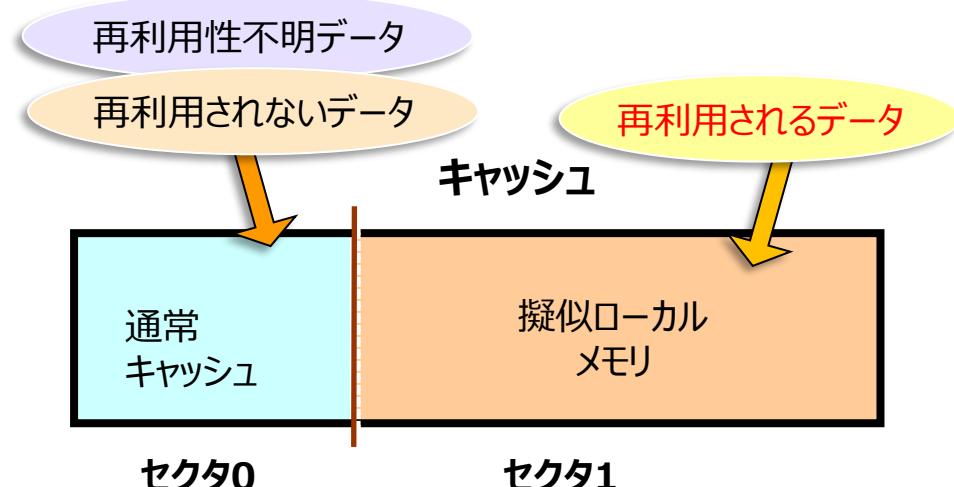
(※1)現在仕様ではL1Dキャッシュのセクタ数は2です。お客様が簡単に使用できることを優先しこの仕様としました。



■ セクタキャッシュ：疑似ローカルメモリ

ソフトウェアが、データの再利用性に応じて
セクタを使い分けることが可能

- 再利用するデータ → セクタ1を使用
- その他のデータ → セクタ0を使用
- セクタ1上のデータは、他のデータによって
追い出されがない
- 指示行でセクタ1に載せる配列を指定できる



セクタキャッシュ指定のコンパイラ指示行の使用例

```

!OCL SCACHE_ISOLATE_WAY(L2=10)
!OCL SCACHE_ISOLATE_ASSIGN(a)
do j=1,m
  do i=1,n
    a(i) = a(i) + b(i,j)*c(i,j)
  enddo
enddo
!OCL END_SCACHE_ISOLATE_ASSIGN
!OCL END_SCACHE_ISOLATE_WAY

```

旧仕様の指定方法（京, FX100）

```

!OCL CACHE_SECTOR_SIZE(4,10)
!OCL CACHE_SUBSECTOR_ASSIGN(a)
do j=1,m
  do i=1,n
    a(i) = a(i) + b(i,j)*c(i,j)
  enddo
enddo
!OCL END_CACHE_SUBSECTOR
!OCL END_CACHE_SECTOR_SIZE

```

<狙い>

ループ内で配列 b と配列 c のアクセスによって
再利用性のある配列 a がキャッシュから追い出されないようにする

セクタキャッシュの使用方法（2/2）

セクタキャッシュを使用する場合は、以下の最適化制御行を指定します。

最適化指示子 (Fortran)	意味	指定可能な最適化制御行			
		プログラム単位	DOループ単位	文単位	配列代入文単位
SCACHE_ISOLATE_WAY(L2=n1[,L1=n2]) END_SCACHE_ISOLATE_WAY	1次キャッシュと2次キャッシュのセクタ1の最大way数を指示します。	○	×	○	×
SCACHE_ISOLATE_ASSIGN(array1[,array2]…) END_SCACHE_ISOLATE_ASSIGN	キャッシュのセクタ1に載せる配列を指示します。	○	×	○	×
最適化指示子 (C/C++)	意味	指定可能な最適化制御行			
		global行	procedure行	loop行	statement行
scache_isolate_way(L2=n1[,L1=n2]) end_scache_isolate_way	1次キャッシュと2次キャッシュのセクタ1の最大way数を指示します。	×	○	×	○
scache_isolate_assign(array1[,array2]…) end_scache_isolate_assign	キャッシュのセクタ1に載せる配列を指示します。	×	○	×	○

■ 注意事項

- 2次キャッシュはアシスタントコアが常時2wayを使用します。そのため、n1およびn2に指定できる範囲は以下のようになります。
 $0 \leq n1 \leq \text{"2次キャッシュの最大way数 - 2"}$
 $0 \leq n2 \leq \text{"1次キャッシュの最大way数"}$
- アシスタントコアを含むCMGでは、L2キャッシュの一部(2way=1MiB分)をアシスタントコア用に使用します。そのため、アシスタントコアを含むCMGでは、**2次キャッシュの最大way数は14, サイズは7MiBとなります。**
- clangモードではセクタキャッシュ機能を使用できません

A64FX諸元	
CMG数	4
L1Iキャッシュ・サイズ	64KiB / 4way
L1Dキャッシュ・サイズ	64KiB / 4way
L2キャッシュ・サイズ	32MiB / 16way (8MiB / CMG)

配列bのデータはキャッシュから追い出されてしまうため、再利用できません。
そのため浮動小数点ロードL2アクセス待ちが多くなっています。

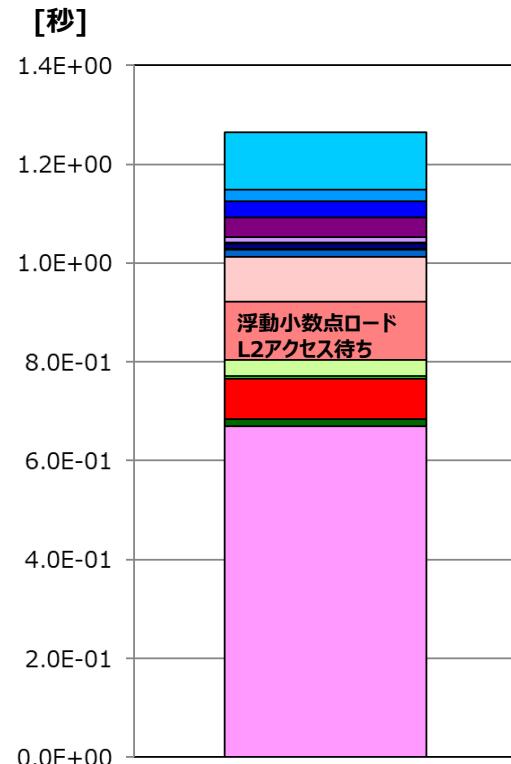
改善前ソース

```

66      parameter(n=8*1024*1024, m=9*512*1024/8)
67      real*8   a(n), b(m), s
68      integer*8 c(n)
69      real*8   dummy1(140),dummy2(140)
70      common /data/a,dummy1,c,dummy2,b
71
    <<< Loop-information Start >>>
    <<< [PARALLELIZATION]
    <<< Standard iteration count: 843
    <<< [OPTIMIZATION]
    <<< SIMD(VL: 8)
    <<< SOFTWARE PIPELINING(IPC: 2.66, ITR: 176, MVE: 4, POL: S)
    <<< PREFETCH(HARD) Expected by compiler :
    <<<   c, a
    <<< Loop-information End >>>
72      1 pp 2v      do i=1,n
73      1 p 2v      a(i) = a(i) + s * b(c(i))
74      1 p 2v      enddo

```

配列サイズ
a: 64MiB
b: 4.5MiB
c: 64MiB



改善前

Cache	L1 miss rate (/Effective instruction)	Load-store instruction	L1D miss	L1D miss rate (/Load-store instruction)	L1D miss demand rate (%) (/L1D miss)	L1D miss hardware prefetch rate (%) (/L1D miss)	L1D miss software prefetch rate (%) (/L1D miss)	L2 miss	L2 miss rate (/Load-store instruction)	L2 miss demand rate (%) (/L2 miss)	L2 miss hardware prefetch rate (%) (/L2 miss)	L2 miss software prefetch rate (%) (/L2 miss)
改善前	0.00	4.76E+09	7.89E+08	0.17	0.89%	99.11%	0.00%	7.34E+08	0.15	0.77%	100.00%	0.00%
改善前	Memory throughput (GB/s)											
改善前	203.11											

メモリアクセス負荷が高い

L2キャッシュミス率が高い

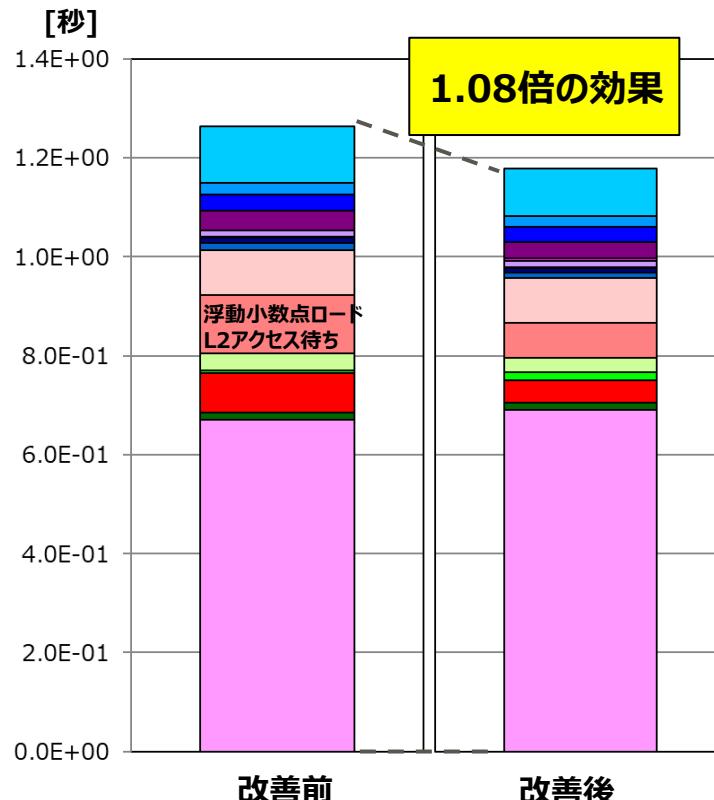
配列bをセクタ1に載せることでキヤッショ効率が向上し、浮動小数点ロードL2アクセス待ちが改善されました。

改善後ソース (最適化制御行チューニング)

```

58      parameter(n=8*1024*1024, m=9*512*1024/8)
59      real*8  a(n), b(m), s
60      integer*8 c(n)
61      real*8  dummy1(140),dummy2(140)
62      common /data/a,dummy1,c,dummy2,b
63
64      !OCL SCACHE_ISOLATE_WAY(L2=10)
65      !OCL SCACHE_ISOLATE_ASSIGN(b)
<<< Loop-information Start >>>
<<< [PARALLELIZATION]
<<< Standard iteration count: 843
<<< [OPTIMIZATION]
<<< SIMD(VL: 8)
<<< SOFTWARE PIPELINING(IPC: 2.66, ITR: 176, MVE: 4, POL: S)
<<< PREFETCH(HARD) Expected by compiler :
<<<   c, a
<<< Loop-information End >>>
66  1 pp 2v      do i=1,n
67  1 p 2v      a(i) = a(i) + s * b(c(i))
68  1 p 2v      enddo
69      !OCL END_SCACHE_ISOLATE_ASSIGN
70      !OCL END_SCACHE_ISOLATE_WAY

```



	L1I miss rate (/Effective instruction)	Load-store instruction	L1D miss	L1D miss rate (/Load-store instruction)	L1D miss demand rate (%) (/L1D miss)	L1D miss hardware prefetch rate (%) (/L1Dmiss)	L1D miss software prefetch rate (%) (/L1Dmiss)	L2 miss	L2 miss rate (/Load-store instruction)	L2 miss demand rate (%) (/L2 miss)	L2 miss hardware prefetch rate (%) (/L2miss)	L2 miss software prefetch rate (%) (/L2miss)
改善前	0.00	4.76E+09	7.89E+08	0.17	0.89%	99.11%	0.00%	7.34E+08	0.15	0.77%	100.00%	0.00%
改善後	0.00	5.19E+09	7.93E+08	0.15	1.19%	98.81%	0.01%	5.99E+08	0.12	1.93%	99.69%	0.00%

	Memory throughput (GB/s)
改善前	203.11
改善後	188.07

L2ミスが減少した

配列b のデータはキャッシュから追い出されてしまうため、再利用できません。
そのため浮動小数点ロードL2アクセス待ちが多くなっています。

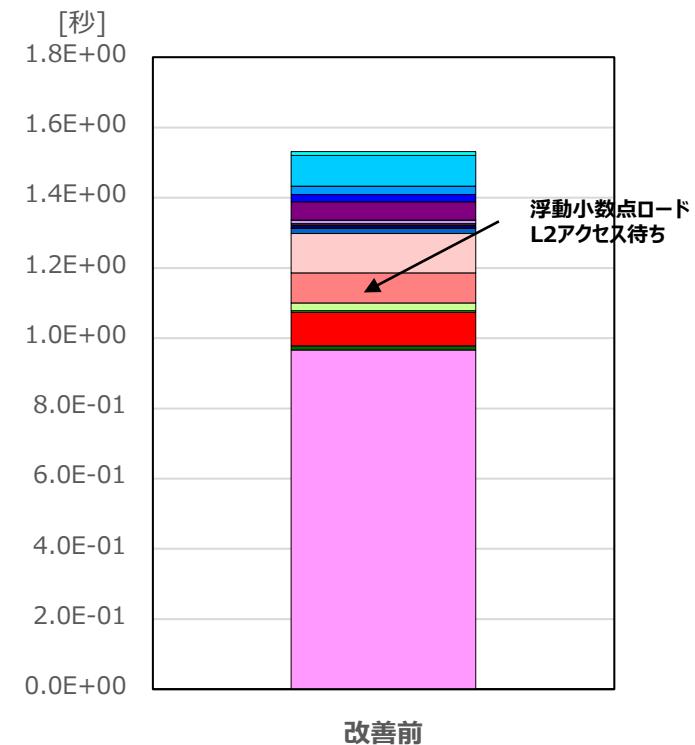
改善前ソース

```

56     void sub(double s) {
57         long long int i;
58
59         #pragma omp parallel for
60         <<< Loop-information Start >>>
61         <<< [OPTIMIZATION]
62         <<< SIMD(VL: 8)
63         <<< SOFTWARE PIPELINING(IPC: 2.33, ITR: 160,
64             MVE: 4, POL: S)
65         <<< PREFETCH(HARD) Expected by compiler :
66             <<< c, a
67             <<< Loop-information End >>>
68
69         p 2v for(i=0;i<n;i++) {
70             p 2v     a[i] = a[i] + s * b[c[i]];
71             p 2v }
72         }

```

配列宣言：サイズ
 double a[8388608]: 64MiB
 double b[589824]: 4.5MiB
 long long int c[8388608]: 64MiB



改善前

Cache	L1I miss rate (/Effective instruction)	Load-store instruction	L1D miss	L1D miss rate (/Load-store instruction)	L1D miss demand rate (%) (/L1D miss)	L1D miss hardware prefetch rate (%) (/L1D miss)	L1D miss software prefetch rate (%) (/L1D miss)	L2 miss	L2 miss rate (/Load-store instruction)	L2 miss demand rate (%) (/L2 miss)	L2 miss hardware prefetch rate (%) (/L2 miss)	L2 miss software prefetch rate (%) (/L2 miss)
改善前	0.00	4.85E+09	7.89E+08	0.16	1.07%	98.92%	0.01%	7.39E+08	0.15	0.78%	100.00%	0.00%

Statistics	Memory throughput (GB/s)
改善前	167.47

メモリアクセス負荷が高い

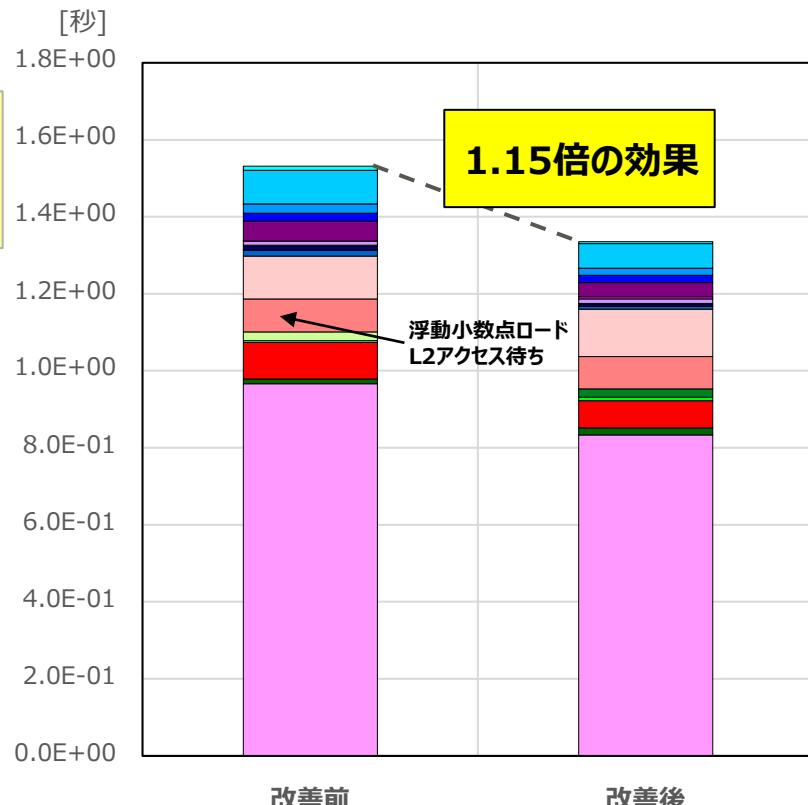
L2キャッシュミス率が高い

配列bをセクタ1に載せることでキャッシュ効率が向上し、浮動小数点ロードL2アクセス待ちが改善されました。

改善後ソース（最適化制御行チューニング）

```

56     void sub(double s){          配列宣言：サイズ
57       long long int i;          double a[8388608]: 64MiB
58       #pragma statement scache_isolate_way L2=10
59       #pragma statement scache_isolate_assign b
60       #pragma omp parallel for
<<< Loop-information Start >>>
<<< [OPTIMIZATION]
<<< SIMD(VL: 8)
<<< SOFTWARE PIPELINING(IPC: 2.33, ITR: 160,
                           MVE: 4, POL: S)
<<< PREFETCH(HARD) Expected by compiler :
<<<   c, a
<<< Loop-information End >>>
61   p 2v  for(i=0;i<n;i++) {
62   p 2v    a[i] = a[i] + s * b[c[i]];
63   p 2v  }
64   #pragma statement end_scache_isolate_assign
65   #pragma statement end_scache_isolate_way
66 }
```



Cache	L1I miss rate (/Effective instruction)	Load-store instruction	L1D miss	L1D miss rate (/Load-store instruction)	L1D miss demand rate (%) (/L1D miss)	L1D miss hardware prefetch rate (%) (/L1D miss)	L1D miss software prefetch rate (%) (/L1D miss)	L2 miss	L2 miss rate (/Load-store instruction)	L2 miss demand rate (%) (/L2 miss)	L2 miss hardware prefetch rate (%) (/L2 miss)	L2 miss software prefetch rate (%) (/L2 miss)
改善前	0.00	4.85E+09	7.89E+08	0.16	1.07%	98.92%	0.01%	7.39E+08	0.15	0.78%	100.00%	0.00%
改善後	0.00	5.03E+09	7.91E+08	0.16	1.23%	98.78%	0.00%	5.48E+08	0.11	1.99%	98.64%	0.00%

Statistics	Memory throughput (GB/s)
改善前	167.47
改善後	155.55

L2ミスが減少した

配列uのデータはキャッシュから追い出されてしまうため、再利用できません。
そのため、メモリ・キャッシュビジー待ちが多くなっています。

改善前ソース

```

167 1 s      do iter = 1, niter
    <<< Loop-information Start >>>
    <<< [PARALLELIZATION]
    <<< Standard iteration cou
    <<< Loop-information End >
168 2 pp      do k=1,n3-2
    <<< Loop-information Start >
    <<< [OPTIMIZATION]
    <<< PREFETCH(HARD) Expected by compiler :
    <<< u, rhs, unew
    <<< Loop-information End >
169 3 p      do j=1,n2-2
    <<< Loop-information Start >
    <<< [OPTIMIZATION]
    <<< SIMD(VL: 8)
    <<< SOFTWARE PIPELINING(IPC: 3.71, ITR: 136, MVE: 9, POL: S)
    <<< PREFETCH(HARD) Expected by compiler :
    <<< u, rhs, unew
    <<< Loop-information End >>>
170 4 p v      do i=1,n1-2
171 4 p v          unew(i,j,k) = &
172 4              ((u(i+1,j,k) + u(i-1,j,k)) * h1sqinv &
173 4                  +(u(i,j+1,k) + u(i,j-1,k)) * h2sqinv &
174 4                  +(u(i,j,k+1) + u(i,j,k-1)) * h3sqinv &
175 4                  -rhs(i,j,k)) * hhhinv
176 4 p v      end do
177 3 p      end do
178 2 p      end do
179 1       end do

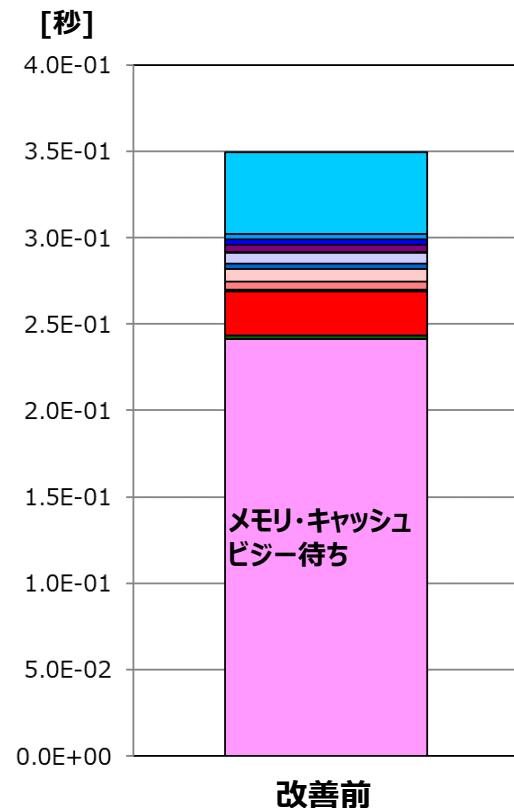
```

n1=452
n2=52
n3=322

配列サイズ
unew: 60.5MB
u: 60.5MB
rhs: 60.5MB

配列uのi,j次元に再利用性
があるため、配列uをオン
キャッシュにしたい。

メモリアクセス負荷が高い



改善前

	Memory throughput (GB/s)
改善前	215.62

Cache	L2 miss	L2 miss rate (/Load-store instruction)	L2 miss demand rate (%) (/L2 miss)	L2 miss hardware prefetch rate (%) (/L2 miss)	L2 miss software prefetch rate (%) (/L2 miss)
改善前	2.46E+08	0.15	2.57%	98.19%	0.00%

配列uのk次元の一部をセクタ1に載せることで、キャッシュ効率が向上し、メモリ・キャッシュビジー待ちが改善されました。

改善後ソース

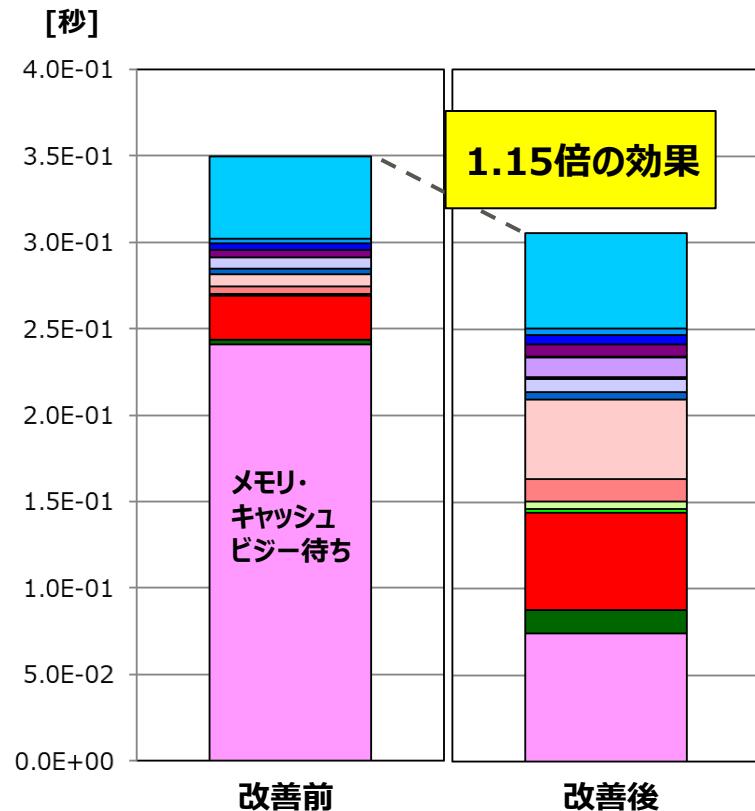
```

166      !OCL SCACHE_ISOLATE_WAY(L2=13)
167      !OCL SCACHE_ISOLATE_ASSIGN(u)
168 1 s      do iter = 1, niter
    <<< Loop-information Start >>>
    <<< [PARALLELIZATION]
    <<< Standard iteration count: 2
    <<< Loop-information End >>>
169 2 pp      do k=1,n3-2
    <<< Loop-information Start >>>
    <<< [OPTIMIZATION]
    <<< PREFETCH(HARD) Expected by compiler :
    <<< u, rhs, unew
    <<< Loop-information End >>>
170 3 p      do j=1,n2-2
    <<< Loop-information Start >>>
    <<< [OPTIMIZATION]
    <<< SIMD(VL: 8)
    <<< SOFTWARE PIPELINING(IPC: 3.71, ITR: 136, MVE: 9, POL: S)
    <<< PREFETCH(HARD) Expected by compiler :
    <<< u, rhs, unew
    <<< Loop-information End >>>
171 4 p v      do i=1,n1-2
172 4 p v          unew(i,j,k) = &
173 4          ((u(i+1,j,k) + u(i-1,j,k)) * h1sqinv &
174 4          +(u(i,j+1,k) + u(i,j-1,k)) * h2sqinv &
175 4          +(u(i,j,k+1) + u(i,j,k-1)) * h3sqinv &
176 4          -rhs(i,j,k)) * hhhinv
177 4 p v      end do
178 3 p      end do
179 2 p      end do
180 1       end do
181      !OCL END_SCACHE_ISOLATE_ASSIGN
182      !OCL END_SCACHE_ISOLATE_WAY

```

配列uの再利用性
が向上

	L2 miss	L2 miss rate (/Load-store instruction)	L2 miss demand rate (%) (/L2 miss)	L2 miss hardware prefetch rate (%) (/L2 miss)	L2 miss software prefetch rate (%) (/L2 miss)
改善前	2.46E+08	0.15	2.57%	98.19%	0.00%
改善後	1.95E+08	0.10	13.43%	87.39%	0.00%



L2ミスが減少

	Memory throughput (GB/s)
改善前	215.62
改善後	205.39

配列uのデータはキャッシュから追い出されてしまうため、再利用できません。
そのため、メモリ・キャッシュビジー待ちが多くなっています。

改善前ソース

```

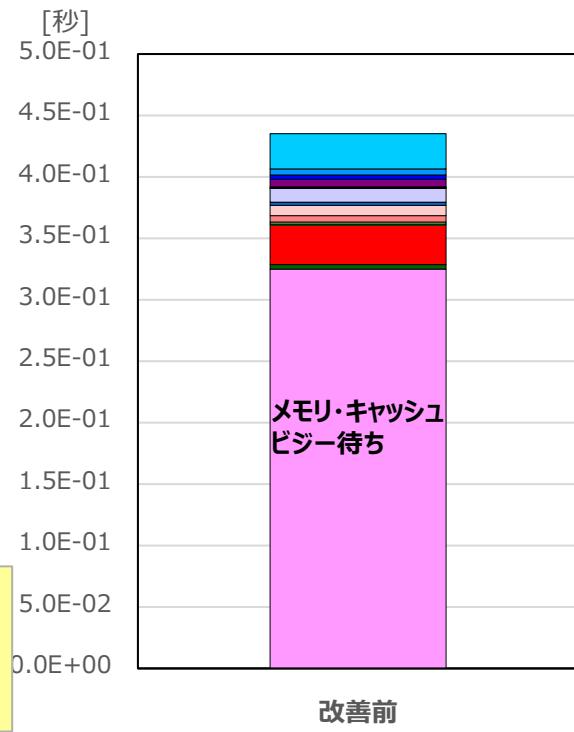
107     for (iter=0; iter<niter; iter++){
108         #pragma omp parallel for private(i,j,k)
109         p         for(k=1;k<=n3-2; k++) {
110             <<< Loop-information Start >>>
111             <<< [OPTIMIZATION]
112             <<< PREFETCH(HARD) Expected by compiler :
113                 (unknown)
114             <<< Loop-information End >>>
115             for(j=1;j<=n2-2;j++) {
116                 <<< Loop-information Start >>>
117                 <<< [OPTIMIZATION]
118                 <<< SIMD(VL: 8)
119                 <<< SOFTWARE PIPELINING(IPC: 2.12, ITR: 120, MVE: 8, POL: S)
120                 <<< PREFETCH(HARD) Expected by compiler :
121                     (unknown)
122                     <<< Loop-information End >>>
123                     p         v         for(i=1;i<=n1-2;i++) {
124                         p         v         unew[k][j][i] =
125                             ((u[k][j][i+1] + u[k][j][i-1]) * h1sqinv
126                             +(u[k][j+1][i] + u[k][j-1][i]) * h2sqinv
127                             +(u[k+1][j][i] + u[k-1][j][i]) * h3sqinv
128                             -rhs[k][j][i]) * hhhinv;
129                     p         v         }
130                 }
131             }
132         }
133     }

```

n1=452
n2=52
n3=322

配列サイズ
unew: 60.5MB
u: 60.5MB
rhs: 60.5MB

配列uのi,j次元に再利用性
があるため、配列uをオン
キャッシュにしたい。



メモリアクセス負荷が高い

Cache	L2 miss	L2 miss rate (/Load-store instruction)	L2 miss demand rate (%) (/L2 miss)	L2 miss hardware prefetch rate (%) (/L2 miss)	L2 miss software prefetch rate (%) (/L2 miss)
改善前	2.46E+08	0.15	2.38%	98.32%	0.00%

配列uのk次元の一部をセクタ1に載せることで、キャッシュ効率が向上し、メモリ・キャッシュビギー待ちが改善されました。

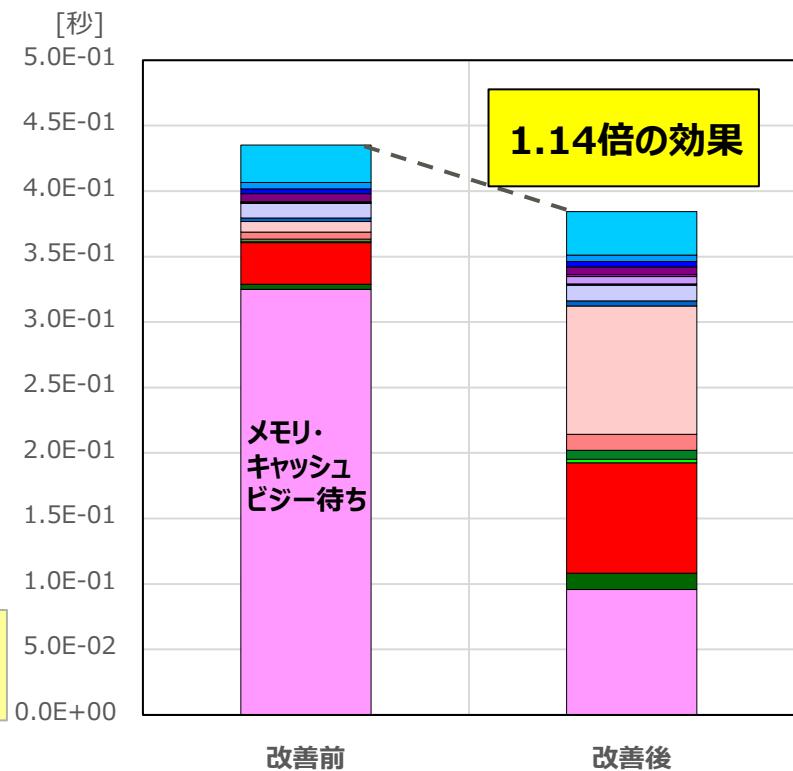
改善後ソース

```

107 #pragma statement scache_isolate_way L2=13
108 #pragma statement scache_isolate_assign u
109   for (iter=0; iter<niter; iter++){
110     #pragma omp parallel for private(i,j,k)
111     p       for(k=1;k<=n3-2; k++){
112       <<< Loop-information Start >>>
113       <<< [OPTIMIZATION]
114       <<< PREFETCH(HARD) Expected by compiler :
115       <<< (unknown)
116       <<< Loop-information End >>>
117       p       for(j=1;j<=n2-2;j++){
118       <<< Loop-information Start >>>
119       <<< [OPTIMIZATION]
120       <<< SIMD(VL: 8)
121       <<< SOFTWARE PIPELINING(IPC: 2.12, ITR: 120, MVE: 8, POL: S)
122       <<< PREFETCH(HARD) Expected by compiler :
123       <<< (unknown)
124       <<< Loop-information End >>>
125     p       v       for(i=1;i<=n1-2;i++){
126       v       unew[k][j][i] =
127         ((u[k][j][i+1] + u[k][j][i-1]) * h1sqinv
128         +(u[k][j+1][i] + u[k][j-1][i]) * h2sqinv
129         +(u[k+1][j][i] + u[k-1][j][i]) * h3sqinv
130         -rhs[k][j][i]) * hhhinv;
131     p       v       }
132   p       }
133   }
134 #pragma statement end_scache_isolate_assign
135 #pragma statement end_scache_isolate_way

```

配列uの再利用性
が向上



L2ミスが減少

Cache	L2 miss	L2 miss rate (/Load-store instruction)	L2 miss demand rate (%) (/L2 miss)	L2 miss hardware prefetch rate (%) (/L2 miss)	L2 miss software prefetch rate (%) (/L2 miss)
改善前	2.46E+08	0.15	2.38%	98.32%	0.00%
改善後	1.99E+08	0.12	13.90%	86.99%	0.00%

ループ交換

- ループ交換とは
- ループ交換（改善前）
- ループ交換のチューニング内容
- ループ交換の効果（ソースチューニング）

ループ交換とは

ループ交換とは、多重ループに対してループの順番を入れ替えることにより、データのアクセス効率を向上させる手法です。

Fortran言語の場合、配列は列方向に連續に確保されます。そのため、以下のようにループの順序を変更し連續アクセスにすることで、高速に動作します。

(C言語では、配列は行方向に連續に確保されるため、Fortranとは逆の順序になります)

改善前ソース例

```
do j=1,n1  
  do i=1,n2  
    a(j,i) = b(j,i) * c(j,i)  
  enddo  
enddo
```

配列a, b, cはストライドアクセス
キャッシュ効率が悪い

改善後ソース例

```
do i=1,n2  
  do j=1,n1  
    a(j,i) = b(j,i) * c(j,i)  
  enddo  
enddo
```

ループ順序を変更することで
連續アクセスになり
キャッシュ効率が改善



■ ポイント

- 最内ループの繰返し数が小さい場合は、ソフトウェアパイプラインングが行われなくなる場合があるため注意が必要です。
ただし、最内ループを SIMD 長に固定できる場合は、その外側のループでソフトウェアパイプラインングされます。
- ストアする配列とロードする配列でアクセス方向が異なる場合は、ストアする配列を連續アクセスにする方が性能が向上します。

■ 注意事項

- clangモードではループ交換機能を使用できません

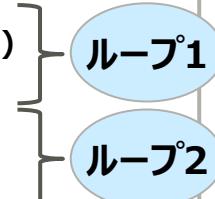
改善前ソース

```

do j=1,n1
  do i=1,n2
    a(i) = s1 + c(j,i) / (s1 + s2 / d(j,i))
  enddo
  do i=2,n2
    b(j,i) = a(i) / (s2 + s1 / a(i-1))
  enddo
enddo

```

配列b, c, dはストライド
アクセスなので
キャッシュ効率が悪い



① 配列aを2次元配列にする

```

do j=1,n1
  do i=1,n2
    a(j,i) = s1 + c(j,i) / (s1 + s2 / d(j,i))
  enddo
  do i=2,n2
    b(j,i) = a(j,i) / (s2 + s1 / a(j,i-1))
  enddo
enddo

```

ループ分割の阻害
要因である配列aの
依存がなくなる

② ループ¹とループ²を分割する

```

do j=1,n1
  do i=1,n2
    a(j,i) = s1 + c(j,i) / (s1 + s2 / d(j,i))
  enddo
enddo

do j=1,n1
  do i=2,n2
    b(j,i) = a(j,i) / (s2 + s1 / a(j,i-1))
  enddo
enddo

```



③ ループ¹を交換する

```

do i=1,n2
  do j=1,n1
    a(j,i) = s1 + c(j,i) / (s1 + s2 / d(j,i))
  enddo
enddo

-----  

do j=2,n2
  do j=1,n1
    b(j,i) = a(j,i) / (s2 + s1 / a(j,i-1))
  enddo
enddo

```

配列b, c, dが連続アクセスになり
キャッシュ効率が改善される

配列b、配列c、配列dはストライドアクセスのためキャッシュの利用効率が悪く、浮動小数点ロードアクセス待ちが多くなっています。

```

    改善前ソース
45      real*8 a(n2),b(n1,n2),c(n1,n2),d(n1,n2)
46      real*8 s1,s2
47      integer n1,n2
48      !$omp parallel
49      !$omp do private(a)
50 1 p      do j=1,n1
51      <<< Loop-information Start
52      <<< [OPTIMIZATION]
53      <<< SIMD(VL: 8)
54      <<< SOFTWARE PIPELINING(IPC: 1.96, ITR: 112, MVE: 3, POL: S)
55      <<< Loop-information End >>>
56 2 p 2v      do i=1,n2
57 2 p 2v          a(i) = s1 + c(j,i) / (s1 + s2 / d(j,i))
58 2 p 2v      enddo
59      <<< Loop-information Start >>>
60      <<< [OPTIMIZATION]
61      <<< SIMD(VL: 8)
62      <<< SOFTWARE PIPELINING(IPC: 2.27, ITR: 96, MVE: 2, POL: S)
63      <<< Loop-information End >>>
64 2 p 2v      do i=2,n2
65 2 p 2v          b(j,i) = a(i) / (s2 + s1 / a(i-1))
66 2 p 2v      enddo
67 1 p      enddo
68      !$omp end do
69      !$omp end parallel

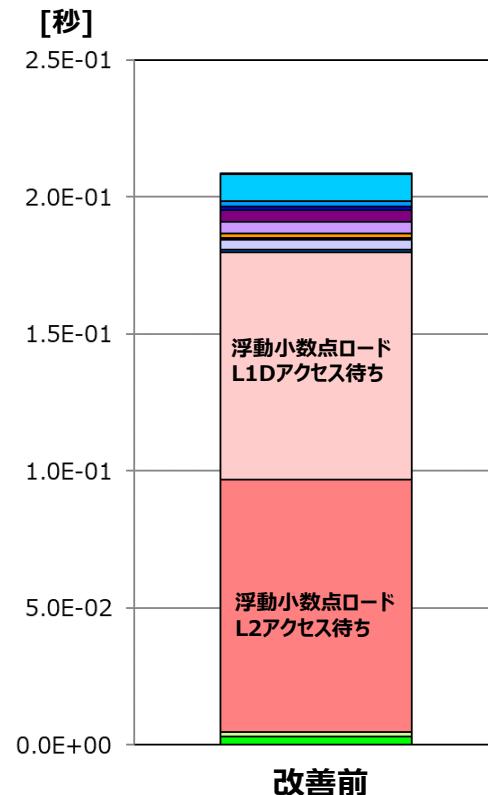
```

配列b、c、dは
ストライドアクセスなので
キャッシュ効率が悪い

ループ1

ループ2

L1Dミスが大きい



Cache	L1I miss rate (/Effective instruction)	Load-store instruction	L1D miss	L1D miss rate (/Load-store instruction)	L1D miss demand rate (%) (/L1D miss)	L1D miss hardware prefetch rate (%) (/L1D miss)	L1D miss software prefetch rate (%) (/L1D miss)	L2 miss	L2 miss rate (/Load-store instruction)	L2 miss demand rate (%) (/L2 miss)	L2 miss hardware prefetch rate (%) (/L2 miss)	L2 miss software prefetch rate (%) (/L2 miss)
改善前	0.00	3.60E+08	3.89E+08	1.08	98.28%	1.72%	0.00%	1.07E+04	0.00	62.54%	46.87%	0.00%

ループ交換により、配列を連続アクセスすることでキャッシュ効率が向上し、浮動小数点ロードアクセス待ちが改善されました。

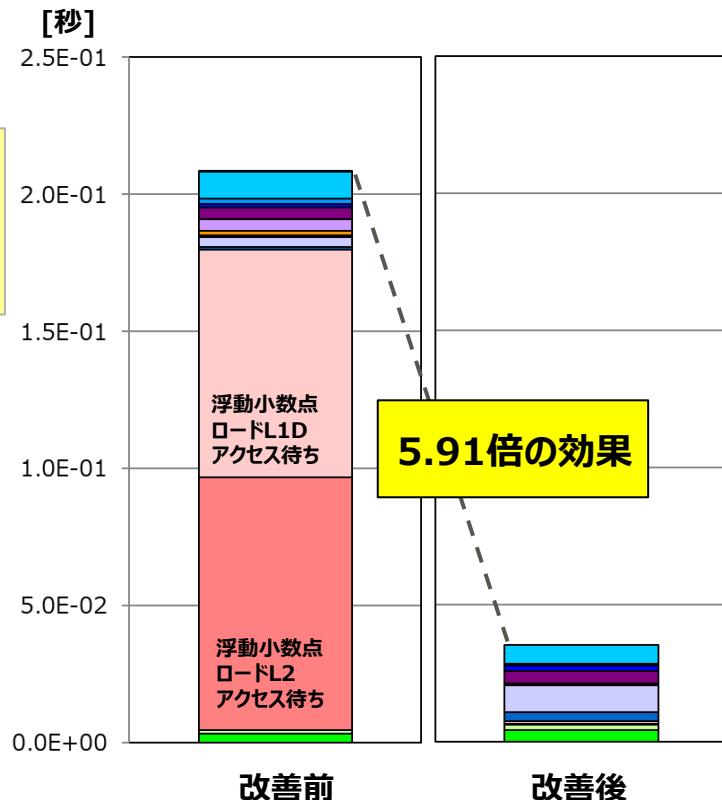
```

    改善後ソース
45      real*8 a(n1,n2),b(n1,n2),c(n1,n2),d(n1,n2)
46      real*8 s1,s2
47      integer n1,n2
48      !$omp parallel
49      !$omp do
50      1 p      do i=1,n2
51      <<< Loop-information Start >>>
52      <<< [OPTIMIZATION]
53      <<< SIMD(VL: 8)
54      <<< SOFTWARE PIPELINING(IPC: 2.13, ITR: 112, MVE: 2, POL: S)
55      <<< Loop-information End >>>
56      2 p 2v      do j=1,n1
57      2 p 2v          a(j,i) = s1 + c(j,i) / (s1 + s2 / d(j,i))
58      2 p 2v          enddo
59      1 p      enddo
60      !$omp end do
61      !$omp do
62      1 p      do i=2,n2
63      <<< Loop-information Start >>>
64      <<< [OPTIMIZATION]
65      <<< SIMD(VL: 8)
66      <<< SOFTWARE PIPELINING(IPC: 2.04, ITR: 96, MVE: 2, POL: S)
67      <<< Loop-information End >>>
68      2 p 2v      do j=1,n1
69      2 p 2v          b(j,i) = a(j,i) / (s2 + s1 / a(j,i-1))
70      2 p 2v          enddo
71      1 p      enddo
72      !$omp end do
73      !$omp end parallel

```

チューニング内容

- ①配列aを2次元配列にする
- ②ループ1とループ2を分割する
- ③ループ交換する



L1Dミス数が大幅に減少

Cache	L1I miss rate (/Effective instruction)	Load-store instruction	L1D miss	L1D miss rate (/Load-store instruction)	L1D miss demand rate (%) (/L1D miss)	L1D miss hardware prefetch rate (%) (/L1D miss)	L1D miss software prefetch rate (%) (/L1D miss)	L2 miss	L2 miss rate (/Load-store instruction)	L2 miss demand rate (%) (/L2 miss)	L2 miss hardware prefetch rate (%) (/L2 miss)	L2 miss software prefetch rate (%) (/L2 miss)
改善前	0.00	3.60E+08	3.89E+08	1.08	98.28%	1.72%	0.00%	1.07E+04	0.00	62.54%	46.87%	0.00%
改善後	0.00	1.94E+08	2.15E+07	0.11	9.28%	90.72%	0.00%	8.66E+03	0.00	17.98%	87.84%	0.00%

配列b、配列c、配列dはストライドアクセスのためキャッシュの利用効率が悪く、浮動小数点ロードアクセス待ちが多くなっています。

改善前ソース

```

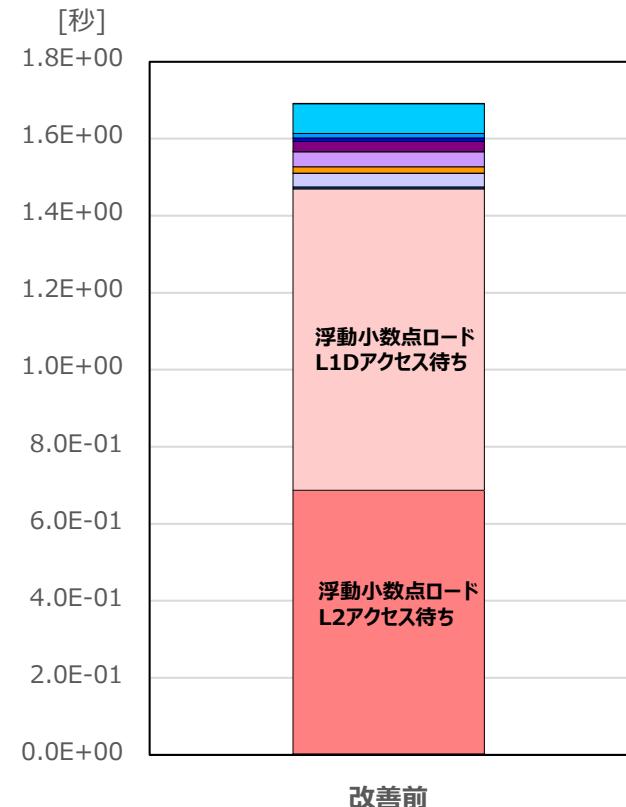
33 void sub(int n,int m,double s1,double s2) {
34     double a[n2];
35     int i,j;
36
37     <<< Loop-information Start >>>
38     <<< [OPTIMIZATION]
39     <<< PREFETCH(HARD) Expected by compiler :
40     <<< a
41     <<< Loop-information End >>>
42     for(j=0;j<n1;j++) {
43         <<< Loop-information Start >>>
44         <<< [OPTIMIZATION]
45         <<< SIMD(VL: 8)
46         <<< SOFTWARE PIPELINING(IPC: 1.92, ITR: 112, MVE: 3, POL: S)
47         <<< PREFETCH(HARD) Expected by compiler :
48         <<< a
49         <<< Loop-information End >>>
50         for(i=0;i<n2;i++) {
51             2v         a[i] = s1 + c[i][j] / (s1 + s2 / d[i][j]);
52         }
53         <<< Loop-information Start >>>
54         <<< [OPTIMIZATION]
55         <<< SIMD(VL: 8)
56         <<< SOFTWARE PIPELINING(IPC: 2.18, ITR: 96, MVE: 2, POL: S)
57         <<< PREFETCH(HARD) Expected by compiler :
58         <<< a
59         <<< Loop-information End >>>
60         for(i=1;i<n2;i++) {
61             2v         b[i][j] = a[i] / (s2 + s1 / a[i-1]);
62         }
63     }
64 }
```

配列b、c、dはストライドアクセスなのでキャッシュ効率が悪い

ループ1

ループ2

L1Dミスが大きい



Cache	L1I miss rate (/Effective instruction)	Load-store instruction	L1D miss	L1D miss rate (/Load-store instruction)	L1D miss demand rate (%) (/L1D miss)	L1D miss hardware prefetch rate (%) (/L1D miss)	L1D miss software prefetch rate (%) (/L1D miss)	L2 miss	L2 miss rate (/Load-store instruction)	L2 miss demand rate (%) (/L2 miss)	L2 miss hardware prefetch rate (%) (/L2 miss)	L2 miss software prefetch rate (%) (/L2 miss)
改善前	0.00	1.87E+09	3.89E+08	0.21	97.94%	2.06%	0.00%	2.18E+04	0.00	71.23%	38.77%	0.00%

ループ交換により、配列を連続アクセスすることでキャッシュ効率が向上し、浮動小数点ロードアクセス待ちが改善されました。

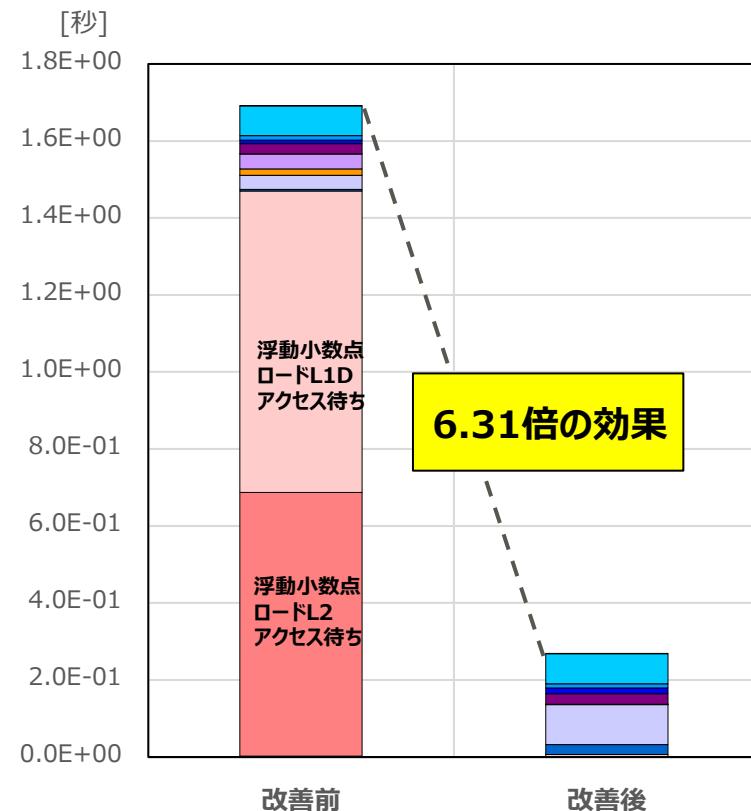
改善後ソース

```

34     void sub(int n1,int n2,double s1,double s2)
35     {
36         double a[m][n];
37         int i,j;
38
39         <<< Loop-information Start >>>
40         <<< [OPTIMIZATION]
41         <<< PREFETCH(HARD) Expected by compiler :
42             d, c, a
43         <<< Loop-information End >>>
44             for(i=0;i<n2;i++) {
45                 <<< Loop-information Start >>>
46                 <<< [OPTIMIZATION]
47                 <<< SIMD(VL: 8)
48                 <<< SOFTWARE PIPELINING(IPC: 1.88, ITR: 96, MVE: 2, POL: S)
49                 <<< PREFETCH(HARD) Expected by compiler :
50                     c, d, a
51                 <<< Loop-information End >>>
52
53                 2v             for(j=0;j<n1;j++) {
54                     2v             a[i][j] = s1 + c[i][j] / (s1 + s2 / d[i][j]);
55                 2v             }
56
57                 <<< Loop-information Start >>>
58                 <<< [OPTIMIZATION]
59                 <<< PREFETCH(HARD) Expected by compiler :
60                     a, b
61                 <<< Loop-information End >>>
62                     for(i=0;i<n2;i++) {
63                         <<< Loop-information Start >>>
64                         <<< [OPTIMIZATION]
65                         <<< SIMD(VL: 8)
66                         <<< SOFTWARE PIPELINING(IPC: 2.09, ITR: 96, MVE: 2, POL: S)
67                         <<< PREFETCH(HARD) Expected by compiler :
68                             a, b
69                         <<< Loop-information End >>>
70                         2v             for(j=0;j<n1;j++) {
71                             2v             b[i][j] = a[i][j] / (s2 + s1 / a[i-1][j]);
72                         2v             }
73
74                 }
    
```

チューニング内容

- ①配列aを2次元配列にする
- ②ループ1とループ2を分割する
- ③ループ交換する



L1Dミス数が大幅に減少

Cache	L1I miss rate (/Effective instruction)	Load-store instruction	L1D miss	L1D miss rate (/Load-store instruction)	L1D miss demand rate (%) (/L1D miss)	L1D miss hardware prefetch rate (%) (/L1D miss)	L1D miss software prefetch rate (%) (/L1D miss)	L2 miss	L2 miss rate (/Load-store instruction)	L2 miss demand rate (%) (/L2 miss)	L2 miss hardware prefetch rate (%) (/L2 miss)	L2 miss software prefetch rate (%) (/L2 miss)
改善前	0.00	1.87E+09	3.89E+08	0.21	97.94%	2.06%	0.00%	2.18E+04	0.00	71.23%	38.77%	0.00%
改善後	0.00	1.87E+09	2.27E+07	0.01	14.23%	85.77%	0.00%	3.63E+04	0.00	42.14%	62.84%	0.00%

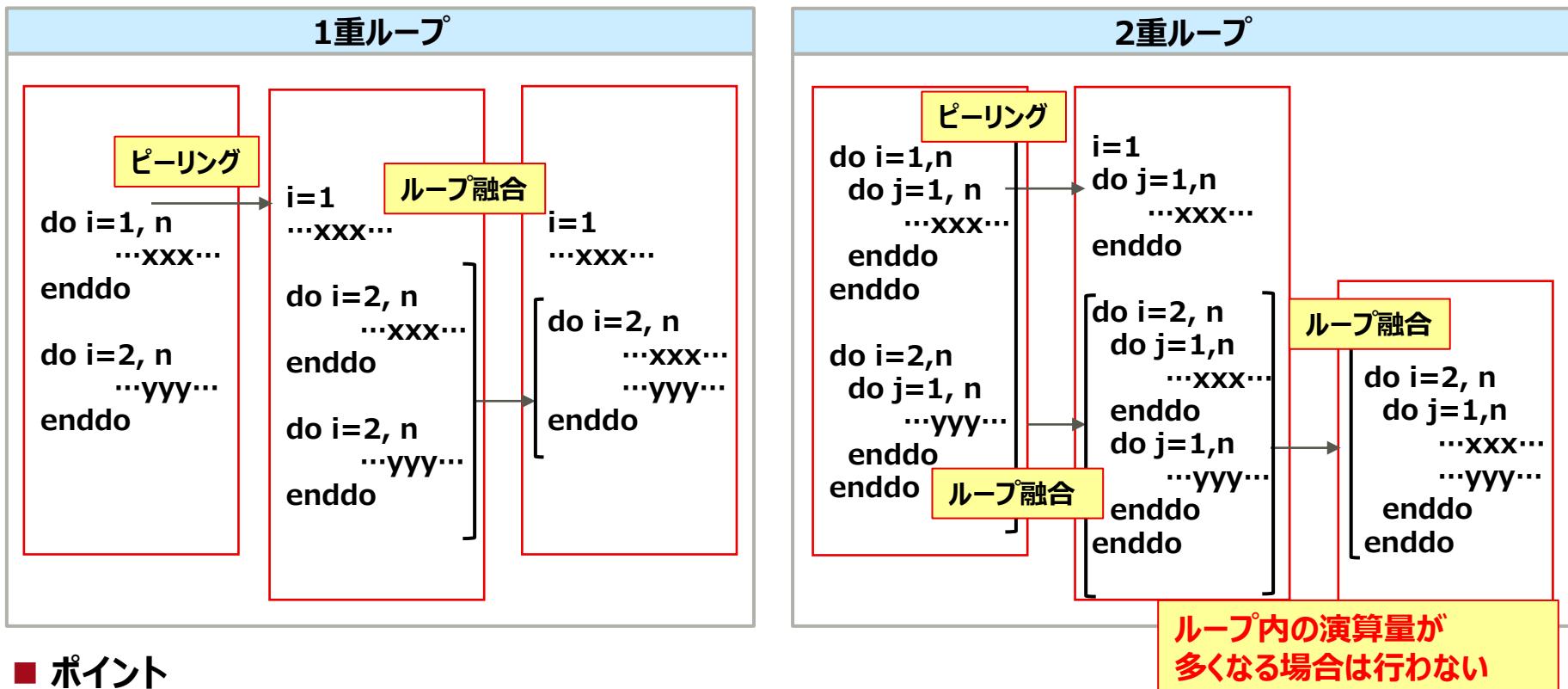
ループ融合

- ループ融合とは
- ループ融合（改善前）
- ループ融合の効果（ソースチューニング）

ループ融合とは

ループ融合とは、以下の効果を目的にループを連結させることです。

- ・ データの局所化：配列の再利用を行う
- ・ 命令レベルの並列化向上：ループ内の命令数を増やし、命令レベルの並列性を向上させる



■ ポイント

- ループ長が同じループの結合はコンパイラによって自動で行われます。
- ループ内の演算が多くなりすぎると、ソフトウェアパイプラインングが促進されなくなる場合があります。

ループ¹の繰返し数が多く、配列データがキャッシュに載りきらないため、ループ²で再利用できません。そのため、メモリ・キャッシュビジー待ちが多くなっています。

改善前ソース

```

42 1 pp v      do j=1,m-1
     <<< Loop-information Start >>>
     <<< [OPTIMIZATION]
     <<< COLLAPSED
     <<< Loop-information End >>>
43 2 p          do i=1,n
44 2 p v        s1 = s1 + a(i,j) * (s3 * b(i,j) + c(i,j) * (s2 + s3 * d(i,j) ))
45 2 p v        enddo
46 1 p v        enddo
47
     <<< Loop-information Start >>>
     <<< [PARALLELIZATION]
     <<< Standard iteration count: 572
     <<< [OPTIMIZATION]
     <<< COLLAPSED
     <<< SIMD(VL: 8)
     <<< SOFTWARE PIPELINING(IPC: 2.30, ITR: 96, MVE: 2, POL: S)
     <<< PREFETCH(HARD) Expected by compiler :
     <<< b, a, d, c, e
     <<< Loop-information End >>>
48 1 pp 2v      do j=1,m
     <<< Loop-information Start >>>
     <<< [OPTIMIZATION]
     <<< COLLAPSED
     <<< Loop-information End >>>
49 2 p 2        do i=1,n
50 2 p 2v      e(i,j) = s2 * (a(i,j) + b(i,j) * (s3 + c(i,j) * d(i,j)))
51 2 p 2v      enddo
52 1 p         enddo

```

m = 50
n = 150000
配列の型 : real*8

配列データの合計：約200MB
配列データがキャッシュに載りきらない

配列アクセスがキャッシュミス



改善前

L1キャッシュミス率とL2キャッシュミス率はストリームアクセスの理論値である0.24になっています。ただし、ループ¹,²の両方でミスが発生しています。つまりループ¹でキャッシュに載せたデータをループ²で使用できません。

Cache	L1 miss rate (/Effective instruction)	Load-store instruction	L1D miss	L1D miss rate (/Load-store instruction)	L1D miss demand rate (%) (/L1D miss)	L1D miss hardware prefetch rate (%) (/L1D miss)	L1D miss software prefetch rate (%) (/L1D miss)	L2 miss	L2 miss rate (/Load-store instruction)	L2 miss demand rate (%) (/L2 miss)	L2 miss hardware prefetch rate (%) (/L2 miss)	L2 miss software prefetch rate (%) (/L2 miss)
改善前	0.00	8.57E+08	2.09E+08	0.24	0.64%	99.35%	0.01%	2.09E+08	0.24	0.66%	99.50%	0.00%

ループ融合によってキャッシュ効率が向上し、メモリ・キャッシュビジー待ちが改善されました。

改善後ソース（ソースチューニング）

```

42 1 pp v      do j=1,m-1
    <<< Loop-information Start >>>
    <<< [OPTIMIZATION]
    <<< COLLAPSED
    <<< Loop-information End >>>
43 2 p          do i=1,n
44 2 p v        s1 = s1 + a(i,j) * (s3 * b(i,j) + c(i,j) * (s2 + s3 * d(i,
45 2 p v        e(i,j) = s2 * (a(i,j) + b(i,j) * (s3 + c(i,j) * d(i,j)))
46 2 p v        enddo
47 1 p v        enddo
:
49 1 s          do j=m,m
    <<< Loop-information Start >>>
    <<< [PARALLELIZATION]
    <<< Standard iteration count: 364
    <<< [OPTIMIZATION]
    <<< SIMD(VL: 8)
    <<< SOFTWARE PIPELINING
    <<< PREFETCH(HARD) Exp. component
    <<< b, a, d, c, e
    <<< Loop-information End >>>
50 2 pp 2v      do i=1,n
51 2 p 2v      e(i,j) = s2 * (a(i,j) + b(i,j) * (s3 + c(i,j) * d(i,j)))
52 2 p 2v      enddo
53 1 p          enddo

```

ループ融合

配列アクセスがキャッシュヒットする

ピーリング

ループ融合させるために切り出す（ピーリング）

L1Dミス数とL2ミス数が大幅に減少

ループ融合のイメージ

ループ融合のイメージ

```

do j=1, m-1
  do i=1, n
    ...xxx...
  enddo
enddo

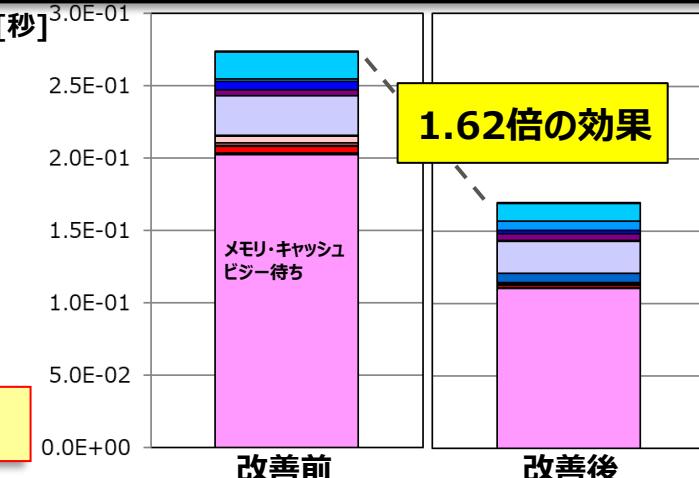
do j=1, m
  do i=1, n
    ...yyy...
  enddo
enddo

do j= m,m
  do i=1, n
    ...yyy...
  enddo
enddo

```

ループ融合

ピーリング



Cache	L1I miss rate (/Effective instruction)	Load-store instruction	L1D miss	L1D miss rate (/Load-store instruction)	L1D miss demand rate (%) / (L1D miss)	L1D miss hardware prefetch rate (%) / (L1D miss)	L1D miss software prefetch rate (%) / (L1D miss)	L2 miss	L2 miss rate (/Load-store instruction)	L2 miss demand rate (%) / (L2 miss)	L2 miss hardware prefetch rate (%) / (L2 miss)	L2 miss software prefetch rate (%) / (L2 miss)
改善前	0.00	8.57E+08	2.09E+08	0.24	0.64%	99.35%	0.01%	2.09E+08	0.24	0.66%	99.50%	0.00%
改善後	0.00	4.92E+08	1.18E+08	0.24	0.36%	99.63%	0.00%	1.17E+08	0.24	0.41%	99.75%	0.00%

ループ¹の繰返し数が多く、配列データがキャッシュに載りきらないため、ループ²で再利用できません。そのため、メモリ・キャッシュビギー待ちが多くなっています。

改善前ソース

```

39      #pragma omp parallel for
40    p   for (j=0;j<m-1;j++) {
41      <<< Loop-information Start >>>
42      <<< [OPTIMIZATION]
43      <<< SIMD(VL: 8)
44      <<< PREFETCH(HARD) Expected by compiler :
45      <<< (unknown)
46      <<< Loop-information End >>>
47    p   8v   for (i=0;i<n;i++) {
48    p     *s1 = *s1 + a[j][i] * (*s3 * b[j][i] + c[j][i] * (*s2 + *s3 * d[j][i]));
49    p   8v   }
50    p   }
51    p   }

#pragma omp parallel for
<<< Loop-information Start >>>
<<< [OPTIMIZATION]
<<< SIMD(VL: 8)
<<< SOFTWARE PIPELINING(IPC: 2.30, ITR: 96, MVE: 2, POL: S)
<<< PREFETCH(HARD) Expected by compiler :
<<< (unknown)
<<< Loop-information End >>>
52    p   for (j=0;j<m;j++) {
53      <<< Loop-information Start >>>
54      <<< [OPTIMIZATION]
55      <<< SIMD(VL: 8)
56      <<< PREFETCH(HARD) Expected by compiler :
57      <<< (unknown)
58      <<< Loop-information End >>>
59    p   2v   for (i=0;i<n;i++) {
60    p     e[j][i] = *s2 * (a[j][i] + b[j][i]* (*s3 + c[j][i] * d[j][i]));
61    p   2v   }
62    p   }
63    p   }

```

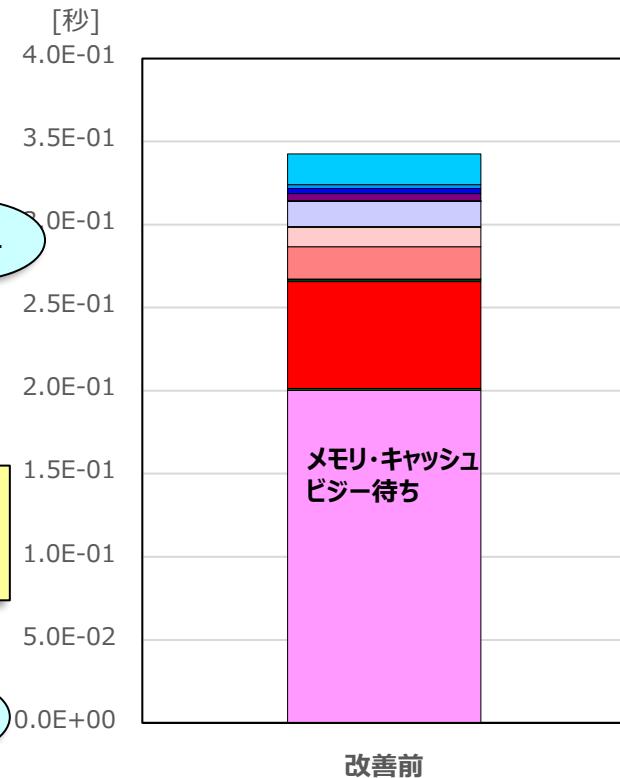
m = 50
n = 150000
配列の型 : double

配列データの合計：約200MB
配列データがキャッシュに載りきらない

配列アクセスがキャッシュミス

ループ¹

ループ²



L1キャッシュミス率とL2キャッシュミス率はストリームアクセスの理論値に近い0.21になっています。ただし、ループ¹,²の両方でミスが発生しています。つまりループ¹でキャッシュに載せたデータをループ²で使用できません。

Cache	L1 miss rate (/Effective instruction)	Load-store instruction	L1D miss	L1D miss rate (/Load-store instruction)	L1D miss demand rate (%) (/L1D miss)	L1D miss hardware prefetch rate (%) (/L1D miss)	software prefetch rate (%) (/L1D miss)	L2 miss	L2 miss rate (/Load-store instruction)	L2 miss demand rate (%) (/L2 miss)	hardware prefetch rate (%) (/L2 miss)	software prefetch rate (%) (/L2 miss)
改善前	0.00	9.87E+08	2.10E+08	0.21	9.82%	90.20%	-0.02%	2.10E+08	0.21	5.83%	95.03%	0.00%

ループ融合によってキャッシュ効率が向上し、メモリ・キャッシュビジー待ちが改善されました。

改善後ソース (ソースチューニング)

```

40 p   for (j=0;j<m-1;j++) {
        <<< Loop-information Start >>>
        <<< [OPTIMIZATION]
        SIMD(VL: 8)
        <<< SOFTWARE PIPELINING(IPC: 1.98, ITR: 192, MVE: 2, POL: S)
        <<< PREFETCH(HARD) Expected by compiler :
        <<< ...
        <<< Loop-information End >>>
41 p   8v   for (i=0;i<n;i++) {
42 p   8v     *s1 = *s1 + a[j][i] * (*s3 * b[j][i] + c[j][i] * (*s2 + *s3 * d[j][i] ));
43 p   8v     e[j][i] = *s2 * (a[j][i] + b[j][i]* (*s3 + c[j][i] * d[j][i] ));
44 p   8v   }
45 p   }
...
48 p   for (j=m-1;j<m;j++) {
        <<< Loop-information Start >>>
        <<< [OPTIMIZATION]
        SIMD(VL: 8)
        <<< SOFTWARE PIPELINING(IPC: 3.83, ITR: 176, MVE: 3, POL: S)
        <<< PREFETCH(HARD) Expected by compiler :
        <<< ...
        <<< Loop-information End >>>
49 p   2v   for (i=0;i<n;i++) {
50 p   2v     e[j][i] = *s2 * (a[j][i] + b[j][i]* (*s3 + c[j][i] * d[j][i] ));
51 p   2v   }
52 p   }

```

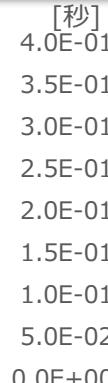
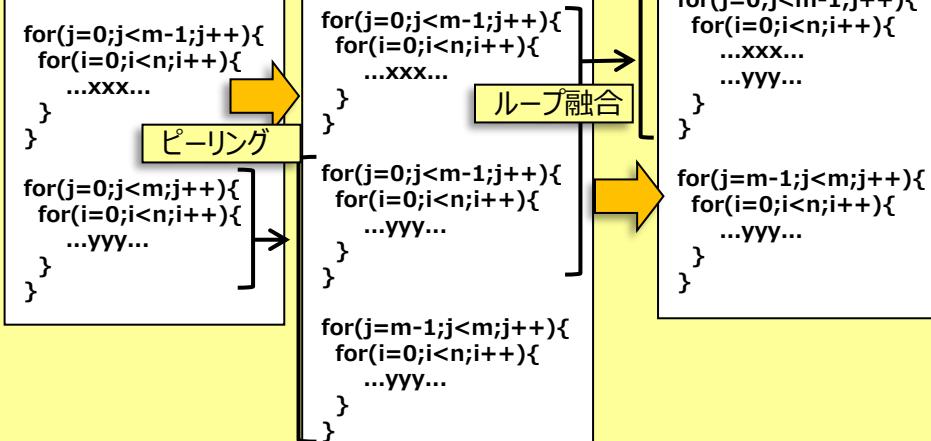
ループ融合

配列アクセスがキャッシュヒットする

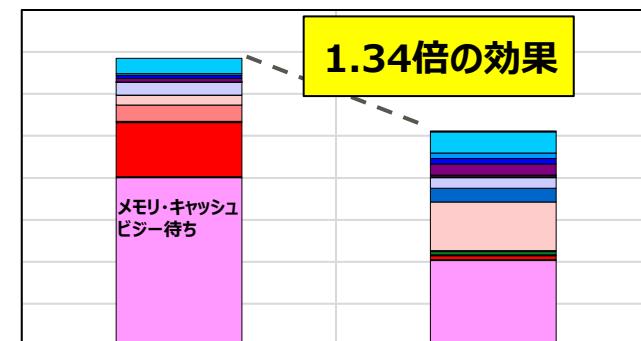
ピーリング

ループ融合させるために切り出す (ピーリング)

ループ融合のイメージ



1.34倍の効果



L1Dミス数とL2ミス数が大幅に減少

改善前

改善後

Cache	L1I miss rate (/Effective instruction)	Load-store instruction	L1D miss	L1D miss rate (/Load-store instruction)	L1D miss demand rate (%) (/L1D miss)	L1D miss hardware prefetch rate (%) (/L1D miss)	L1D miss software prefetch rate (%) (/L1D miss)	L2 miss	L2 miss rate (/Load-store instruction)	L2 miss demand rate (%) (/L2 miss)	L2 miss hardware prefetch rate (%) (/L2 miss)	L2 miss software prefetch rate (%) (/L2 miss)
改善前	0.00	9.87E+08	2.10E+08	0.21	9.82%	90.20%	-0.02%	2.10E+08	0.21	5.83%	95.03%	0.00%
改善後	0.00	1.94E+09	1.19E+08	0.06	2.55%	97.54%	-0.09%	1.18E+08	0.06	1.46%	98.73%	0.00%

配列マージ（インダイレクトアクセスの効率改善）

- 配列マージとは
- 配列マージ（改善前）
- 配列マージの効果（ソースチューニング）

配列マージとは、同一ループ内でアクセスパターンが共通の配列が複数ある場合に、それらを1つの配列に融合することです。データアクセスを連續化して、キャッシュ効率を向上させます。

改善前ソース	改善後ソース
<pre> parameter(n=1000000) real*8 a(n), b(n), c(n) integer d(n+10) : do iter = 1, 100 do i = 1 , n a(d(i)) = b(d(i)) + scalar * c(d(i)) enddo enddo : 異なるキャッシュラインのアクセス (配列dの値が連続値でない場合) </pre> <p style="text-align: center;">(L1Dキャッシュ) </p>	<pre> parameter(n=1000000) real*8 abc(3, n) integer d(n+10) : do iter = 1, 100 do i = 1 , n abc(1, d(i)) = abc(2, d(i)) + scalar * abc(3, d(i)) enddo enddo : 同じキャッシュラインのアクセス </pre> <p style="text-align: center;">(L1Dキャッシュ) </p>

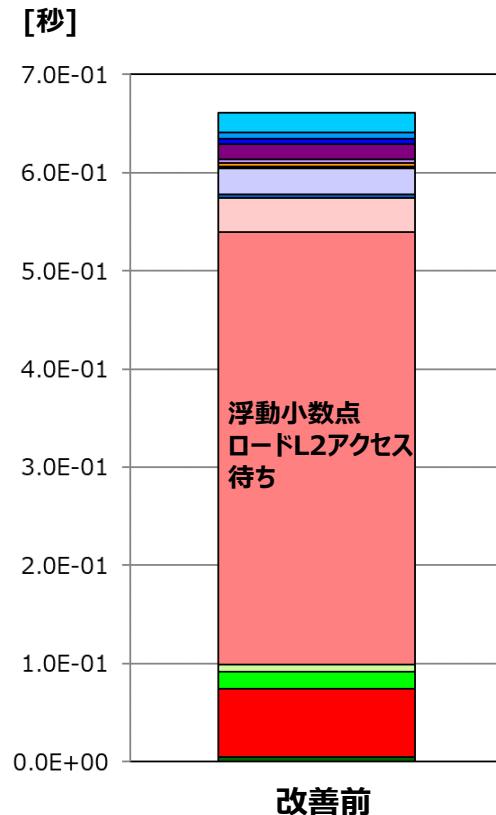
L1Dミス率が高いことから、キャッシュの利用効率が悪く（インダイレクトアクセス）、浮動小数点ロードL2アクセス待ちが多くなっています。

改善前ソース

```

1      parameter(n=2*1000*1000/8)
2      real*8 a(n),b(n),c(n),e(n),f(n),s
3      integer d(n)
4
5      1 s  s      call sub(a,b,c,d,e,f,s,n)
6
7
8      subroutine sub(a,b,c,d,e,f, s, n)
9      real*8 a(n),b(n),c(n),e(n),f(n),s
10     integer d(n), ii
11
12
13
14      !$omp parallel do schedule (static,96)
15      <<< Loop-information Start >>>
16      <<< [OPTIMIZATION]
17      <<< SIMD(VL: 8)
18      <<< SOFTWARE PIPELINING(IPC: 1.07, ITR: 80, MVE: 3, POL: S)
19      <<< PREFETCH(HARD) Expected by compiler :
20      <<< d
21      <<< Loop-information End >>>
22
23      1 p 2v      do i = 1 , n
24      1 p 2v      ii = d(i)
25      1 p 2v      a(ii) = s / ( s + f(ii) / ( s + e(ii) / ( b(ii) + s / c(ii)))) )
26      1 p 2v      enddo
27
28      !$omp end parallel do
29
30
31
32
33
34

```



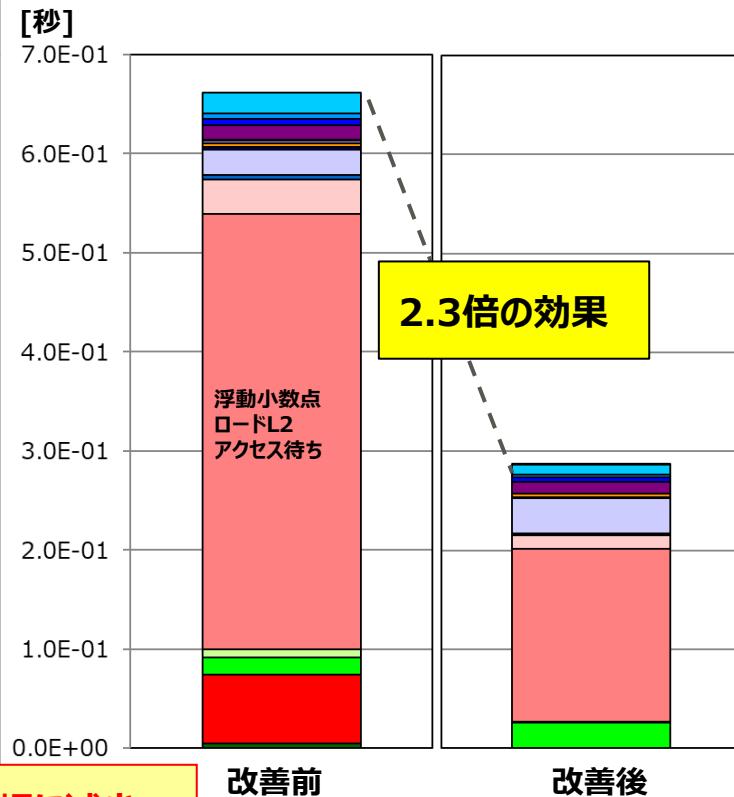
Cache	L1I miss rate (/Effective instruction)	Load-store instruction	L1D miss	L1D miss rate (/Load-store instruction)	L1D miss demand rate (%) (/L1D miss)	L1D miss hardware prefetch rate (%) (/L1D miss)	L1D miss software prefetch rate (%) (/L1D miss)	L2 miss	L2 miss rate (/Load-store instruction)	L2 miss demand rate (%) (/L2 miss)	L2 miss hardware prefetch rate (%) (/L2 miss)	L2 miss software prefetch rate (%) (/L2 miss)
改善前	0.00	6.44E+08	1.27E+09	1.97	99.98%	0.01%	0.00%	4.82E+07	0.07	60.81%	50.33%	0.00%

リストアクセスの配列を配列マージすることで、キャッシュ効率が向上し、浮動小数点ロードL2アクセス待ちが改善されました。

改善後ソース（ソースチューニング）

```

1      parameter(n=2*1000*1000/8)
2      real*8 abcef(5,n),s
3      integer d(n)
4
5      1 s s      call sub(abcef,d,s,n)
6
7      :
8
9      subroutine sub(abcef,d, s, n)
10     real*8 abcef(5,n),s
11     integer d(n), ii
12
13
14
15
16
17
18
19
20
21
22
23
24
25
26
27
28
29
30
31
32
33
34
      !$omp parallel do schedule (static,96)
      <<< Loop-information Start >>>
      <<< [OPTIMIZATION]
      <<< SIMD(VL: 8)
      <<< SOFTWARE PIPELINING(IPC: 1.09, ITR: 80, MVE: 3, POL: S)
      <<< PREFETCH(HARD) Expected by compiler :
      <<< d
      <<< Loop-information End >>>
      1 p 2v      do i = 1 , n
      1 p 2v          ii = d(i)
      1 p 2v          abcef(1,ii) = s / (s + abcef(5,ii) / (s + abcef(4,ii)
      1 p 2v          * / (abcef(2,ii) + s / abcef(3,ii))))
      1 p 2v      enddo
      !$omp end parallel do
  
```



L1Dミス率が大幅に減少

2.3倍の効果

Cache	L1I miss rate (/Effective instruction)	Load-store instruction	L1D miss	L1D miss rate (/Load-store instruction)	L1D miss demand rate (%) (/L1D miss)	L1D miss hardware prefetch rate (%) (/L1D miss)	L1D miss software prefetch rate (%) (/L1D miss)	L2 miss	L2 miss rate (/Load-store instruction)	L2 miss demand rate (%) (/L2 miss)	L2 miss hardware prefetch rate (%) (/L2 miss)	L2 miss software prefetch rate (%) (/L2 miss)
改善前	0.00	6.44E+03	1.27E+09	1.97	99.98%	0.01%	0.00%	4.82E+07	0.07	60.81%	50.33%	0.00%
改善後	0.00	3.19E+03	2.97E+08	0.93	99.68%	0.32%	0.00%	1.58E+04	0.00	63.74%	54.78%	0.00%

L1Dミス率が高いことから、キャッシュの利用効率が悪く（インダイレクトアクセス）、浮動小数点ロードL2アクセス待ちが多くなっています。

改善前ソース

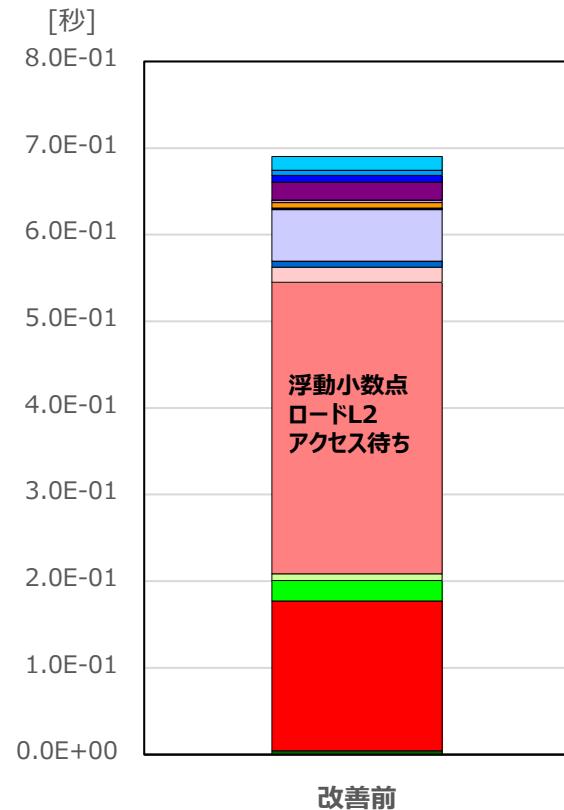
```

24 void sub(double (* restrict a),double (* restrict b),
           double (* restrict c),int (* restrict d),double (* restrict e),
           double (* restrict f),double s,int n) {
25   int i,ii;
26
27   #pragma omp parallel for schedule (static,96)
28   <<< Loop-information Start >>>
29   <<< [OPTIMIZATION]
30   <<< SIMD(VL: 8)
31   <<< SOFTWARE PIPELINING(IPC: 1.33, ITR: 112, MVE: 4, POL: S)
32   <<< PREFETCH(HARD) Expected by compiler :
33   <<< (unknown)
34   <<< Loop-information End >>>
35   p 2v for(i=0;i<n;i++) {
36   p 2v   ii = d[i];
37   p 2v   a[ii] = s / (s + f[ii] / (s + e[ii] / (b[ii] + s / c[ii])));
38   p 2v }
39 }
```

配列宣言

```

double a[250000];
double b[250000];
double c[250000];
double e[250000];
double f[250000];
```



改善前

Cache	L1I miss rate (/Effective instruction)	Load-store instruction	L1D miss	L1D miss rate (/Load-store instruction)	L1D miss demand rate (%) (/L1D miss)	L1D miss hardware prefetch rate (%) (/L1D miss)	L1D miss software prefetch rate (%) (/L1D miss)	L2 miss	L2 miss rate (/Load-store instruction)	L2 miss demand rate (%) (/L2 miss)	L2 miss hardware prefetch rate (%) (/L2 miss)	L2 miss software prefetch rate (%) (/L2 miss)
改善前	0.00	6.58E+08	1.27E+09	1.93	99.94%	0.06%	0.00%	4.24E+07	0.06	56.24%	56.98%	0.00%

リストアクセスの配列を配列マージすることで、キャッシュ効率が向上し、浮動小数点ロードL2アクセス待ちが改善されました。

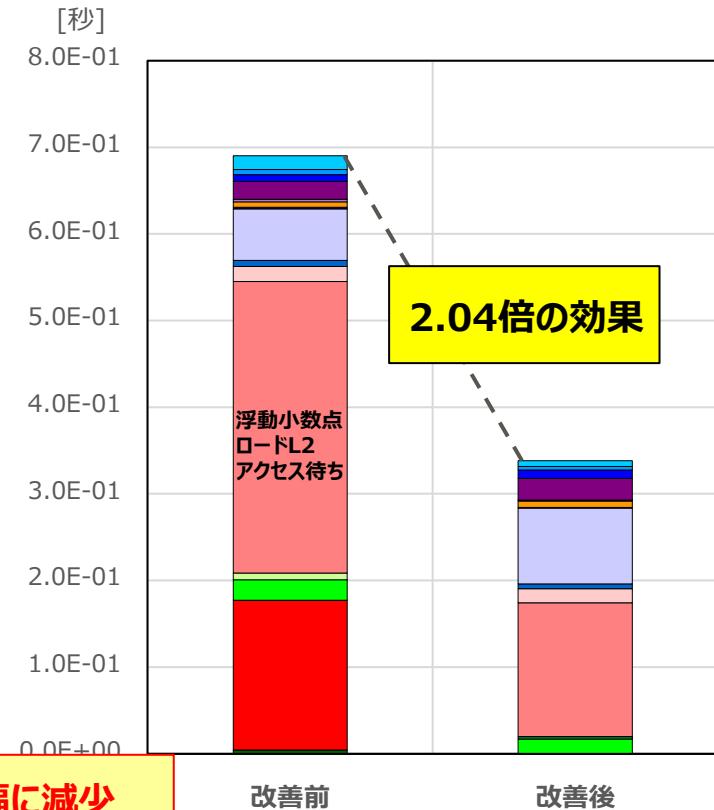
改善後ソース（ソースチューニング）

```

24 void sub(double (* restrict abcef)[5],int (* restrict d),double s,int n) {
25     int i,ii;
26
27 #pragma omp parallel for schedule (static,96)
28 <<< Loop-information Start >>>
29 <<< [OPTIMIZATION]
30 <<< SIMD(VL: 8)
31 <<< SOFTWARE PIPELINING(IPC: 1.36, ITR: 112, MVE: 4, POL: S)
32 <<< PREFETCH(HARD) Expected by compiler :
33 <<< (unknown)
34 <<< Loop-information End >>>
35 p 2v  for(i=0;i<n;i++) {
36 p 2v    ii = d[i];
37 p 2v    abcef[ii][0] = s / (s + abcef[ii][4] / ( s + abcef[ii][3] /
38 p 2v      ( abcef[ii][1] + s / abcef[ii][2] ) )
39 p 2v    }
40 }
```

配列宣言

```
double a[250000];
double b[250000];
double c[250000];
double e[250000];
double f[250000];
```



L1Dミス率が大幅に減少

改善前

改善後

Cache	L1 miss rate (/Effective instruction)	Load-store instruction	L1D miss	L1D miss rate (/Load-store instruction)	L1D miss demand rate (%) (/L1D miss)	L1D miss hardware prefetch rate (%) (/L1D miss)	L1D miss software prefetch rate (%) (/L1D miss)	L2 miss	L2 miss rate (/Load-store instruction)	L2 miss demand rate (%) (/L2 miss)	L2 miss hardware prefetch rate (%) (/L2 miss)	L2 miss software prefetch rate (%) (/L2 miss)
改善前	0.00	6.58E+08	1.27E+09	1.93	99.94%	0.06%	0.00%	4.24E+07	0.06	56.24%	56.98%	0.00%
改善後	0.00	3.67E+08	2.94E+08	0.80	99.69%	0.30%	0.00%	1.41E+04	0.00	76.45%	61.05%	0.00%

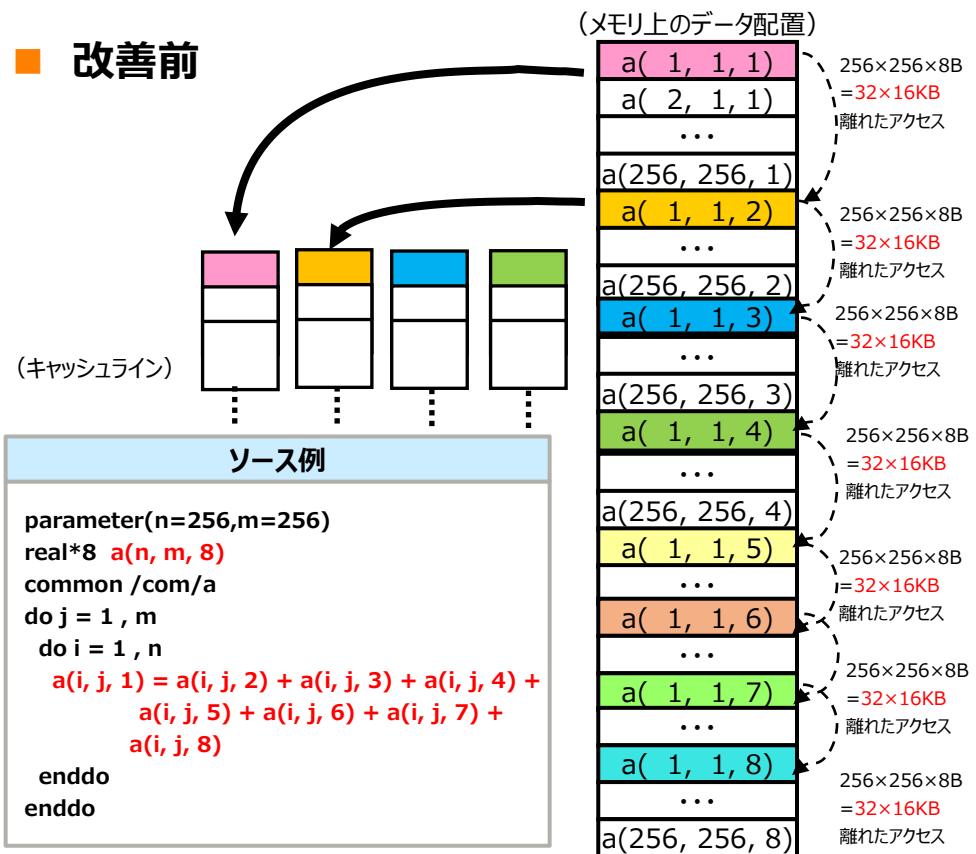
配列次元移動

- 配列次元移動とは
- 配列次元移動（改善前）
- 配列次元移動の効果（ソースチューニング）
- 配列次元移動の効果（翻訳オプションチューニング）

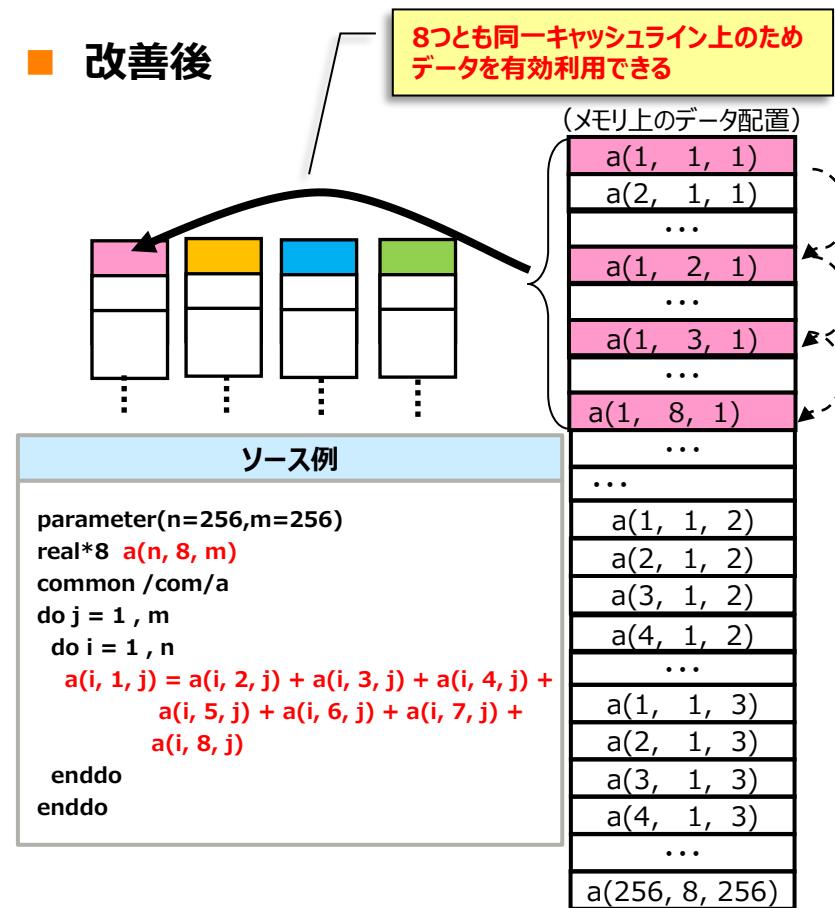
配列次元移動とは

配列次元移動とは、配列のアクセス次元を変更し、キャッシュの利用率を上げるチューニングです。次の例のように、変化する次元を内側に移動することにより、キャッシュの利用効率を向上させることができます。

改善前



改善後



→ キャッシュへの格納

→ メモリへのアクセス順番

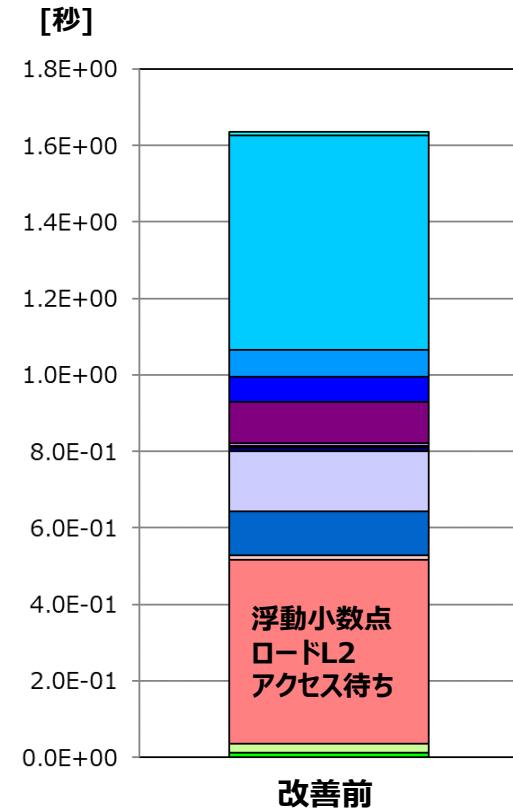
配列aの最内ループでのアクセスはデータ局所性が低いため、浮動小数点ロードL2アクセス待ちが発生しています。

改善前ソース

```

16      real(8)::a(N,M,8)
:
<<< Loop-information Start >>>
<<< [OPTIMIZATION]
<<< PREFETCH(HARD) Expected by compiler :
<<< a
<<< Loop-information End >>>
21  2 p           do j=1,M
<<< Loop-information Start >>>
<<< [OPTIMIZATION]
<<< SIMD(VL: 8)
<<< SOFTWARE PIPELINING(IPC: 2.87, ITR: 56,
                           MVE: 2, POL: S)
<<< PREFETCH(HARD) Expected by compiler :
<<< a
<<< Loop-information End >>>
22  3 p v         do i=1,N
23  3 p v         a(i,j,1)=a(i,j,2)+a(i,j,3)+a(i,j,4)+&
24  3                   a(i,j,5)+a(i,j,6)+a(i,j,7)+a(i,j,8)
25  3 p v         enddo
26  2 p           enddo

```



Cache

	L1 miss rate (/Effective instruction)	Load-store instruction	L1D miss	L1D miss rate (/Load-store instruction)	L1D miss demand rate (%) (/L1D miss)	L1D miss hardware prefetch rate (%) (/L1D miss)	L1D miss software prefetch rate (%) (/L1D miss)	L2 miss	L2 miss rate (/Load-store instruction)	L2 miss demand rate (%) (/L2 miss)	L2 miss hardware prefetch rate (%) (/L2 miss)	L2 miss software prefetch rate (%) (/L2 miss)
改善前	0.00	2.12E+10	8.20E+08	0.04	78.62%	21.38%	0.00%	5.45E+03	0.00	87.63%	20.43%	0.00%

配列次元移動を行うことでデータの局所性が向上され、浮動小数点L2アクセス待ちが改善されました。

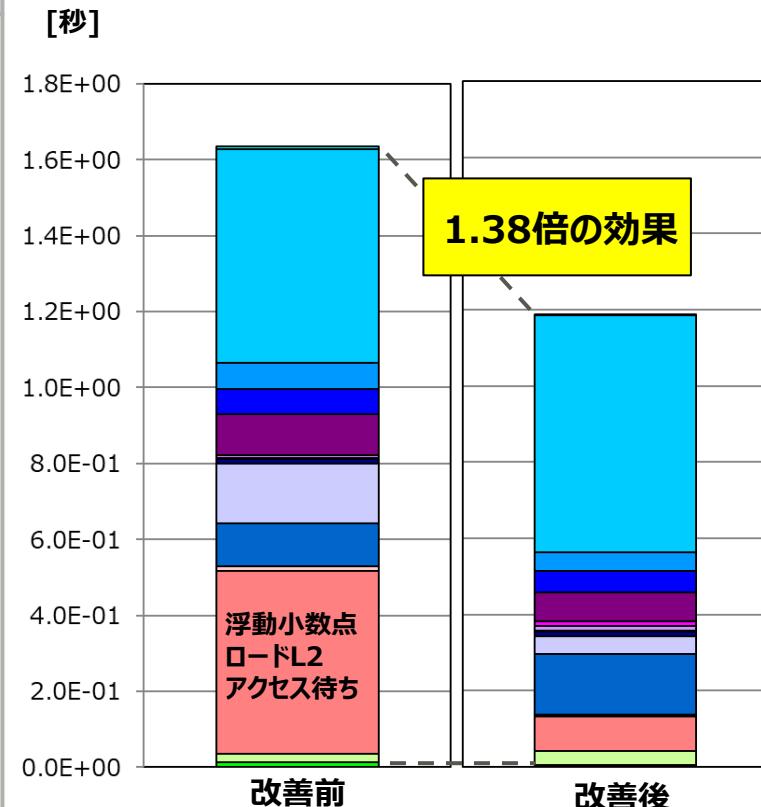
改善後ソース（ソースチューニング）

```

16      real(8)::a(N,8,M)
        <<< Loop-information Start >>
        <<< [OPTIMIZATION]
        <<< PREFETCH(HARD) Expect
        <<< a
        <<< Loop-information End >>>
21  2 p      do j=1,M
        <<< Loop-information Start >>>
        <<< [OPTIMIZATION]
        <<< SIMD(VL: 8)
        <<< SOFTWARE PIPELINING(IPC: 2.87, ITR: 56,
                                MVE: 2, POL: S)
        <<< PREFETCH(HARD) Expected by compiler :
        <<< a
        <<< Loop-information End >>>
22  3 p v      do i=1,N
23  3 p v          a(i,1,j)=a(i,2,j)+a(i,3,j)+a(i,4,j)+&
24  3                  a(i,5,j)+a(i,6,j)+a(i,7,j)+a(i,8,j)
25  3 p v      enddo
26  2 p      enddo

```

N=96
M=100



Cache

	L1I miss rate (/Effective instruction)	Load-store instruction	L1D miss	L1D miss rate (/Load-store instruction)	L1D miss demand rate (%) (/L1D miss)	L1D miss hardware prefetch rate (%) (/L1D miss)	L1D miss software prefetch rate (%) (/L1D miss)	L2 miss	L2 miss rate (/Load-store instruction)	L2 miss demand rate (%) (/L2 miss)	L2 miss hardware prefetch rate (%) (/L2 miss)	L2 miss software prefetch rate (%) (/L2 miss)
改善前	0.00	2.12E+10	8.20E+08	0.04	78.62%	21.38%	0.00%	5.45E+03	0.00	87.63%	20.43%	0.00%
改善後	0.00	2.11E+10	7.39E+07	0.00	99.99%	0.01%	0.00%	4.12E+03	0.00	80.75%	32.43%	0.00%

配列aの最内ループでのアクセスはデータ局所性が低いため、浮動小数点ロードL2アクセス待ちが発生しています。

改善前ソース

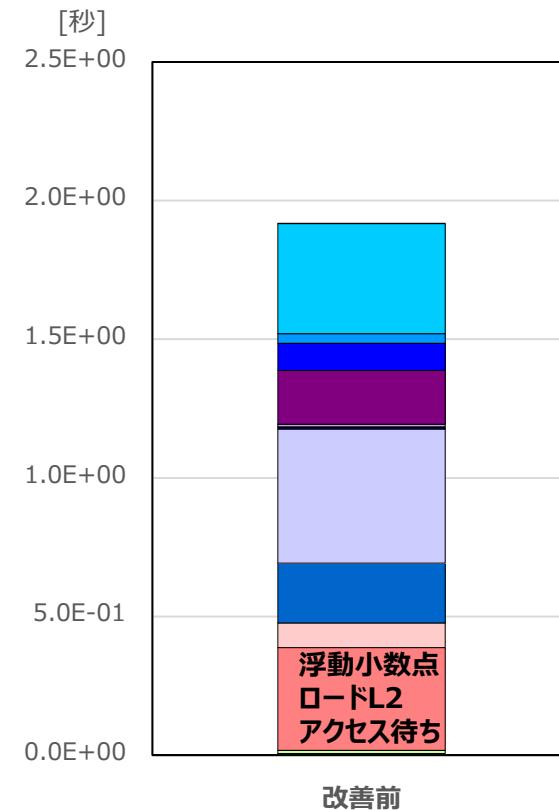
```

36 void sub(int N,int M,int ITER, double (* restrict a)[M][N])
37 {
38     int i,j,k;
39     #pragma omp parallel private(i,j,k)
40     {
41         for(k=0; k<ITER; k++)
42         {
43             #pragma omp for nowait
44             <<< Loop-information Start >>>
45             :
46             <<< Loop-information End >>>
47             p       for(j=0; j<M; j++)
48             p       {
49                 <<< Loop-information Start >>>
50                 :
51                 <<< Loop-information End >>>
52                 p       v   for(i=0; i<N; i++)
53                 p       v   {
54                     a[0][j][i]=a[1][j][i]+a[2][j][i]+a[3][j][i] +
55                     a[4][j][i]+a[5][j][i]+a[6][j][i]+a[7][j][i];
56                 }
57             }
58         return;
59     }

```

N=96
M=100

配列宣言
double a[8][M][N];



Cache	L1 miss rate (/Effective instruction)	Load-store instruction	L1D miss	L1D miss rate (/Load-store instruction)	L1D miss demand rate (%) (/L1D miss)	L1D miss hardware prefetch rate (%) (/L1D miss)	L1D miss software prefetch rate (%) (/L1D miss)	L2 miss	L2 miss rate (/Load-store instruction)	L2 miss demand rate (%) (/L2 miss)	L2 miss hardware prefetch rate (%) (/L2 miss)	L2 miss software prefetch rate (%) (/L2 miss)
改善前	0.00	2.18E+10	8.20E+08	0.04	82.13%	17.87%	0.00%	1.49E+04	0.00	81.54%	38.68%	0.00%

配列次元移動を行うことでデータの局所性が向上され、浮動小数点L2アクセス待ちが改善されました。

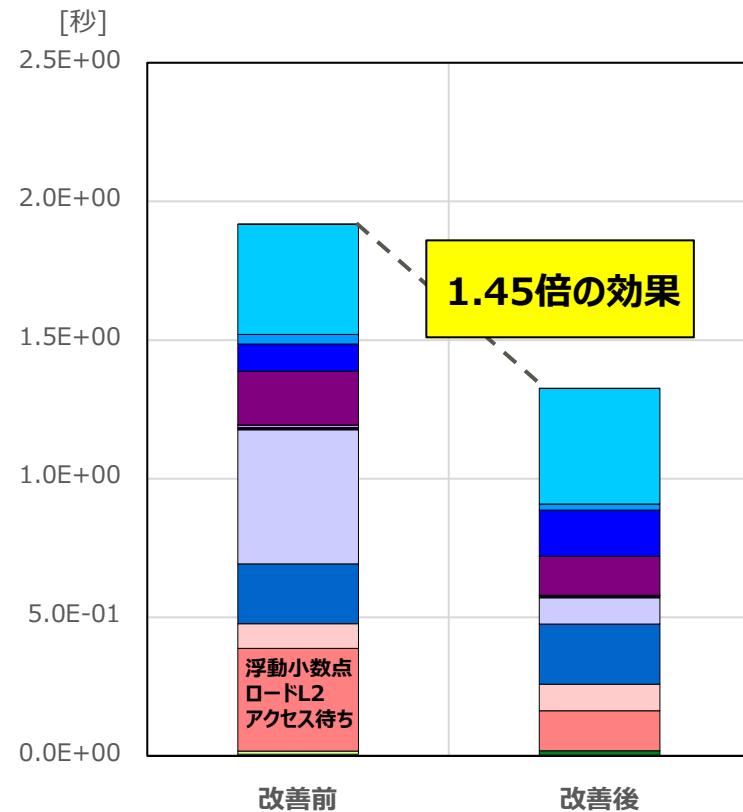
改善後ソース（ソースチューニング）

```

36     void sub(int N,int M,int ITER, double (* restrict a)[8][N]) {
37         int i,j,k;
38         #pragma omp parallel private(i,j,k)
39         {
40             for(k=0; k<ITER; k++) {
41                 #pragma omp for nowait
42                 <<< Loop-information Start >>>
43                 :
44                 <<< Loop-information End >>>
45                 p     for(j=0; j<M; j++) {
46                     <<< Loop-information Start >>>
47                     :
48                     <<< Loop-information End >>>
49                     p     v     for(i=0; i<N; i++) {
50                         p     v     a[j][0][i]=a[j][1][i]+a[j][2][i]+a[j][3][i]+
51                             a[j][4][i]+a[j][5][i]+a[j][6][i]+a[j][7][i];
52                     }
53                 }
54             return;
55         }
56     }

```

N=96
M=100
配列宣言
double a[8][M][N];



Cache	L1I miss rate (/Effective instruction)	Load-store instruction	L1D miss	L1D miss rate (/Load-store instruction)	L1D miss demand rate (%) (/L1D miss)	L1D miss hardware prefetch rate (%) (/L1D miss)	L1D miss software prefetch rate (%) (/L1D miss)	L2 miss	L2 miss rate (/Load-store instruction)	L2 miss demand rate (%) (/L2 miss)	L2 miss hardware prefetch rate (%) (/L2 miss)	L2 miss software prefetch rate (%) (/L2 miss)
改善前	0.00	2.18E+10	8.20E+08	0.04	82.13%	17.87%	0.00%	1.49E+04	0.00	81.54%	38.68%	0.00%
改善後	0.00	2.09E+10	8.63E+07	0.00	99.97%	0.02%	0.00%	1.56E+04	0.00	86.83%	39.15%	0.00%

以下の翻訳オプション（Fortran固有）を指定することで、ソースチューニングと同等の効果を得ることができます。

翻訳オプション	機能説明
-Karray_subscript	4次元以上の割付け配列および最終次元の要素が10以下で、最終次元以外の要素が100以上の4次元以上の配列に対して配列の次元移動を行うことを指示します。 -Karray_subscript_element=100, -Karray_subscript_elementlast=10 および -Karray_subscript_rank=4 も同時に有効になります。
-Karray_subscript_element=N ($2 \leq N \leq 2,147,483,647$)	配列の次元移動の対象となる配列の最終次元以外の要素数がN以上であることを指示します。本オプションは、-Karray_subscript オプションが有効な場合に意味があります。ただし、割付け配列では、本オプションの意味はありません。
-Karray_subscript_elementlast=N ($2 \leq N \leq 2,147,483,647$)	配列の次元移動の対象となる配列の最終次元の要素数がN以下であることを指示します。本オプションは、-Karray_subscript オプションが有効な場合に意味があります。ただし、割付け配列では、本オプションの意味はありません。
-Karray_subscript_rank=N ($2 \leq N \leq 30$)	配列の次元移動の対象となる配列の次元数がN次元以上であることを指示します。本オプションは、-Karray_subscript オプションが有効な場合に意味があります。

■ 使用例（改善前ソース時）

```
$frtpx -Kfast,parallel sample.f90
```

-Karray_subscript,array_subscript_rank=2,array_subscript_element=2

■ 注意事項

- 対象配列を使うソース全てにオプション指定が必要です。
- 移動の効果はプログラムによって異なります。
- 正しくない使い方をした場合には計算結果が異なる場合があります。

アンロールアンドジャム

- ・ アンロールアンドジャムとは
- ・ アンロールアンドジャム（改善前）
- ・ アンロールアンドジャムの効果（最適化制御行チューニング）

アンロールアンドジャムとは、多重ループの外側のループをアンローリングにより n 重に展開し、さらに展開された内側のループを融合するものです。

最適化前ソース	最適化後ソース
<pre> !OCL UNROLL_AND_JAM_FORCE(2) DO J=1, 128 DO I=1, 128 A(I, J) = B(I, J) + B(I, J+1) ... END DO END DO : </pre>	<pre> DO J=1, 128, 2 DO I=1, 128 A(I, J) = B(I, J) + B(I, J+1) A(I, J+1) = B(I, J+1) + B(I, J+2) ... END DO END DO </pre>

外側ループのアンローリング

内側ループの融合
(共通式の除去)

アンロールアンドジャムを適用することで、共通式の除去が促進され、実行性能が向上する場合があります。

ただし、データストリーム数の増加やデータのアクセス順序の変化は、キャッシュの利用効率を低下させ、実行性能が低下する場合があります。

以下の最適化制御行を指定します。

最適化指示子 (Fortran)	意味	指定可能な最適化制御行			
		プログラム単位	DOループ単位	文単位	配列代入文単位
UNROLL_AND_JAM[(n)]	最適化の効果があると判断したループに対してアンロールアンドジャムを有効にします。nは展開数(多深度)を表す2~100の10進数です。	○	○	×	×
UNROLL_AND_JAM_FORCE[(n)]	アンロールアンドジャムを有効にします。nは展開数(多深度)を表す2~100の10進数です。	×	○	×	×
NOUNROLL_AND_JAM	アンロールアンドジャムを無効にします。	○	○	×	×

最適化指示子 (C/C++)	意味	指定可能な最適化制御行			
		global行	procedure行	loop行	statement行
unroll_and_jam[(n)]	最適化の効果があると判断したループに対してアンロールアンドジャムを有効にします。nは展開数(多深度)を表す2~100の10進数です。	○	○	○	×
unrool_and_jam_force[(n)]	アンロールアンドジャムを有効にします。nは展開数(多深度)を表す2~100の10進数です。	×	×	○	×
nounroll_and_jam	アンロールアンドジャムを無効にします。	○	○	○	×

■ 注意事項

- UNROLL_AND_JAM指示子は最適化の効果が期待できない場合や、回転を跨いだデータ依存がある場合は最適化を行いません。
- UNROLL_AND_JAM_FORCE指示子を誤って指定した（繰り返しを跨いだデータ依存がある）場合、実行結果は保証されません。
- 最内ループはアンロールアンドジャムの対象なりません。
- 最内ループの繰り返し回数が小さい場合、データストリーム数の増加やデータのアクセス順序の変化によっては、キャッシュの利用効率を低下させ、実行性能が低下する場合があります。
- キャッシュミスが増加する場合はプリフェッチで補正する良くなる場合があります。

以下の翻訳オプションを指定することで、最適化制御行チューニングと同等の効果を得ることができます。

翻訳オプション	機能説明
-K{ unroll_and_jam[=N] nounroll_and_jam } $2 \leq N \leq 100$	<p>アンロールアンドジャムの最適化を行つかどうかを指示します。Nにはループ展開数の上限を2~100の範囲で指定できます。Nの指定を省略した場合、コンパイラが自動的に最良な値を決定します。デフォルトは-Kounroll_and_jamです。</p> <p>アンロールアンドジャムを適用することで、共通式の除去が促進され、実行性能が向上する場合があります。ただし、データストリーム数の増加やデータのアクセス順序の変化は、キャッシュの利用効率を低下させ、実行性能が低下する場合があります。</p> <p>また、アンロールアンドジャムの最適化による実行性能への影響はループ単位で異なります。このため、本最適化は、-Kunroll_and_jam[=N]オプションでプログラム全体へ適用せず、最適化指示子UNROLL_AND_JAMや最適化指示子UNROLL_AND_JAM_FORCEでループ単位に適用することを推奨します。</p>

■使用例（改善前ソース時）

```
$ frtpx -Kfast,parallel sample.f90 -Kunroll_and_jam
$ fccpx -Kfast,parallel sample.c -Kunroll_and_jam
```

◆注意事項

clangモードではアンロールアンドジャムの最適化機能を使用できません

L1Dミス率が高いことからキャッシュの利用効率が悪く、浮動小数点ロードL2アクセス待ちが多くなっています。

改善前ソース

```

11      DATA_TYPE,dimension(IMAX,JMAX,KMAX)::a
12      DATA_TYPE,dimension(JMAX,KMAX)::b
13      DATA_TYPE,dimension(JMAX,IMAX)::c \
:
15 1 pp    do k=1, KMAX
16 1
17 2 p     do j=1, JMAX - 3
    <<< Loop-information Start >>>
    <<< [OPTIMIZATION]
    <<< SIMD(VL: 8)
    <<< SOFTWARE PIPELINING(IPC: 1.85, ITR: 80, MVE: 2, POL: S)
    <<< PREFETCH(HARD) Expected by compiler :
    <<< a
    <<< Loop-information End >>>
18 3 p 2v   do i =1, IMAX
19 3 p 2v     a(i,j,k) = a(i,j+1,k) + a(i,j+2,k) + a(i,j+3,k) &
20 3           + (b(j+2,k) / c(j,i))
21 3 p 2v   end do
22 2 p
23 1 p   end do

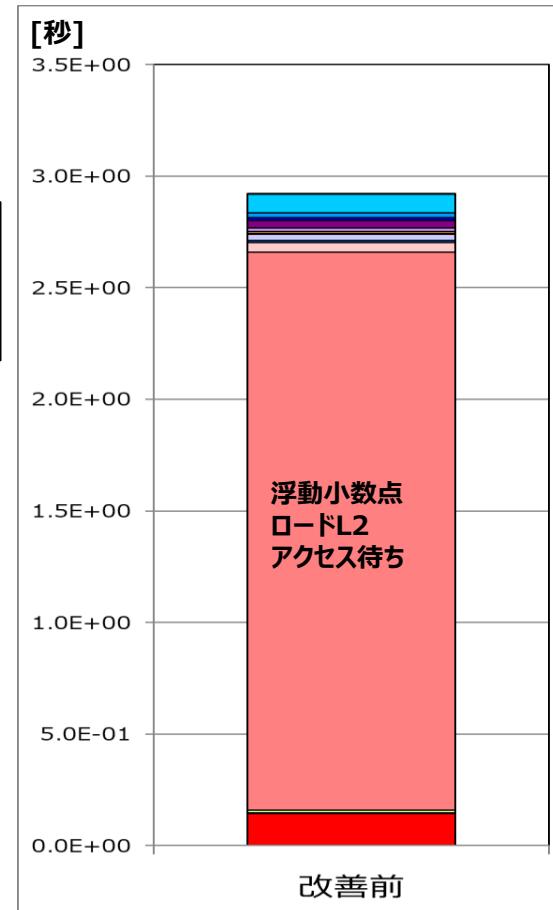
```

```

#define DATA_TYPE real(kind=8)
#define IMAX 512
#define JMAX 512
#define KMAX 128

```

配列cはストライドアクセスのため、
キャッシュの利用効率が悪い



Cache	L1I miss rate (/Effective instruction)	Load-store instruction	L1D miss	L1D miss rate (/Load-store instruction)	L1D miss demand rate (%) (/L1D miss)	L1D miss hardware prefetch rate (%) (/L1D miss)	L1D miss software prefetch rate (%) (/L1D miss)	L2 miss	L2 miss rate (/Load-store instruction)	L2 miss demand rate (%) (/L2 miss)
改善前	0.00	3.39E+09	3.82E+09	1.13	99.50%	0.50%	0.00%	1.15E+08	0.03	46.67%

キャッシュの利用効率の向上により、L1Dミス数が減少して浮動小数点ロードL2アクセス待ちが改善されました。

改善後ソース（最適化制御行チューニング）

```

11      DATA_TYPE,dimension(IMAX,JMAX,KMAX)::a
12      DATA_TYPE,dimension(JMAX,KMAX)::b
13      DATA_TYPE,dimension(JMAX,IMAX)::c
14
15      1 pp      do k=1, KMAX
16      1         !ocl unroll_and_jam_force(8)
17      2 p      do j=1, JMAX - 3
18      3 p 2v      <<< Loop-information Start >>>
19      3 p 2v      <<< [OPTIMIZATION]
20      3 p 2v      <<< SIMD(VL: 8)
21      3 p 2v      <<< SOFTWARE PIPELINING
22      2 p      <<< PREFETCH(HARD) Exp
23      1 p      <<< a
24      1         PREFETCH(SOFT) : 32
25      1         SEQUENTIAL : 32
26      1         a: 32
27      1         SPILLS :
28      1             GENERAL : SPILL 0 FILL 4
29      1             SIMD&FP : SPILL 0 FILL 0
30      1             SCALABLE : SPILL 0 FILL 0
31      1             PREDICATE : SPILL 0 FILL 0
32      1         <<< Loop-information End >>>
33      3 p 2v      do i =1, IMAX
34      3 p 2v          a(i,j,k) = a(i,j+1,k) + a(i,j+2,k) + a(i,j+3,k) &
35      3 p 2v                  + (b(j+2,k) / c(j,i))
36      3 p 2v          end do
37      2 p      end do
38      1 p      end do

```

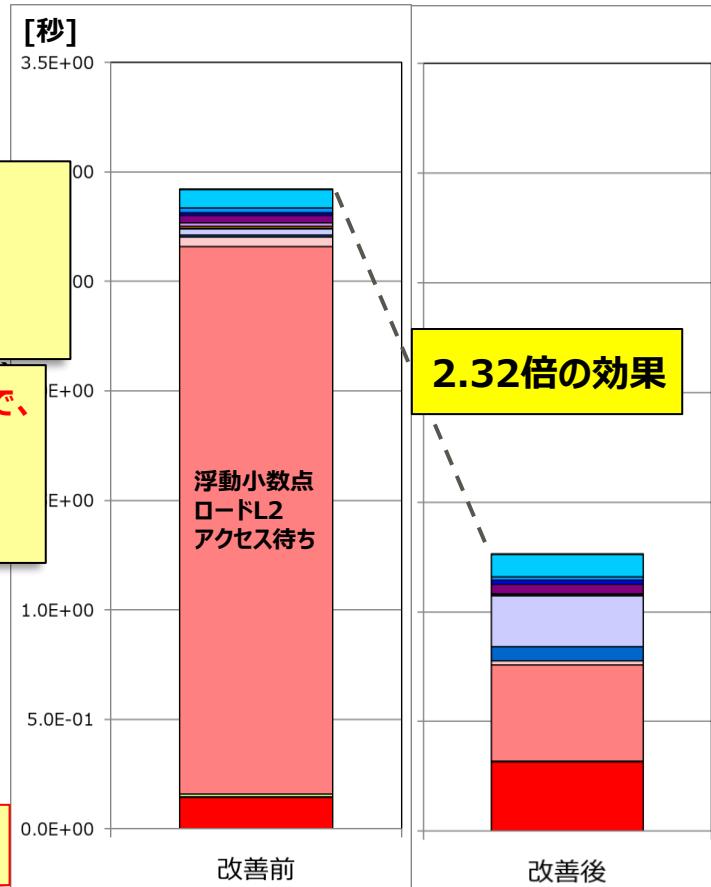
```

#define DATA_TYPE
real(kind=8)
#define IMAX 512
#define JMAX 512
#define KMAX 128

```

外側ループでアンロールすることで、配列aの共通式の除去による命令の削減のほか、配列aおよび配列cのキャッシュ利用効率が向上した

L1Dミスが大幅に減少



Cache	L1I miss rate (/Effective instruction)	Load-store instruction	L1D miss	L1D miss rate (/Load-store instruction)	L1D miss demand rate (%) (/L1D miss)	L1D miss hardware prefetch rate (%) (/L1D miss)	L1D miss software prefetch rate (%) (/L1D miss)	L2 miss	L2 miss rate (/Load-store instruction)	L2 miss demand rate (%) (/L2 miss)
改善前	0.00	3.39E+09	3.82E+09	1.13	99.50%	0.50%	0.00%	1.15E+08	0.03	46.67%
改善後	0.00	2.68E+09	7.37E+08	0.27	89.34%	2.55%	8.12%	1.15E+08	0.04	44.44%

L1Dミス率が高いことからキャッシュの利用効率が悪く、浮動小数点ロードL2アクセス待ちが多くなっています。

改善前ソース

```

69 p      #pragma omp parallel for
70   for (k=0;k<KMAX;k++){
    <<< Loop-information Start >>>
    <<< [OPTIMIZATION]
    <<< PREFETCH(HARD) Expected by compiler :
    <<< (unknown)
    <<< Loop-information End >>>
    for (j=0;j<JMAX-3;j++){
    <<< Loop-information Start >>>
    <<< [OPTIMIZATION]
    <<< SIMD(VL: 8)
    <<< SOFTWARE PIPELINING(IPC: 2.09, ITR: 80, MVE: 2, POL: S)
    <<< PREFETCH(HARD) Expected by compiler :
    <<< (unknown)
    <<< Loop-information End >>>
72 p      2v      for (i=0;i<IMAX;i++){
73 p      2v      a[k][j][i] = a[k][j+1][i] + a[k][j+2][i] + a[k][j+3][i]
74 p      2v          + (b[k][j+2] / c[i][j]);
75 p      2v      }
76 p      }
77 p      }

```

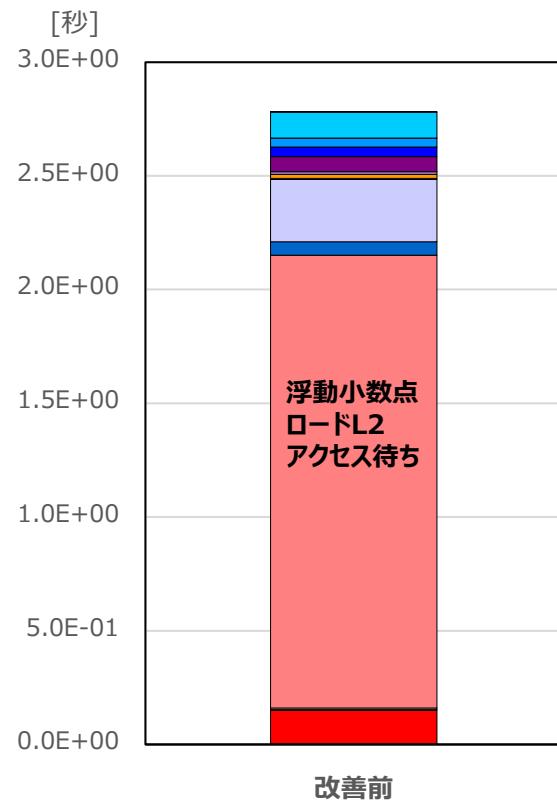
```

#define IMAX 512
#define JMAX 512
#define KMAX 128

配列宣言
double
a[KMAX][JMAX][IMAX],
b[KMAX][JMAX],
c[IMAX][JMAX];

```

配列cはストライドアクセスのため、
キャッシュの利用効率が悪い



改善前

Cache	L1I miss rate (/Effective instruction)	Load-store instruction	L1D miss	L1D miss rate (/Load-store instruction)	L1D miss demand rate (%) (/L1D miss)	L1D miss hardware prefetch rate (%) (/L1D miss)	L1D miss software prefetch rate (%) (/L1D miss)	L2 miss	L2 miss rate (/Load-store instruction)	L2 miss demand rate (%) (/L2 miss)
改善前	0.00	3.50E+09	3.68E+09	1.05	99.73%	0.26%	0.00%	1.16E+08	0.03	59.11%

キャッシュの利用効率の向上により、L1Dミス数が減少して浮動小数点ロードL2アクセス待ちが改善されました。

改善後ソース（最適化制御行チューニング）

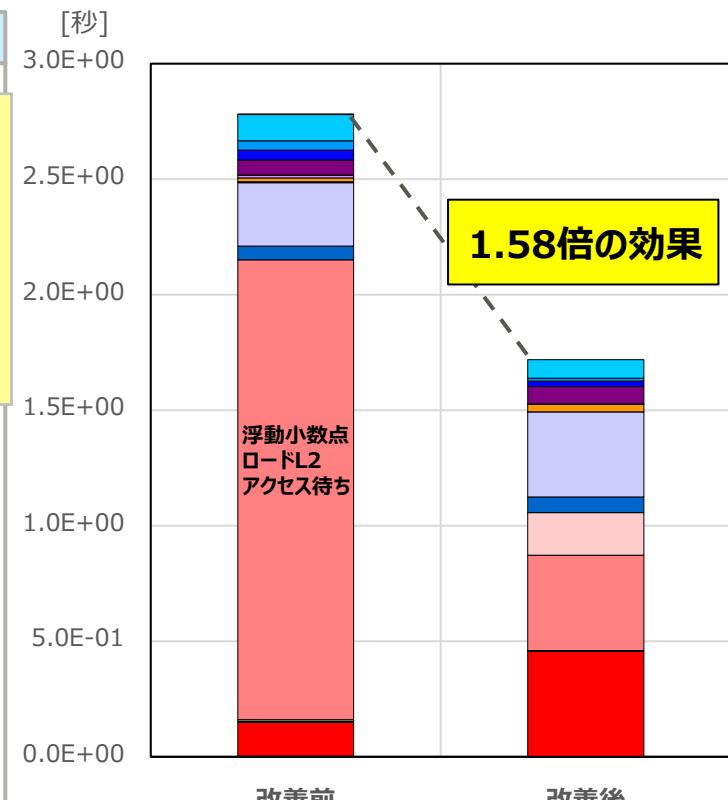
```

74     for (k=0;k<KMAX;k++){
75         #pragma loop unroll_and_jam_force 8
76         <<< Loop-information Start >>>
77         <<< [OPTIMIZATION]
78         <<< PREFETCH(HARD) Expected by computer :
79             (unknown)
80         <<< Loop-information End >>>
81         for (j=0;j<JMAX-3;j++){
82             <<< Loop-information Start >>>
83             <<< [OPTIMIZATION]
84             <<< SOFTWARE PIPELINING(IPC: 0.53,
85             <<< PREFETCH(HARD) Expected by computer :
86                 (unknown)
87             <<< PREFETCH(SOFT) : 32
88             <<< SEQUENTIAL : 32
89                 (unknown): 32
90             <<< SPILLS :
91                 GENERAL : SPILL 0
92                 SIMD&FP : SPILL 0
93                 SCALABLE : SPILL 0 FILL 0
94                 PREDICATE : SPILL 0 FILL 0
95             <<< Loop-information End >>>
96             2v         for (i=0;i<IMAX;i++){
97                 2v             a[k][j][i] = a[k][j+1][i] + a[k][j+2][i] + a[k][j+3][i]
98                     + (b[k][j+2] / c[i][j]);
99             }           }
100         }
101     }

```

外側ループでアンロールすることで、配列aの共通式の除去による命令の削減のほか、配列aおよび配列cのキャッシュ利用効率が向上した

L1Dミスが大幅に減少



Cache	L1 miss rate (/Effective instruction)	Load-store instruction	L1D miss	L1D miss rate (/Load-store instruction)	L1D miss demand rate (%) (/L1D miss)	L1D miss hardware prefetch rate (%) (/L1D miss)	L1D miss software prefetch rate (%) (/L1D miss)	L2 miss	L2 miss rate (/Load-store instruction)	L2 miss demand rate (%) (/L2 miss)
改善前	0.00	3.50E+09	3.68E+09	1.05	99.73%	0.26%	0.00%	1.16E+08	0.03	59.11%
改善後	0.00	2.42E+09	6.24E+08	0.26	82.30%	3.11%	14.60%	1.15E+08	0.05	39.17%

データアクセス待ち（レイテンシの隠蔽）

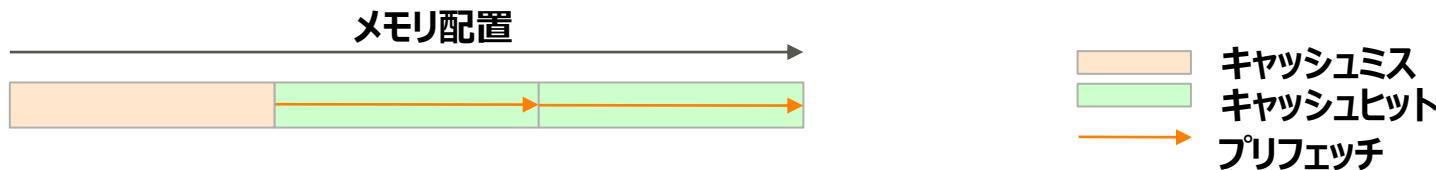
- ・ インダイレクトアクセスプリフェッヂ
- ・ 連続アクセスでないデータへのソフトウェアプリフェッヂの利用

インダイレクトアクセスプリフェッチ

- プリフェッチとは
- インダイレクトアクセスプリフェッチ（最適化制御行）
- インダイレクトアクセスプリフェッチの効果（翻訳オプションチューニング）
- インダイレクトアクセスプリフェッチ（改善前）
- インダイレクトアクセスプリフェッチの効果（最適化制御行チューニング）

プリフェッчとは

プリフェッчとは、実行される命令によってデータが必要になる前に、キャッシュにデータをロードしておくことで、パフォーマンスを向上させる仕組みです。

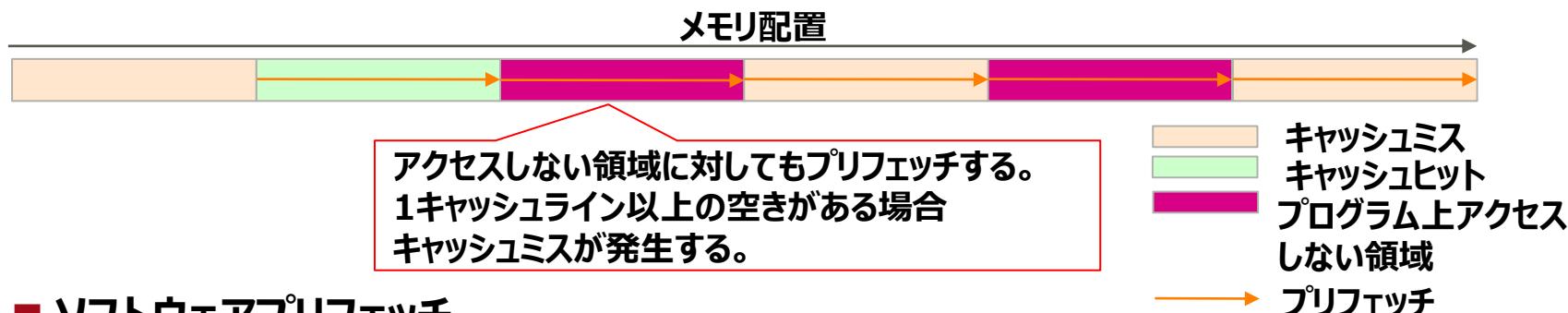


プリフェッчにはハードウェアプリフェッチとソフトウェアプリフェッチがあります。

■ ハードウェアプリフェッチ

プログラムのメモリアクセスの規則性からハードウェアがデータアクセスを予測しプリフェッチします。

プログラム上アクセスしない領域がある場合にもプリフェッチしてしまうため、そのようなプログラムではキャッシュ効率が非常に悪くなります。その場合はソフトウェアプリフェッチを利用します。



■ ソフトウェアプリフェッチ

ソフトウェア（コンパイラ）がプログラムを解析し、`prefetch`命令を生成することでプリフェッチします。

以下の最適化制御行を指定します。

最適化指示子 (Fortran)	意味	指定可能な最適化制御行			
		プログ ラム 単位	DO ループ 単位	文単位	配列代 入文 単位
PREFETCH	コンパイラの自動 prefetch 機能を有効にします。	○	○	×	×

最適化指示子 (C/C++)	意味	指定可能な最適化制御行			
		global 行	proce dure 行	loop行	state ment 行
prefetch	コンパイラの自動 prefetch 機能を有効にします。	○	○	○	×

◆補足

prefetch最適化指示子は、以下の翻訳オプションを指定したものと等価です。

-Kprefetch_sequential,prefetch_stride,prefetch_indirect,prefetch_conditional,
prefetch_cache_level=all

◆注意事項

翻訳オプション-Kprefetch_sequential、-Kprefetch_stride、
-Kprefetch_indirectまたは-Kprefetch_conditionalが有効な場合のプリフェッ
チでは、ループのキャッシュ効率、分岐の有無または添字の複雑さによって、実行性能が
低下することがあります。

以下の翻訳オプションを指定することで、最適化制御行チューニングと同等の効果を得ることができます。

翻訳オプション	機能説明
-Kprefetch间接	ループ内で使用される間接的にアクセス（リストアクセス）される配列データに対して、 prefetch命令を使用したオブジェクトを生成するかどうかを指示します。 本オプションは、-O1 オプション以上が有効な場合に意味があります。 デフォルトは、-Kprefetch_noindirect です。

■ 使用例（改善前ソース時）

```
$ frtpx -Kfast,parallel sample.f90 -Kprefetch间接  
$ fccpx -Kfast,parallel sample.c -Kprefetch间接
```

◆ 注意事項

プリフェッчが実施されますが、ループのキャッシュ効率やIF構文の有無、
添字の複雑さによっては、意図している効果が得られない場合があります。
clangモードではインダイレクトプリフェッчの最適化機能を使用できません

インダイレクトアクセス（リストアクセス）の場合、推奨オプションではプリフェッチが生成されないためメモリアクセスのレイテンシが見えています。
そのため、浮動小数点ロードL2アクセス待ちが多くなっています。

改善前ソース

```

<<< Loop-information Start >>>
<<< [PARALLELIZATION]
<<< Standard iteration count: 572
<<< [OPTIMIZATION]
<<< SOFTWARE PIPELINING(IPC: 3.50, ITR: 18, MVE: 3, POL: S)
<<< PREFETCH(HARD) Expected by compiler :
<<< e, d, a
<<< Loop-information End >>>
52 1 pp 2      do i = 1 , n
53 1 p 2      a(i) = b(d(i)) + scalar * c(e(i))
54 1 p 2      enddo

```

配列b,cが
インダイレクトアクセス



L1Dミスd m率および、L2ミスdm率が高くプリフェッチが効いていない。
メモリループットには余裕があるため性能向上の可能性がある。

Cache	L1I miss rate (/Effective instruction)	Load-store instruction	L1D miss	L1D miss rate (/Load-store instruction)	L1D miss demand rate (%) (/L1D miss)	L1D miss hardware prefetch rate (%) (/L1Dmiss)	L1D miss software prefetch rate (%) (/L1Dmiss)	L2 miss	L2 miss rate (/Load-store instruction)	L2 miss demand rate (%) (/L2 miss)	L2 miss hardware prefetch rate (%) (/L2miss)	L2 miss software prefetch rate (%) (/L2miss)
改善前	0.00	9.01E+09	3.77E+09	0.42	98.23%	1.77%	0.00%	4.29E+08	0.05	54.19%	75.86%	0.00%

prefetch指示子を指定することでインダイレクトアクセス（リストアクセス）に対してプリフェッチを生成します。その結果、浮動小数点ロードモリL2アクセス待ちが改善されました。

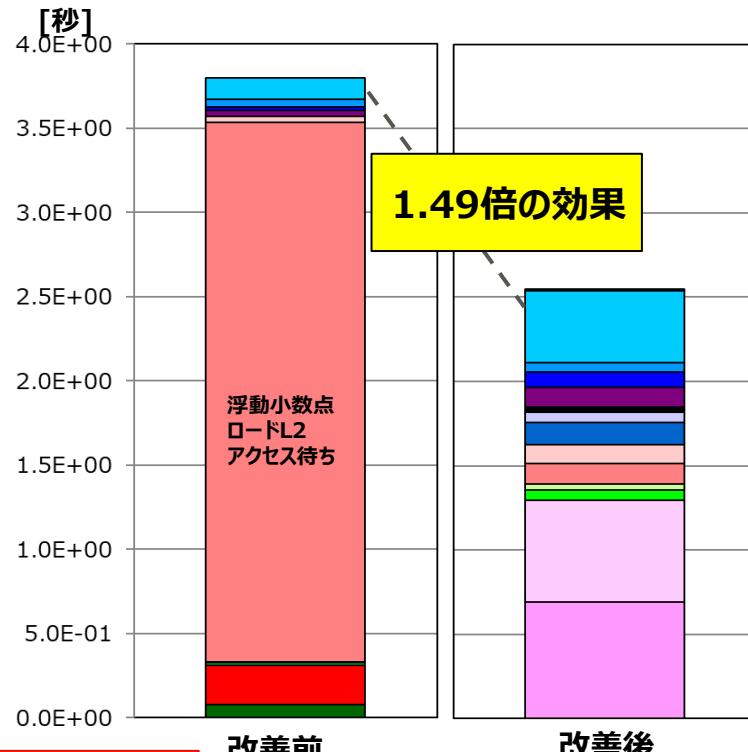
改善後ソース

```

51      !ocl prefetch
      <<< Loop-information Start >>>
      <<< [PARALLELIZATION]
      <<< Standard iteration count: 572
      <<< [OPTIMIZATION]
      <<< PREFETCH(HARD) E>
      <<< e, d, a
      <<< PREFETCH(SOFT). 8
      <<< INDIRECT : 8
      <<< c: 4, b: 4
      <<< Loop-information End >>>
52    1 pp 2      do i = 1 , n
53    1 p 2      a(i) = b(d(i)) + scalar * c(e(i))
54    1 p 2      enddo

```

インダイレクトアクセス
(配列b, c)に対するプリフェッチが生成された



インダイレクトアクセス(配列b, c)に対してプリフェッチ命令を
生成することでL1Dミスdm率、L2 ミスdm率が減った。

	Memory throughput (GB/s)
改善前	32.67
改善後	183.71

Cache	L1I miss rate (/Effective instruction)	Load-store instruction	L1D miss	L1D miss rate (/Load-store instruction)	L1D miss demand rate (%) (/L1D miss)	L1D miss hardware prefetch rate (%) (/L1D miss)	L1D miss software prefetch rate (%) (/L1D miss)	L2 miss	L2 miss rate (/Load-store instruction)	L2 miss demand rate (%) (/L2 miss)	L2 miss hardware prefetch rate (%) (/L2 miss)	L2 miss software prefetch rate (%) (/L2 miss)
改善前	0.00	9.01E+09	3.77E+09	0.42	98.23%	1.77%	0.00%	4.29E+08	0.05	54.19%	75.86%	0.00%
改善後	0.00	1.66E+10	3.82E+09	0.23	2.94%	2.94%	94.12%	1.77E+09	0.11	2.02%	6.20%	91.78%

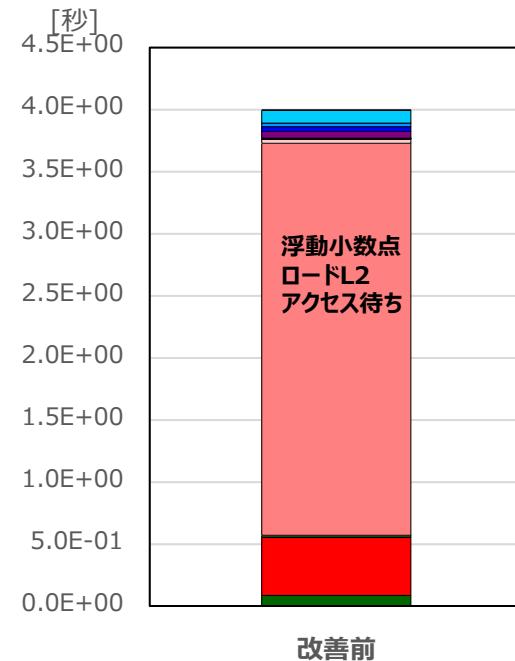
インダイレクトアクセス（リストアクセス）の場合、推奨オプションではプリフェッチが生成されないためメモリアクセスのレイテンシが見えています。
そのため、浮動小数点ロードL2アクセス待ちが多くなっています。

改善前ソース

```

53     void sub(double * restrict a, double * restrict b,
54             double * restrict c, int * restrict d, int * restrict e,
55             double scalar, int n){
56         int i;
57 #pragma omp parallel for
58 <<< Loop-information Start >>>
59 <<< [OPTIMIZATION]
60 <<< SOFTWARE PIPELINING(IPC: 2.83, ITR: 12, MVE: 2, POL: S)
61 <<< PREFETCH(HARD) Expected by compiler :
62 <<< (unknown)
63 <<< Loop-information End >>>
64 p 2   for (i = 0; i < n; i++){
65 p 2     a[i] = b[d[i]] + scalar * c[e[i]];
66 p 2   }
67 }
```

配列b,cが
インダイレクトアクセス



L1Dミスd m率および、L2ミスdm率が高くプリフェッチが効いていない。
メモリループットには余裕があるため性能向上の可能性がある。

Statistics	Memory throughput (GB/s)
改善前	31.24

Cache	L1 miss rate (/Effective instruction)	Load-store instruction	L1D miss	L1D miss rate (/Load-store instruction)	L1D miss demand rate (%) (/L1D miss)	L1D miss hardware prefetch rate (%) (/L1D miss)	L1D miss software prefetch rate (%) (/L1D miss)	L2 miss	L2 miss rate (/Load-store instruction)	L2 miss demand rate (%) (/L2 miss)	L2 miss hardware prefetch rate (%) (/L2 miss)	L2 miss software prefetch rate (%) (/L2 miss)
改善前	0.00	8.65E+09	3.78E+09	0.44	98.01%	1.99%	0.00%	4.30E+08	0.05	47.96%	100.00%	0.00%

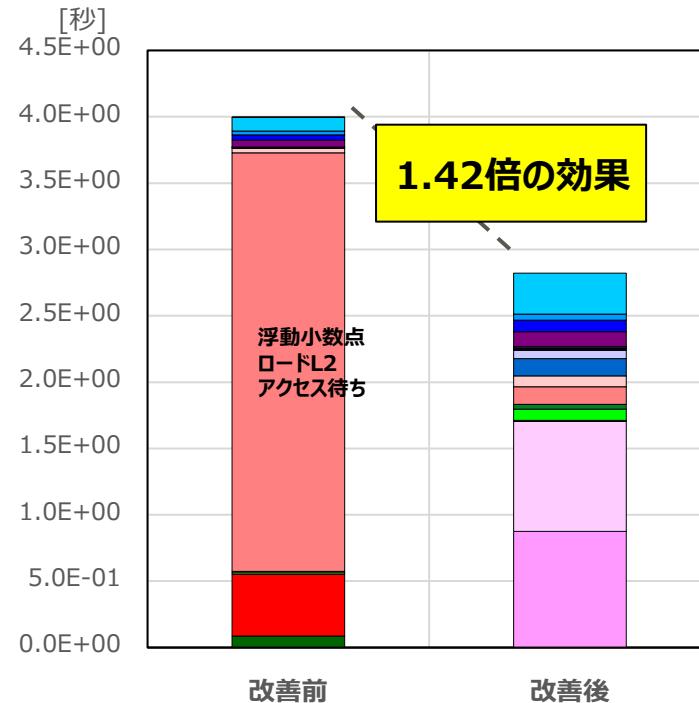
prefetch指示子を指定することでインダイレクトアクセス（リストアクセス）に対してプリフェッチを生成します。その結果、浮動小数点ロードによるL2アクセス待ちが改善されました。

改善後ソース

```

53 void sub(double * restrict a, double * restrict b,
          double * restrict c, int * restrict d,
          int * restrict e, double scalar, int n){
54     int i;
56     #pragma loop prefetch
57     #pragma omp parallel for
<<< Loop-information Start >>>
<<< [OPTIMIZATION]
<<< PREFETCH(HARD) Expected by compiler :
<<< (unknown)
<<< PREFETCH(SOFT) : 8
<<< INDIRECT : 8
<<< (unknown): 8
<<< Loop-information End >>>
58 p 2   for (i = 0; i < n; i++){
59 p 2     a[i] = b[d[i]] + scalar * c[e[i]];
60 p 2   }
61 }
```

インダイレクトアクセス
(配列b, c)に対するプリフェッチが生成された



インダイレクトアクセス(配列b, c)に対してプリフェッチ命令を
生成することでL1Dミスdm率、L2 ミスdm率が減った。

	Statistics	Memory throughput (GB/s)
改善前		31.24
改善後		159.90

Cache	L1 miss rate (/Effective instruction)	Load-store instruction	L1D miss	L1D miss rate (/Load-store instruction)	L1D miss demand rate (%) (/L1D miss)	L1D miss hardware prefetch rate (%) (/L1D miss)	L1D miss software prefetch rate (%) (/L1D miss)	L2 miss	L2 miss rate (/Load-store instruction)	L2 miss demand rate (%) (/L2 miss)	L2 miss hardware prefetch rate (%) (/L2 miss)	L2 miss software prefetch rate (%) (/L2 miss)
改善前	0.00	8.65E+09	3.78E+09	0.44	98.01%	1.99%	0.00%	4.30E+08	0.05	47.96%	100.00%	0.00%
改善後	0.00	1.70E+10	3.83E+09	0.22	2.99%	2.94%	94.08%	1.70E+09	0.10	1.80%	6.41%	91.79%

連続アクセスでないデータへのソフトウェアプリフェッヂの利用

- ・ 連続アクセスでないデータへのソフトウェアプリフェッヂの利用
- ・ 連続アクセスでないデータへのソフトウェアプリフェッヂの利用（改善前）
- ・ 連続アクセスでないデータへのソフトウェアプリフェッヂの利用①（最適化制御行チューニング）
- ・ 連続アクセスでないデータへのソフトウェアプリフェッヂの利用②[推奨]（最適化制御行チューニング）

非連続または連続な部分の短いデータに対するアクセスでは、ハードウェアプリフェッチはキャッシュミスが発生しやすくなります。その場合は、ソフトウェアプリフェッチを利用してすることで性能が向上します。

ソフトウェアプリフェッチを利用するためには、後述の**PREFETCH_READ指示子**, **PREFETCH_WRITE指示子**を使用するか、下記の翻訳時オプションを利用します。

翻訳時オプション	機能説明
-Kprefetch_sequential=auto	ループ内で使用される連続的にアクセスされる配列データに対して、ハードウェアプリフェッチを利用するか、prefetch命令を出力するかをコンパイラが自動的に選択します。本オプションは、-O1オプション以上が有効な場合に意味があります。-O2オプション以上が有効な場合、デフォルトは-Kprefetch_sequential=autoです。
-Kprefetch_sequential=soft	ループ内で使用される連続的にアクセスされる配列データに対して、ハードウェアプリフェッチを利用せずに、prefetch命令を出力します。本オプションは、-O1オプション以上が有効な場合に意味があります。
-Kprefetch_nosequential	ループ内で使用される連続的にアクセスされる配列データに対して、prefetch命令を使用せずにオブジェクトを生成します。-O0または-O1オプションが有効な場合、デフォルトは -Kprefetch_nosequentialです。

PREFETCH_READ指示子, PREFETCH_WRITE指示子は次のように指定します。

最適化指示子	意味	指定可能な最適化制御行			
		プログ ラム 単位	DO ループ 単位	文単位	配列代 入文 単位
PREFETCH_READ(name[,level={1 2}] [,strong={0 1}])	参照されているデータに対してprefetch命令を生成することを指示します。 nameは配列要素名,levelはプリフェッチを行うキャッシュレベル、strongはstrong prefetchをするかどうかです。	○	×	○	×
PREFETCH_WRITE(name[,level={1 2}] [,strong={0 1}])	定義されているデータに対してprefetch命令を生成することを指示します。	○	×	○	×

改善前

```
do j=1,n
  do i=1, isize
    a(i,j) = b(i,j) + scalar * c(i,j)
  enddo
enddo
```



改善後① (ソフトウェアプリフェッチ要素指定)

```
do j=1,n
  do i=1, isize
    !OCL PREFETCH_WRITE(a(i,j+1),level=1)
    !OCL PREFETCH_READ(b(i,j+1),level=1)
    !OCL PREFETCH_READ(c(i,j+1),level=1)
    a(i,j) = b(i,j) + scalar * c(i,j)
  enddo
enddo
```

推奨

改善後② (ソフトウェアプリフェッチベクトル指定)

```
do j=1,n
  !OCL PREFETCH_WRITE(a(1:isize,j+1),level=1)
  !OCL PREFETCH_READ(b(1:isize,j+1),level=1)
  !OCL PREFETCH_READ(c(1:isize,j+1),level=1)
  do i=1, isize
    a(i,j) = b(i,j) + scalar * c(i,j)
  enddo
enddo
```

配列要素をベクトル指定します。
プリフェッチ命令を複数同時に生成できるため
要素ごとの指定よりも性能が向上します。



ビルトインプリフェッチ関数は次のように使用します。仕様についてはGNU C、C++コンパイラのホームページをご参考ください。

改善前

```
for (j = 0; j < n; j++) {  
    for (i = 0; i < isize; i++) {  
        a[j][i] = b[j][i] + scalar * c[j][i];  
    }  
}
```



改善後① (ビルトインプリフェッチ要素指定)

```
for (j = 0; j < n; j++) {  
    for (i = 0; i < isize; i++) {  
        __builtin_prefetch(&a[j+1][0], 1, 3);  
        __builtin_prefetch(&b[j+1][0], 0, 3);  
        __builtin_prefetch(&c[j+1][0], 0, 3);  
        __builtin_prefetch(&a[j+1][32], 1, 3);  
        :  
        a[j][i] = b[j][i] + scalar * c[j][i];  
    }  
}
```

プリフェッチ回数が少なく性能が向上します



改善後② (ビルトインプリフェッチベクトル指定)

```
for (j = 0; j < n; j++) {  
    __builtin_prefetch(&a[j+1][0], 1, 3);  
    __builtin_prefetch(&a[j+1][32], 1, 3);  
    __builtin_prefetch(&a[j+1][64], 1, 3);  
    __builtin_prefetch(&b[j+1][0], 0, 3);  
    __builtin_prefetch(&b[j+1][32], 0, 3);  
    __builtin_prefetch(&b[j+1][64], 0, 3);  
    __builtin_prefetch(&c[j+1][0], 0, 3);  
    :  
    for (i = 0; i < isize; i++) {  
        a[j][i] = b[j][i] + scalar * c[j][i];  
    }  
}
```

推奨

最内ループは、繰返し数が少なく、かつ繰返し数よりも配列サイズが大きいため、連續的ではありません。通常のプリフェッチではプリフェッチの立ち上がりのコストがみえてしまいます。そのため、浮動小数点ロードアクセス待ちが多くなっています。

改善前ソース

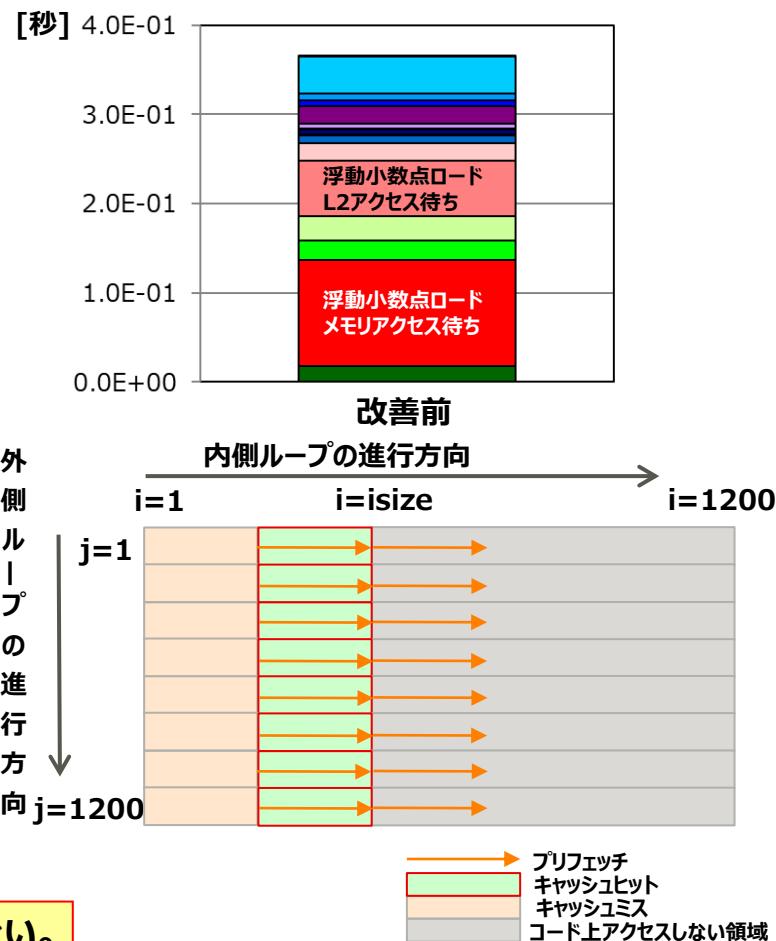
```

42      parameter(n=1200)
43      integer n
44      real*8 a(n,n),b(n,n),c(n,n),scalar
45      common /com/a,b,c
46
47      !$omp parallel do
48      <<< Loop-information :
49      <<< [OPTIMIZATION]
50      <<< PREFETCH(HARD)
51      <<< c, b, a
52      <<< Loop-information
53          do j=1,n
54              <<< Loop-information Start >>>
55              <<< [OPTIMIZATION]
56              <<< SIMD(VL: 8)
57              <<< SOFTWARE PIPELINING(IPC: 3.40, ITR: 128, MVE: 3, POL: S)
58              <<< PREFETCH(HARD) Expected by compiler :
59              <<< c, b, a
60              <<< Loop-information End >>>
61      1 p      do i=1,isize
62          a(i,j) = b(i,j) + scalar * c(i,j)
63      2 p 2v      enddo
64      2 p 2v      enddo
65      1 p      enddo

```

配列の1次元目の要素数 : 1200
ループ繰返し数(isize) : 128
外側ループjがインクリメントされた時に
アクセスの連續性が途切れる

L1Dミスdm率が高くプリフェッチが効いていない。



Cache	L1I miss rate (/Effective instruction)	Load-store instruction	L1D miss	L1D miss rate (/Load-store instruction)	L1D miss demand rate (%) (/L1D miss)	L1D miss hardware prefetch rate (%) (/L1D miss)	L1D miss software prefetch rate (%) (/L1D miss)	L2 miss	L2 miss rate (/Load-store instruction)	L2 miss demand rate (%) (/L2 miss)	L2 miss hardware prefetch rate (%) (/L2 miss)	L2 miss software prefetch rate (%) (/L2 miss)
改善前	0.00	1.75E+09	2.64E+08	0.15	70.94%	29.46%	-0.40%	1.26E+08	0.07	39.28%	83.10%	0.00%

PREFETCH_READ, PREFETCH_WRITE指示子を使用することで、外側ループの配列に対してプリフェッチを生成しプリフェッチの立ち上がりのコストを隠蔽します。その結果、浮動小数点ロードL2アクセス待ちが改善されました。

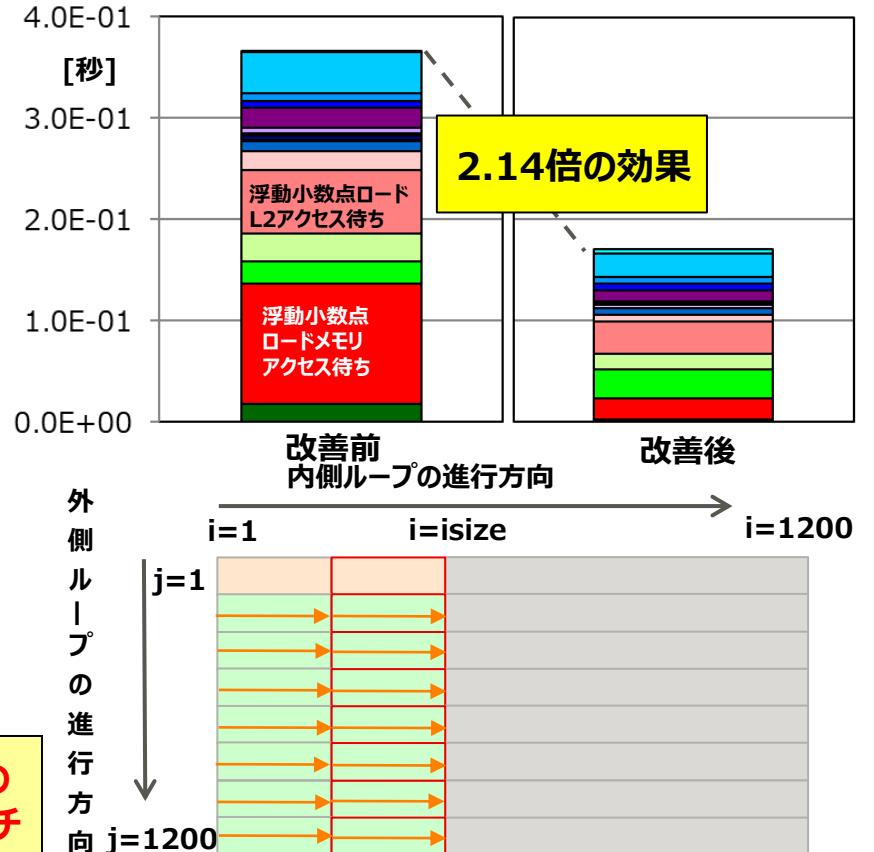
改善後ソース

```

42      parameter(n=1200)
43      integer n
44      real*8 a(n,n),b(n,n),c(n,n),scalar
45      common /com/a,b,c
46
47      !$omp parallel do
<<< Loop-information Start >>>
<<< [OPTIMIZATION]
<<< PREFETCH(HARD) Expected by compiler :
<<<   c, b, a
<<< PREFETCH(SOFT) : 18
<<< SPECIFIED : 18
<<<   a: 6, c: 6, b: 6
<<< Loop-information End >>>
48  1 p      do j=1,n
<<< Loop-information Start >>>
<<< [OPTIMIZATION]
<<< SIMD(VL: 8)
<<< SOFTWARE PIPELINING(IPC: 2.25, ITR: 88, MVE: 6, POL: S)
<<< PREFETCH(HARD) Expected by compiler :
<<<   c, b, a
<<< PREFETCH(SOFT) : 18
<<< SPECIFIED : 18
<<<   c: 6, b: 6, a: 6
<<< Loop-information End >>>
49  2 p      v      do i=1,isize
50  2 p      !OCL PREFETCH_WRITE(a(i,j+1),level=1)
51  2 p      !OCL PREFETCH_READ(b(i,j+1),level=1)
52  2 p      !OCL PREFETCH_READ(c(i,j+1),level=1)
53  2 p      v      a(i,j) = b(i,j) + scalar * c(i,j)
54  2 p      v      enddo
55  1 p      enddo

```

外側ループ1回転先の配列に対してプリフェッチ



L1Dミスdm率が減少

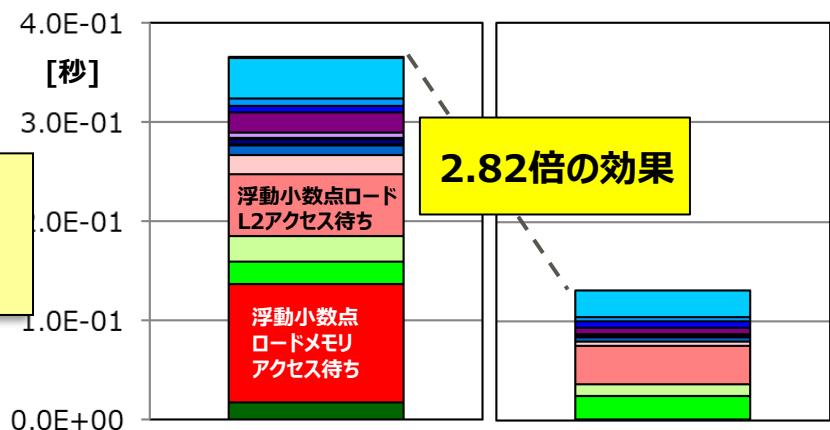
→ プリフェッチ
キャッシュヒット
キャッシュミス
コード上アクセスしない領域

Cache	Load-store instruction	L1D miss	L1D miss rate (/Load-store instruction)	L1D miss demand rate (%) (/L1D miss)	L1D miss hardware prefetch rate (%) (/L1D miss)	L1D miss software prefetch rate (%) (/L1D miss)
改善前	1.75E+09	2.64E+08	0.15	70.94%	29.46%	-0.40%
改善後	1.13E+09	1.90E+08	0.17	4.57%	0.44%	94.99%

外側のループで配列に対するPREFETCH_READ, PREFETCH_WRITE指示子を指定することで、最内ループでの命令数を減らすことができ、さらに性能を向上させることができます。

		改善後ソース	
48	1 p	<<< Loop-information Start >>> <<< [OPTIMIZATION] <<< PREFETCH(HARD) Expected by compiler : <<< c, b, a <<< PREFETCH(SOFT) : 3 <<< SPECIFIED : 3 <<< a: 1, c: 1, b: 1 <<< Loop-information End > do j=1,n <<< Loop-information Start > <<< [OPTIMIZATION] <<< PREFETCH(SOFT) : 4 <<< SPECIFIED : 4 <<< a: 4 <<< Loop-information End >>> !OCL PREFETCH_WRITE(a(1: isize, j+1), level=1)	
49	1 p	4s !OCL PREFETCH_WRITE(a(1: isize, j+1), level=1) <<< Loop-information Start >>> <<< [OPTIMIZATION] <<< PREFETCH(SOFT) : 4 <<< SPECIFIED : 4 <<< b: 4 <<< Loop-information End >>>	
50	1 p	4s !OCL PREFETCH_READ(b(1: isize, j+1), level=1) <<< Loop-information Start >>> <<< [OPTIMIZATION] <<< PREFETCH(SOFT) : 4 <<< SPECIFIED : 4 <<< c: 4 <<< Loop-information End >>>	
51	1 p	4s !OCL PREFETCH_READ(c(1: isize, j+1), level=1) <<< Loop-information Start >>> <<< [OPTIMIZATION] <<< SIMD(VL: 8) <<< SOFTWARE PIPELINING(IPC: 3.40, ITR: 128, MVE: 3, POL: S) <<< PREFETCH(HARD) Expected by compiler : <<< c, b, a <<< Loop-information End >>>	
52	2 p	2v	do i=1, isize
53	2 p	2v	a(i,j) = b(i,j) + scalar * c(i,j)
54	2 p	2v	enddo
55	1 p		enddo

外側ループ 1 回転先の配列に対して配列全体プリフェッチ



L1Dミスdm率が減少

→ プリフェッチ
キャッシュヒット
キャッシュミス
コード上アクセスしない領域

Cache	Load-store instruction	L1D miss	L1D miss rate (/Load-store instruction)	L1D miss demand rate (%) / (L1D miss)	L1D miss hardware prefetch rate (%) / (L1D miss)	L1D miss software prefetch rate (%) / (L1D miss)
改善前	1.75E+09	2.64E+08	0.15	70.94%	29.46%	-0.40%
改善後	9.22E+08	1.93E+08	0.21	23.46%	1.65%	74.90%

最内ループは、繰返し数が少なく、かつ繰返し数よりも配列サイズが大きいため、連續的ではありません。通常のプリフェッチではプリフェッチの立ち上がりのコストがみえてしまいます。そのため、浮動小数点ロードアクセス待ちが多くなっています。

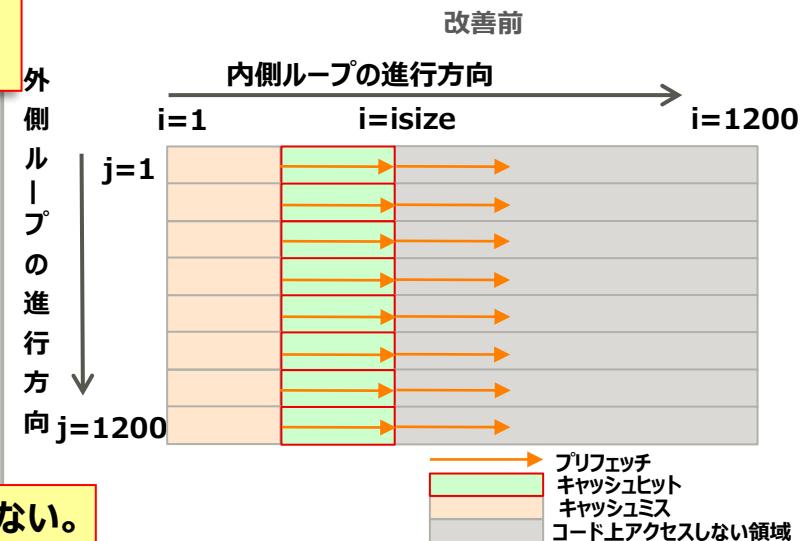
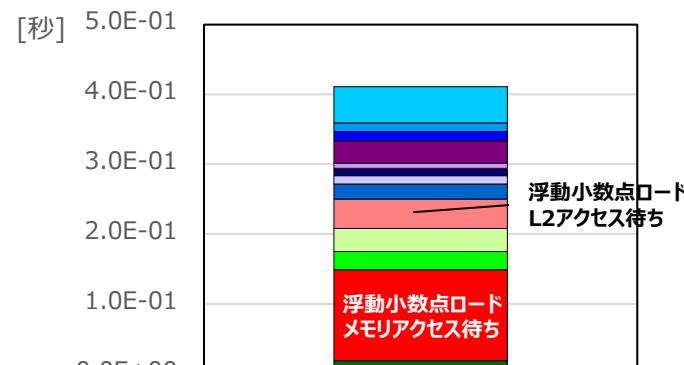
改善前ソース

```

40 void sub(double scalar, int isize){
41     int i, j;
42
43     #pragma omp parallel for
44     <<< Loop-information Start
45     <<< [OPTIMIZATION]
46     <<< PREFETCH(HARD) Expected by compiler :
47     <<< c, b, a
48     <<< Loop-information End >>>
49     p for (j = 0; j < n; j++){
<<< Loop-information Start >>>
45     p 2v for (i = 0; i < isize; i++){
46     p 2v     a[j][i] = b[j][i] + scalar * c[j][i];
47     p 2v }
48     p }
49 }
```

配列の2次元目の要素数 : 1200
 ループ繰返し数(isize) : 128
 外側ループjがインクリメントされた時に
 アクセスの連續性が途切れる

L1Dミスdm率が高くプリフェッチが効いていない。



Cache	L1 miss rate (/Effective instruction)	Load-store instruction	L1D miss	L1D miss rate (/Load-store instruction)	L1D miss demand rate (%) (/L1D miss)	L1D miss hardware prefetch rate (%) (/L1D miss)	L1D miss software prefetch rate (%) (/L1D miss)	L2 miss	L2 miss rate (/Load-store instruction)	L2 miss demand rate (%) (/L2 miss)	L2 miss hardware prefetch rate (%) (/L2 miss)	L2 miss software prefetch rate (%) (/L2 miss)
改善前	0.00	2.18E+09	2.57E+08	0.18	73.77%	26.27%	-0.04%	6.67E+07	0.03	42.40%	82.13%	0.00%

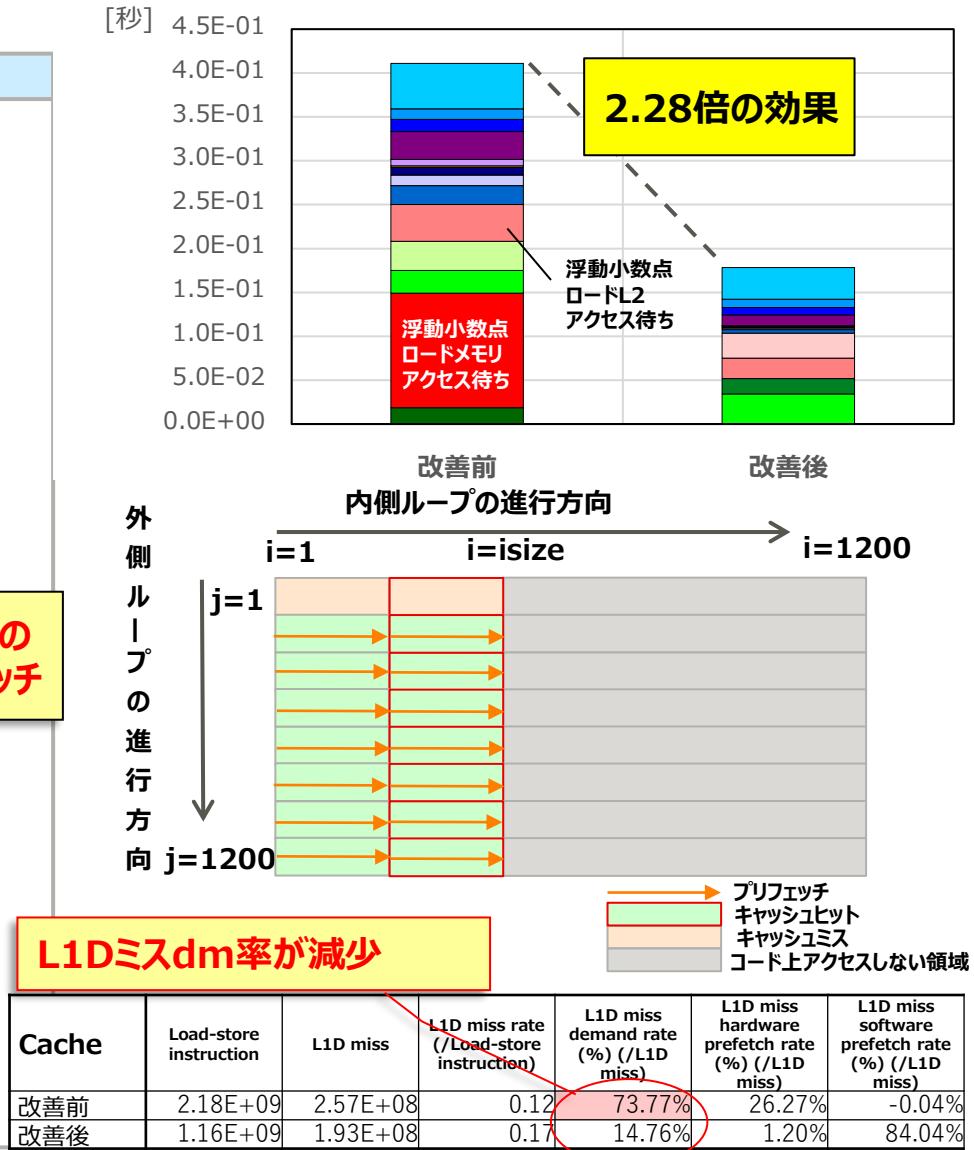
`__builtin_prefetch`関数を使用することで、外側ループの配列に対してプリフェッチを生成しプリフェッチの立ち上がりのコストを隠蔽します。その結果、浮動小数点ロードL2アクセス待ちが改善されました。

改善後ソース

```

41 void sub(double scalar, int isize){
42     int i, j;
43
44     #pragma omp parallel for
45     <<< Loop-information Start >>>
46     <<< [OPTIMIZATION]
47     <<< PREFETCH(HARD) Expected by compiler :
48     <<< c, b, a
49     <<< PREFETCH(SOFT) : 48
50     <<< SPECIFIED : 48
51     <<< a: 16, c: 16, b: 16
52     <<< Loop-information End >>>
53     p   for (j = 0; j < n; j++){
54     <<< Loop-information Start >>>
55     <<< [OPTIMIZATION]
56     <<< SIMD(VL: 8)
57     <<< SOFTWARE PIPELINING(IPC: 2.62, ITR: 56, MVE: 4, POL: S)
58     <<< PREFETCH(HARD) Expected by compiler :
59     <<< c, a, b
60     <<< PREFETCH(SOFT) : 48
61     <<< SPECIFIED : 48
62     <<< b: 16, c: 16, a: 16
63     <<< Loop-information End >>>
64     p   for(i = 0; i < isize; i++){
65     <<< __builtin_prefetch(&a[j+1][0], 1, 3);
66     <<< __builtin_prefetch(&b[j+1][0], 0, 3);
67     <<< __builtin_prefetch(&c[j+1][0], 0, 3);
68     <<< __builtin_prefetch(&a[j+1][32], 1, 3);
69     <<< __builtin_prefetch(&b[j+1][32], 0, 3);
70     <<< __builtin_prefetch(&c[j+1][32], 0, 3);
71     <<< __builtin_prefetch(&a[j+1][64], 1, 3);
72     <<< __builtin_prefetch(&b[j+1][64], 0, 3);
73     <<< __builtin_prefetch(&c[j+1][64], 0, 3);
74     <<< __builtin_prefetch(&a[j+1][96], 1, 3);
75     <<< __builtin_prefetch(&b[j+1][96], 0, 3);
76     <<< __builtin_prefetch(&c[j+1][96], 0, 3);
77     <<< a[j][i] = b[j][i] + scalar * c[j][i];
78     }
79   }
80 }
```

外側ループ1回転先の配列に対してプリフェッチ



外側のループで配列に対する__builtin_prefetch関数を指定することで、最内ループでの命令数を減らすことができ、さらに性能を向上させることができます。 [秒]

改善後ソース

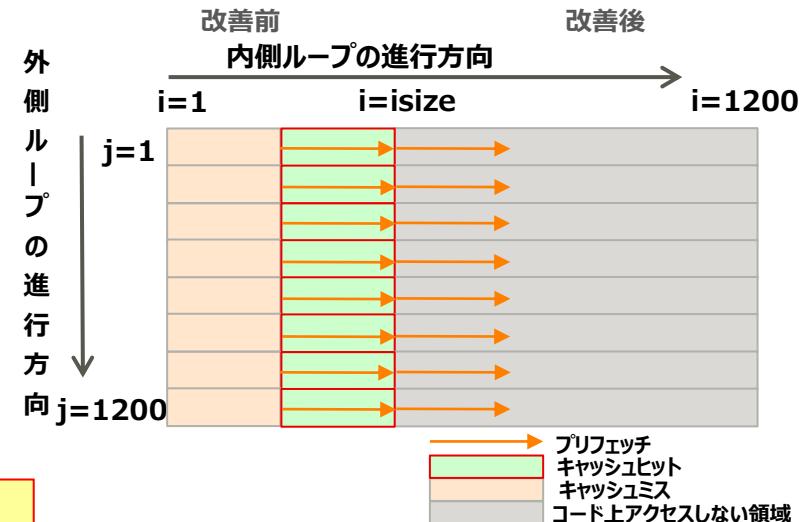
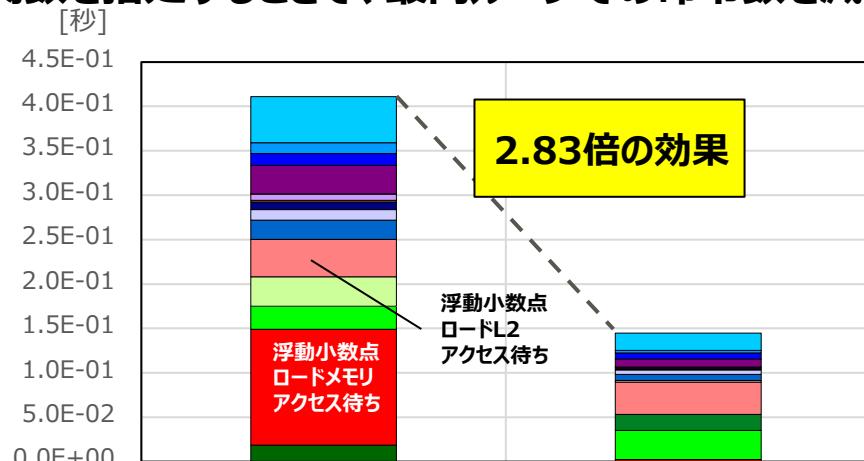
```

41 void sub(double scalar, int isize){
42     int i, j;
43
44     #pragma omp parallel for
45     <<< Loop-information Start
46     <<< [OPTIMIZATION]
47     <<< PREFETCH(HARD)
48     <<< c, b, a
49     <<< PREFETCH(SOFT) : 外側ループ 1 回
50     <<< SPECIFIED : 12
51     <<< a: 4, b: 4, c: 4
52     <<< Loop-information End >>>
53
54     p for (j = 0; j < n; j++){
55         __builtin_prefetch(&a[j+1][0], 1, 3);
56         __builtin_prefetch(&a[j+1][32], 1, 3);
57         __builtin_prefetch(&a[j+1][64], 1, 3);
58         __builtin_prefetch(&a[j+1][96], 1, 3);
59         __builtin_prefetch(&b[j+1][0], 0, 3);
60         __builtin_prefetch(&b[j+1][32], 0, 3);
61         __builtin_prefetch(&b[j+1][64], 0, 3);
62         __builtin_prefetch(&b[j+1][96], 0, 3);
63         __builtin_prefetch(&c[j+1][0], 0, 3);
64         __builtin_prefetch(&c[j+1][32], 0, 3);
65         __builtin_prefetch(&c[j+1][64], 0, 3);
66         __builtin_prefetch(&c[j+1][96], 0, 3);
67     }
68
69     <<< Loop-information Start >>>
70     <<< [OPTIMIZATION]
71     <<< SIMD(VL: 8)
72     <<< SOFTWARE PIPELINING(IPC: 3.25, ITR: 14)
73     <<< PREFETCH(HARD) Expected by compiler :
74     <<< c, b, a
75     <<< Loop-information End >>>
76
77     p 2v for (i = 0; i < isize; i++){
78         2v     a[j][i] = b[j][i] + scalar * c[j][i];
79         2v     }
80     }
81
82     }

```

外側ループ¹回転先の配列に対して配列全体プリフェッチ

L1Dミスdm率が減少



Cache	Load-store instruction	L1D miss	L1D miss rate (/Load-store instruction)	L1D miss demand rate (%) (/L1D miss)	L1D miss hardware prefetch rate (%) (/L1D miss)	L1D miss software prefetch rate (%) (/L1D miss)
改善前	2.18E+09	2.57E+08	0.12	73.77%	26.27%	-0.04%
改善後	1.17E+09	1.94E+08	0.17	23.67%	1.43%	74.90%

データアクセス待ちの改善 (アクセス量の軽減)

- 高速ストア(ZFILL)

高速ストア（ZFILL）

- 高速ストア（ZFILL）とは
- ZFILL（改善前）
- ZFILLの効果（最適化制御行チューニング）

高速ストア (ZFILL) とは

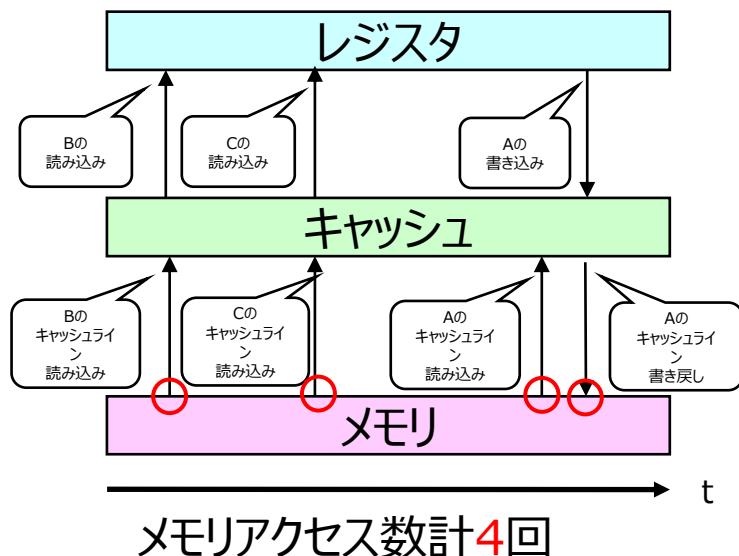
■ 高速ストア (ZFILL) とは

キヤッショ上に書き込み用のキヤッショライン（中身は不定値）を確保する機能です。これにより、メモリからのキヤッショラインの読み込みが削減できるので、メモリループトがネックとなっているプログラムでは性能が改善します。

■ 動作条件

- ストア対象の配列がイタレーション間で依存がないこと
- 定義のある配列の参照が無いこと
- 連続的なメモリアクセスであること

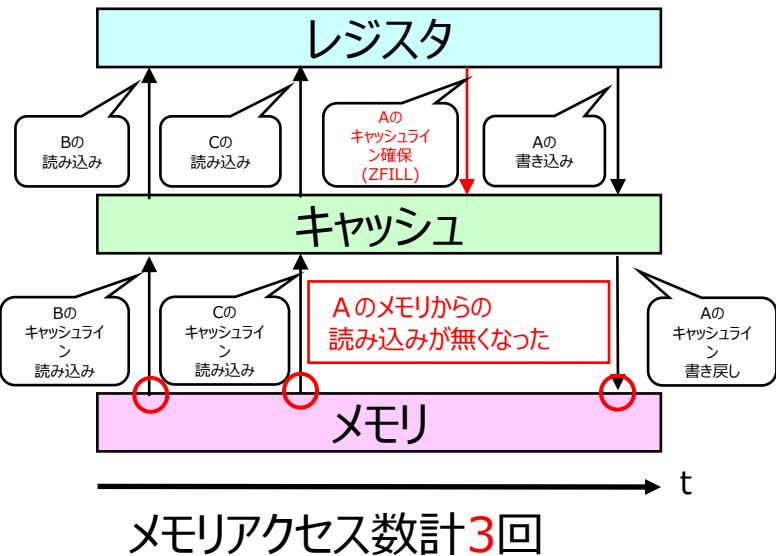
ZFILL未使用時



例)

```
DO I = 1, N
  A(I) = B(I) + C(I)
END DO
```

ZFILL使用時



ZFILLとは(1/2)

以下の最適化制御行を指定します。

最適化指示子 (Fortran)	意味	指定可能な最適化制御行			
		プログラム単位	DOループ単位	文単位	配列代入文単位
ZFILL[(m1)]	ZFILL命令を生成することを指示します。m1はキャッシュライン数を表す1～100の10進数です。	×	○	×	○
NOZFILL	ZFILL命令を生成しないことを指示します。	×	○	×	○

最適化指示子 (C/C++)	意味	指定可能な最適化制御行			
		global行	procedure行	loop行	statement行
zfill[(m1)]	zfill命令を生成することを指示します。m1はキャッシュライン数を表す1～100の10進数です。	×	×	○	×
nozfill	zfill命令を生成しないことを指示します。	×	×	○	×

■ 注意事項

- ZFILL命令は、ループ内でストアされる配列データに対して出力されます。ただし、同一ループ内に参照がある配列、非連続アクセスされる配列またはIF構文配下でストアされる配列に対しては出力されません。
- ZFILL命令が出力された場合、2次キャッシュへのプリフェッチ命令は出力されません。
- ZFILL命令で確保したキャッシュラインは必ずストアするようなループ変形を行うため、以下の最適化が適用できなくなります。これにより、実行性能が低下することがあります。
 - ループアンローリング
 - ループストライピング
- また、以下の場合にも実行性能が低下することがあります。
 - 繰返し数が少ないループ
 - データが1次キャッシュまたは2次キャッシュにある場合

以下の翻訳オプションを指定することで、最適化制御行チューニングと同等の効果を得ることができます。

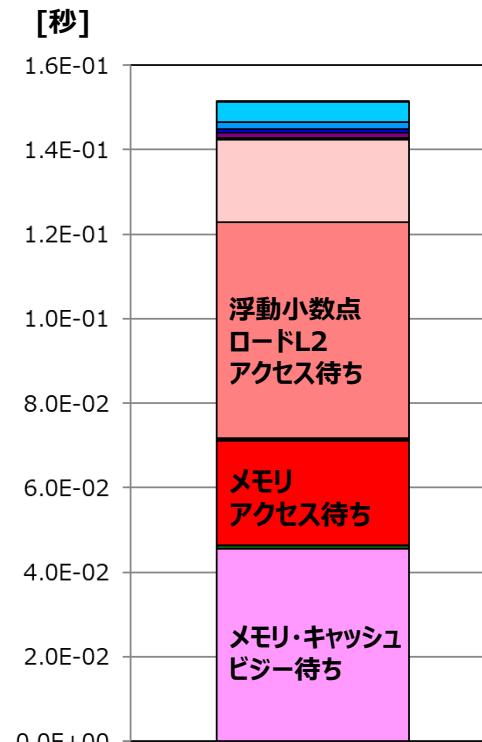
翻訳オプション	機能説明
<pre>-K{ zfill[=N] nozfill } 1 ≤ N ≤ 100</pre>	<p>ループ内で書き込みのみ行う配列データについて、データをメモリからロードすることなく、キャッシュ上に書き込み用のキャッシュラインを確保する命令（ZFILL命令）を生成することを指示します。Nを指定することで、Nキャッシュライン先のデータをZFILL命令の対象とします。</p> <p>Nは1～100の範囲で指定できます。Nの指定を省略した場合、コンパイラが自動的に値を決定します。</p> <p>本オプションは、-O2オプション以上が有効な場合に意味があります。デフォルトは、-Knозfillです。</p>

■使用例（改善前ソース時）

```
$ frtpx -Kfast,parallel sample.f90 -Kzfill
$ fccpx -Kfast,parallel sample.f90 -Kzfill
```

メモリアクセス負荷が高いプログラムのため、メモリスループットがネックになっています。そのため、データアクセス待ちが多くなっています。

改善前ソース	
38	real*8 a(n),b(n),c(n),d
39	<<< Loop-information Start >>>
	<<< [PARALLELIZATION]
	<<< Standard iteration count: 1000
	<<< [OPTIMIZATION]
	<<< SIMD(VL: 8)
	<<< SOFTWARE PIPELINING(IPC: 3.25, ITR: 144, MVE: 4, POL: S)
	<<< PREFETCH(HARD) Expected by compiler :
	<<< c, b, a
	<<< Loop-information End >>>
40	1 pp 2v do i=1,n
41	1 p 2v a(i) = b(i) + c(i)*d
42	1 p 2v enddo



改善前

Cache	L1 miss rate (/Effective instruction)	Load-store instruction	L1D miss	L1D miss rate (/Load-store instruction)	L1D miss demand rate (%) (/L1D miss)	L1D miss hardware prefetch rate (%) (/L1D miss)	L1D miss software prefetch rate (%) (/L1D miss)	L2 miss	L2 miss rate (/Load-store instruction)	L2 miss demand rate (%) (/L2 miss)	L2 miss hardware prefetch rate (%) (/L2 miss)	L2 miss software prefetch rate (%) (/L2 miss)
改善前	0.00	3.90E+08	9.39E+07	0.24	27.78%	72.21%	0.01%	9.38E+07	0.24	12.72%	88.66%	0.00%

	Memory Throughput (GB/s)
改善前	211.59

メモリスループットがネックになっている

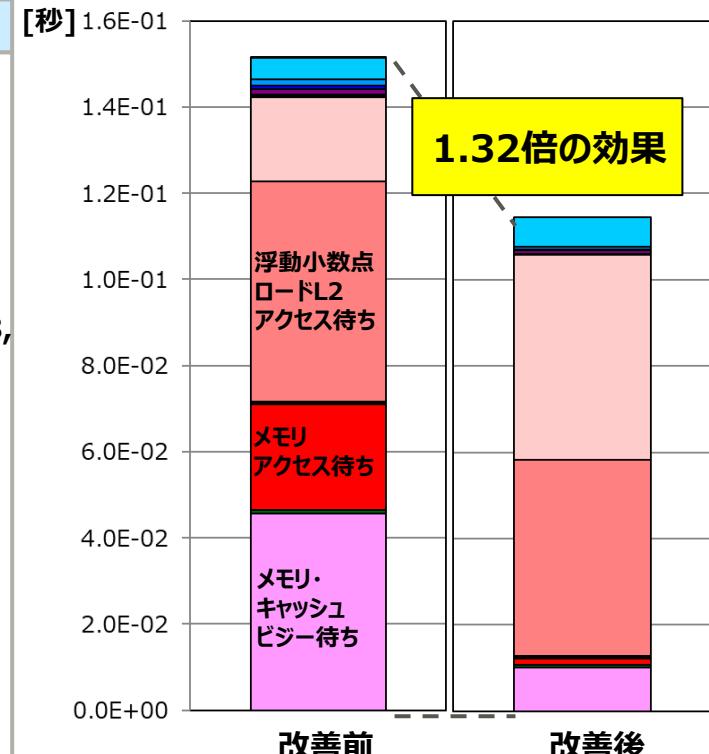
ZFILL指示子を指定することでストア命令によるメモリからのキャッシュラインの読み込みがなくなり、L2ミス数が削減されました。その結果、データアクセス待ちが改善されました。

改善後ソース（最適化制御行チューニング）

```

38      real*8 a(n),b(n),c(n),d
39      !ocl zfill
<<< Loop-information Start >>>
<<< [PARALLELIZATION]
<<< Standard iteration count: 1000
<<< [OPTIMIZATION]
<<< SIMD(VL: 8)
<<< SOFTWARE PIPELINING(IPC: 2.55, ITR: 128,
                           MVE: 2, POL: S)
<<< PREFETCH(HARD) Expected by compiler :
<<<     c, b
<<< PREFETCH(SOFT) : 2
<<< SEQUENTIAL : 2
<<<     a: 2
<<< ZFILL      :
<<<     a
<<< Loop-information End >>>
40      1 pp v    do i=1,n
41      1 p  v    a(i) = b(i) + c(i)*d
42      1 p  v    enddo

```



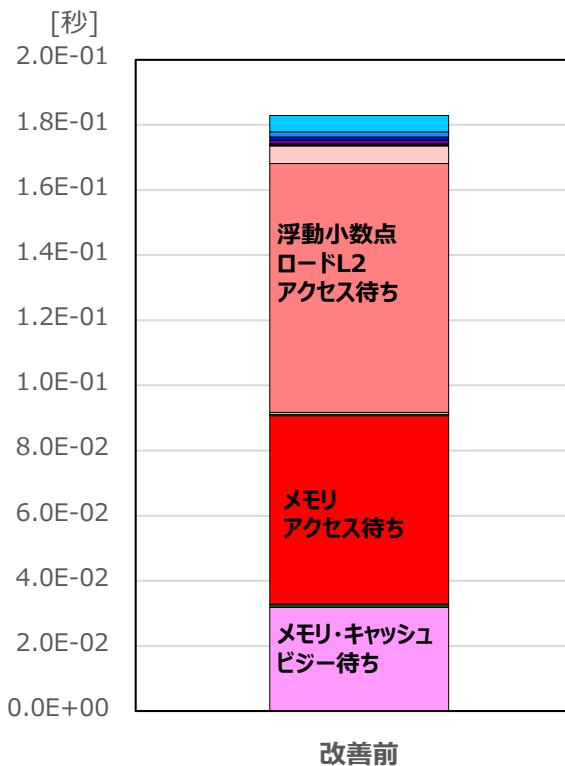
Cache	L1I miss rate (/Effective instruction)	Load-store instruction	L1D miss	L1D miss rate (/Load-store instruction)	L1D miss demand rate (%) (/L1D miss)	L1D miss hardware prefetch rate (%) (/L1D miss)	L1D miss software prefetch rate (%) (/L1D miss)	L2 miss	L2 miss rate (/Load-store instruction)	L2 miss demand rate (%) (/L2 miss)	L2 miss hardware prefetch rate (%) (/L2 miss)	L2 miss software prefetch rate (%) (/L2 miss)
改善前	0.00	3.90E+08	9.39E+07	0.24	27.78%	72.21%	0.01%	9.38E+07	0.24	12.72%	88.66%	0.00%
改善後	0.00	4.38E+08	9.39E+07	0.21	16.80%	49.98%	33.22%	6.25E+07	0.14	1.15%	98.98%	0.00%

	Memory Throughput (GB/s)
改善前	211.59
改善後	209.96

改善後もメモリループトネットではあるが、L2
ミス数が1/3削減された

メモリアクセス負荷が高いプログラムのため、データアクセス待ちが多くなっています。

改善前ソース	
38	void sub(double * restrict a, double * restrict b, double * restrict c, double d, int n){
39	int i;
40	
41	#pragma omp parallel for <<< Loop-information Start >>> <<< [OPTIMIZATION] <<< SIMD(VL: 8) <<< SOFTWARE PIPELINING(IPC: 3.25, ITR: 144, MVE: 4, POL: S) <<< PREFETCH(HARD) Expected by compiler : <<< (unknown) <<< Loop-information End >>>
42	p 2v for (i = 0; i < n; i++){
43	p 2v a[i] = b[i] + c[i]*d;
44	p 2v }
45	}



Cache	L1I miss rate (/Effective instruction)	Load-store instruction	L1D miss	L1D miss rate (/Load-store instruction)	L1D miss demand rate (%) (/L1D miss)	L1D miss hardware prefetch rate (%) (/L1D miss)	L1D miss software prefetch rate (%) (/L1D miss)	L2 miss	L2 miss rate (/Load-store instruction)	L2 miss demand rate (%) (/L2 miss)	L2 miss hardware prefetch rate (%) (/L2 miss)	L2 miss software prefetch rate (%) (/L2 miss)
改善前	0.00	3.95E+08	9.40E+07	0.24	41.24%	58.75%	0.01%	9.39E+07	0.24	19.15%	83.83%	0.00%

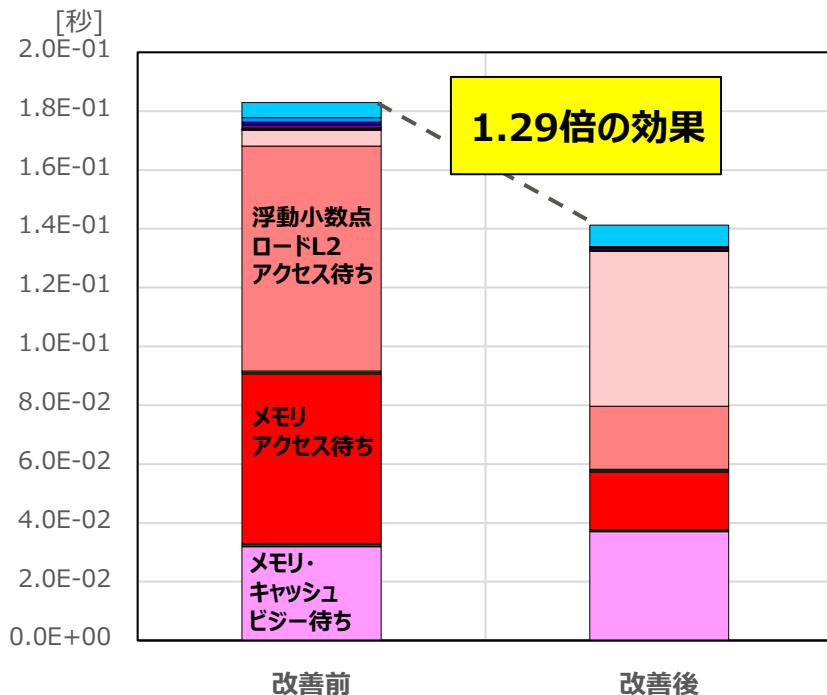
Statistics	Memory throughput (GB/s)
改善前	175.68

ZFILL指示子を指定することでストア命令によるメモリからのキャッシュラインの読み込みがなくなり、L2ミス数が削減されました。その結果、データアクセス待ちが改善されました。

改善後ソース（最適化制御行チューニング）

```

37     void sub(double * restrict a, double * restrict b,
38             double * restrict c, double d, int n){
39             int i;
40
41             #pragma omp parallel for
42             #pragma loop zfill
43             <<< [Loop-information Start >>>
44             <<< [OPTIMIZATION]
45             <<< SIMD(VL: 8)
46             <<< SOFTWARE PIPELINING(IPC: 2.55, ITR: 128,
47                             MVE: 2, POL: S)
48             <<< PREFETCH(HARD) Expected by compiler :
49             <<< (unknown)
50             <<< PREFETCH(SOFT) : 2
51             <<< SEQUENTIAL : 2
52             <<< (unknown): 2
53             <<< ZFILL :
54             <<< (unknown)
55             <<< Loop-information End >>>
56             p   v   for (i = 0; i < n; i++) {
57             p   v   a[i] = b[i] + c[i]*d;
58             p   v   }
59         }
```



1.29倍の効果

Cache	L1 miss rate (/Effective instruction)	Load-store instruction	L1D miss	L1D miss rate (/Load-store instruction)	L1D miss demand rate (%) (/L1D miss)	L1D miss hardware prefetch rate (%) (/L1D miss)	L1D miss software prefetch rate (%) (/L1D miss)	L2 miss	L2 miss rate (/Load-store instruction)	L2 miss demand rate (%) (/L2 miss)	L2 miss hardware prefetch rate (%) (/L2 miss)	L2 miss software prefetch rate (%) (/L2 miss)
改善前	0.00	3.95E+08	9.40E+07	0.24	41.24%	58.75%	0.01%	9.39E+07	0.24	19.15%	83.83%	0.00%
改善後	0.00	4.31E+08	9.40E+07	0.22	30.94%	35.79%	33.27%	6.26E+07	0.15	4.59%	95.79%	0.00%

Statistics	Memory throughput (GB/s)
改善前	175.68
改善後	170.27

L2ミス数が1/3削減された

データアクセス待ち（スラッシングの改善）

- ・ キャッシュスラッシングとは
- ・ 1次元目の配列要素数を増やすパディング
- ・ 2次元目の配列要素数を増やすパディング
- ・ ダミー配列によるパディング
- ・ ダミー配列によるパディング（サイズが異なる配列）
- ・ 配列マージ（スラッシングの改善）
- ・ ループ分割
- ・ ラージページ環境変数によるパディング

キャッシュスラッシングとは

キャッシュスラッシングとは、**キャッシュ上の特定のインデックス（キャッシュ上の位置情報）のデータだけが頻繁に上書きされる現象のことです。**

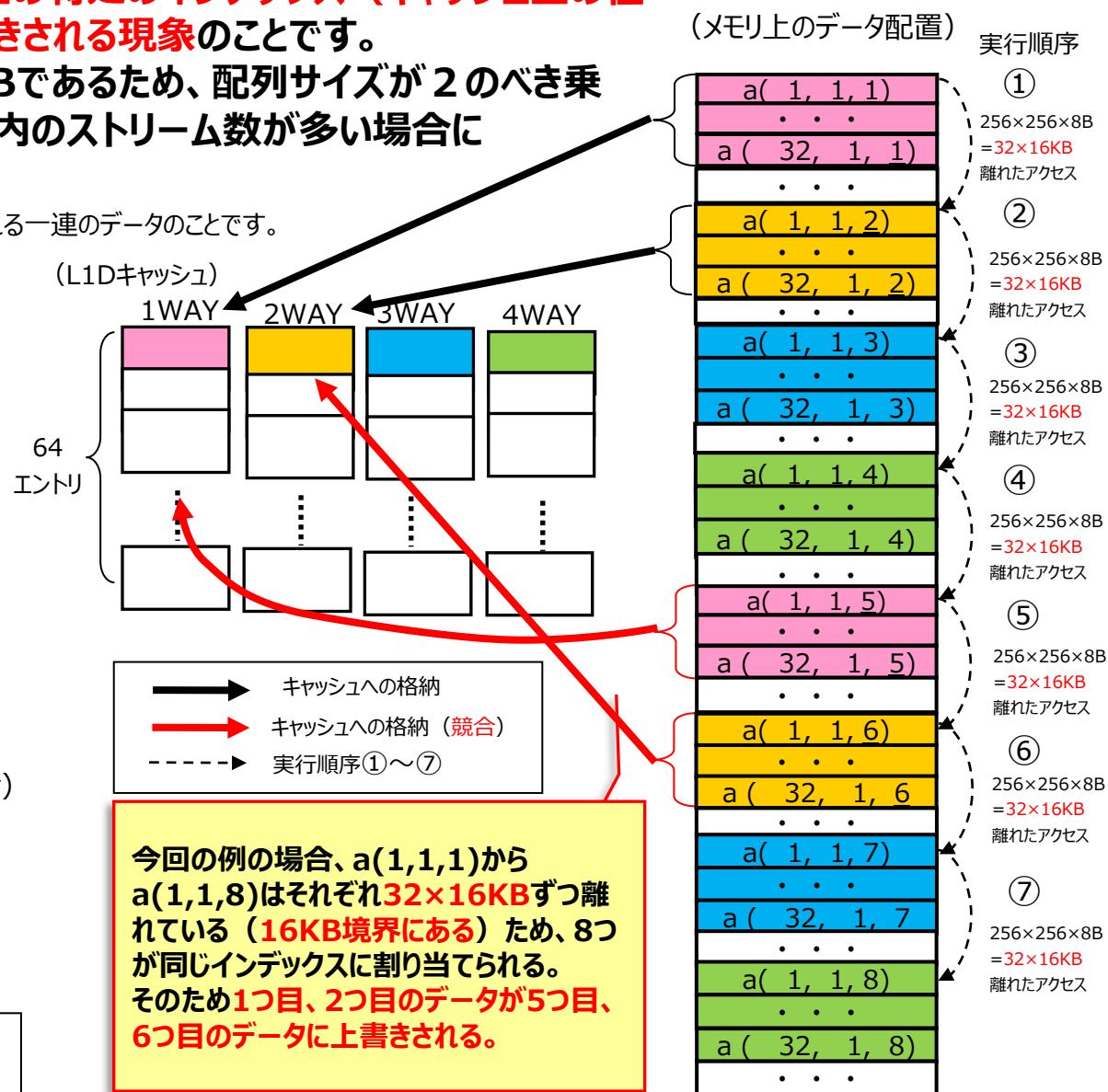
この現象は、1WAYのサイズが16KBであるため、配列サイズが2のべき乗（16KBの倍数）の場合およびループ内のストリーム数が多い場合に発生しやすくなります。

注：ストリームとは、ループの回転に伴って参照定義される一連のデータのことです。

ソース例
<pre>subroutine sub(a, n, m) !n=256, m=256 real*8 a(n,m,8) do j = 1 , m do i = 1 , n a(i,j,8)=a(i,j,1)+a(i,j,2)+a(i,j,3)+a(i,j,4)+a(i,j,5)+a(i,j,6)+a(i,j,7) enddo enddo end</pre>

■ L1Dキャッシュスラッシングの目安
(512bitSVEアクセス命令、連続アクセスの場合)

L1D miss rate (/Load-store instruction)
0.25以上
単精度:64/256 倍精度:64/256



パディング

- ・ パディングとは
- ・ 1次元目の配列要素数を増やすパディング
- ・ 2次元目の配列要素数を増やすパディング
- ・ ダミー配列によるパディング
- ・ ダミー配列によるパディング（サイズが異なる配列）

パディングとは

パディングとは、配列間や配列の中にダミーの領域を挿入することです。

■ 使用条件

同一配列で複数のストリームが存在している。
もしくは
複数の配列が存在している。

■ 狹い

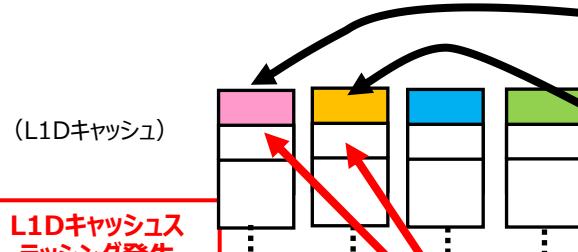
テンポラリ領域を作りアドレスをずらす。

■ 副作用

問題規模変更ごとにパディング量の変更が必要

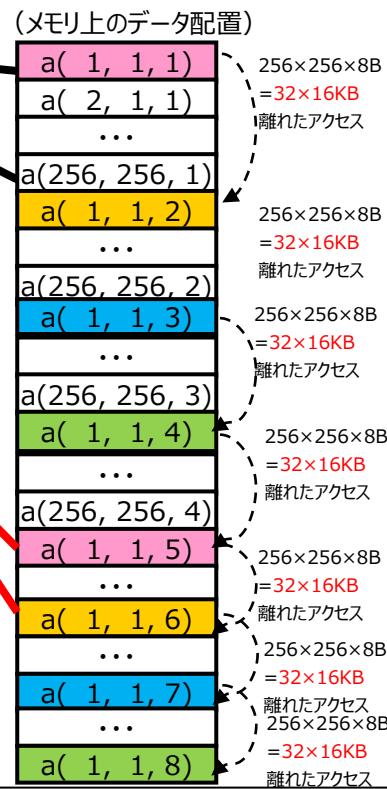
同一配列で複数のストリームが存在している例

■ 改善前

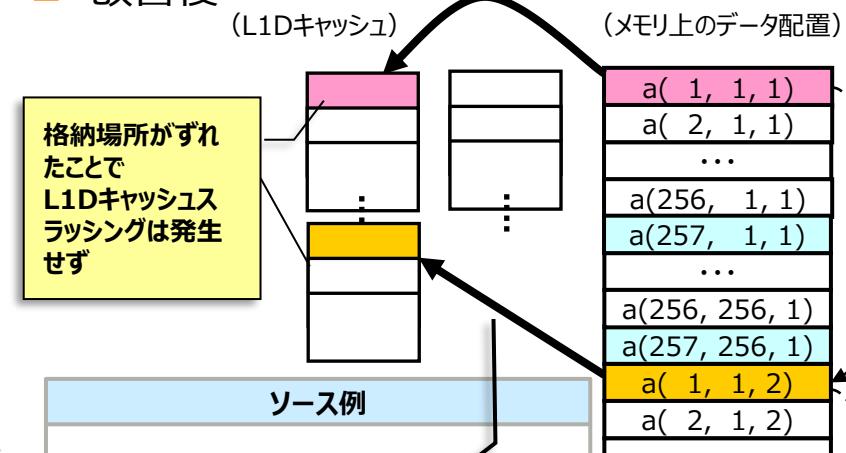


ソース例

```
parameter(n=256,m=256)
real*8 a(n, m, 8)
common /com/a
do j = 1 , m
  do i = 1 , n
    a(i, j, 8) = a(i, j, 1) + a(i, j, 2) + a(i, j, 3) +
      a(i, j, 4) + a(i, j, 5) + a(i, j, 6) +
      a(i, j, 7)
  enddo
enddo
```



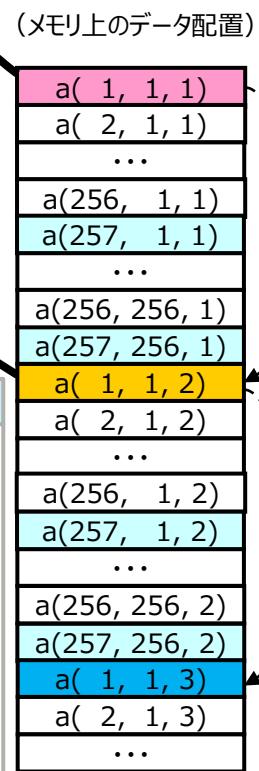
■ 改善後



ソース例

```
parameter(n=257,m=256)
real*8 a(n, m, 8)
common /com/a
do j = 1 , m
  do i = 1 , n
    a(i, j, 8) = a(i, j, 1) + a(i, j, 2) + a(i, j, 3) +
      a(i, j, 4) + a(i, j, 5) + a(i, j, 6) +
      a(i, j, 7)
  enddo
enddo
```

パディングを行うことで キャッシュに格納される場所を離れた位置にする



→ キャッシュへの格納

→ キャッシュへの格納 (競合)

→ メモリへのアクセス順番

1次元目の配列要素数を増やすパディング

- 1次元目の配列要素数を増やすパディング（改善前）
- 1次元目の配列要素数を増やすパディングの効果（ソースチューニング）
- 1次元目の配列要素数を増やすパディングの効果（翻訳オプションチューニング）

配列aのそれぞれのストリーム同士が16KB境界にあるためL1Dキャッシュスラッシングが発生します。そのため、浮動小数点ロードL2アクセス待ちが多くなっています。

改善前ソース

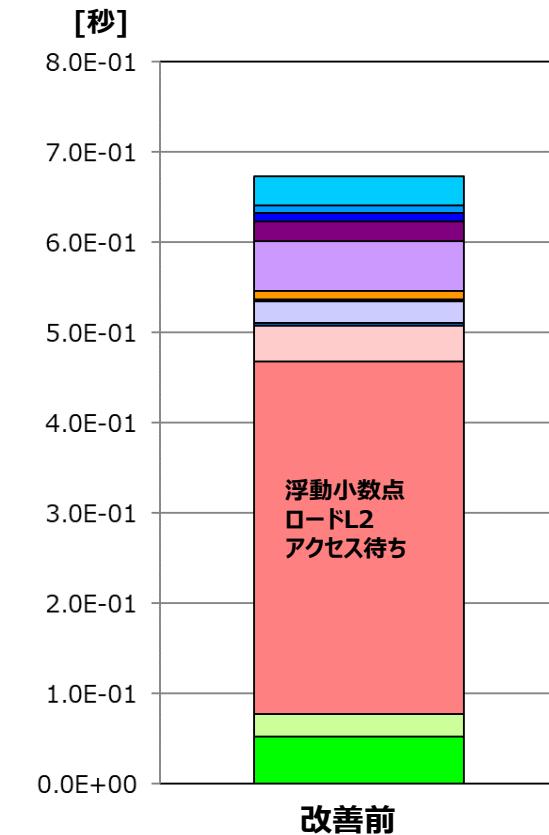
```

42      parameter(n=256,m=256)
43      real*8 a(n, m, 8)
44      common /com/a
        <<< Loop-information Start
        <<< [PARALLELIZATION]
        <<< Standard iteration co
        <<< [OPTIMIZATION]
        <<< COLLAPSED
        <<< SIMD(VL: 8)
        <<< SOFTWARE PIPELINING(IPC: 2.66, ITR: 104,
                                MVE: 7, POL: S)
        <<< PREFETCH(HARD) Expected by compiler :
        <<< a
        <<< Loop-information End >>>
45  1 pp v      do j = 1 , m
        <<< Loop-information Start >>>
        <<< [OPTIMIZATION]
        <<< COLLAPSED
        <<< Loop-information End >>>
46  2 p          do i = 1 , n
47  2 p  v      a(i, j, 8) = a(i, j, 1) + a(i, j, 2) + a(i, j, 3) + &
48  2           a(i, j, 4) + a(i, j, 5) + a(i, j, 6) + a(i, j, 7)
49  2 p  v      enddo
50  1 p          enddo

```

配列サイズ
256×256×8B=
32×16KB
(16KB境界)

同一配列のストリーム



配列が連続アクセスであるにもかかわらずL1Dミスが高い
⇒ L1Dキャッシュスラッシングが発生している

Cache	L1I miss rate (/Effective instruction)	Load-store instruction	L1D miss	L1D miss rate (/Load-store instruction)	L1D miss demand rate (%) / L1D miss	L1D miss hardware prefetch rate (%) / L1D miss	L1D miss software prefetch rate (%) / L1D miss	L2 miss	L2 miss rate (/Load-store instruction)	L2 miss demand rate (%) / L2 miss	L2 miss hardware prefetch rate (%) / L2 miss	L2 miss software prefetch rate (%) / L2 miss
改善前	0.00	3.08E+09	1.30E+09	0.42	67.73%	32.27%	0.00%	1.50E+04	0.00	37.70%	72.22%	0.00%

配列aのそれぞれのストリームの1次元目にパディング(+1)することで、L1Dキャッシュラッシングを回避しました。その結果、浮動小数点ロードL2アクセス待ちが改善されました。

改善後ソース

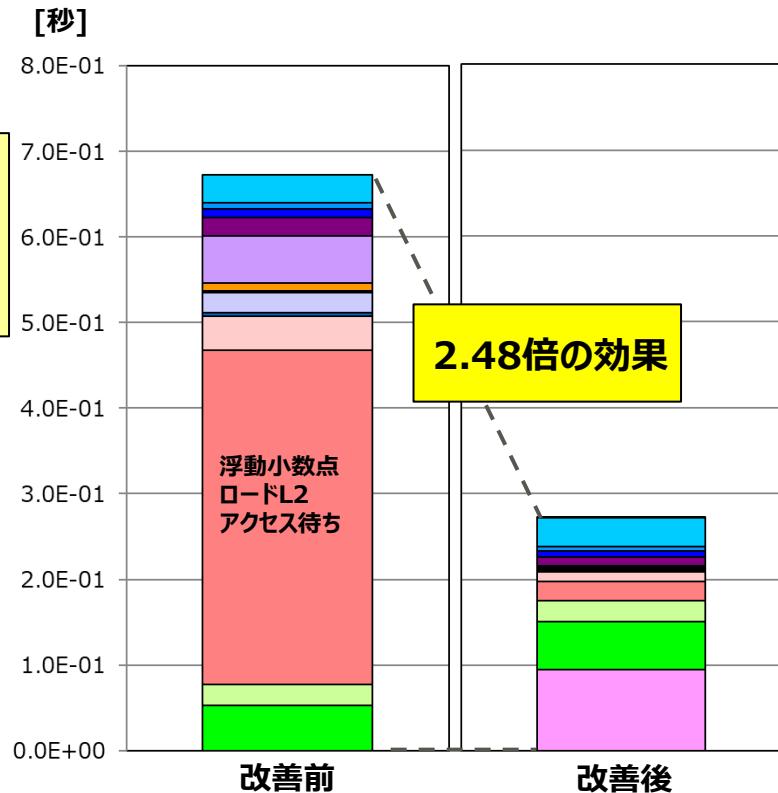
```

42      parameter(n=257,m=256)
43      real*8 a(n, m, 8)
44      common /com/a
45      <<< Loop-information Start
46      <<< [PARALLELIZATION]
47      <<< Standard iteration cou
48      <<< [OPTIMIZATION]
49      <<< COLLAPSED
50      <<< SIMD(VL: 8)
51      <<< SOFTWARE PIPELINING(IPC: 2.66, ITR: 104,
52                                MVE: 7, POL: S)
53      <<< PREFETCH(HARD) Expected by compiler :
54      <<< a
55      <<< Loop-information End >>>
56      1 pp v    do j = 1 , m
57      <<< Loop-information Start >>>
58      <<< [OPTIMIZATION]
59      <<< COLLAPSED
60      <<< Loop-information End >>>
61      2 p      do i = 1 , n
62      2 p v      a(i, j, 8) = a(i, j, 1) + a(i, j, 2) + a(i, j, 3) + &
63                                a(i, j, 4) + a(i, j, 5) + a(i, j, 6) + a(i, j, 7)
64      2 p v      enddo
65      1 p      enddo

```

nを+1することで
16KB境界からずらす

L1Dミスが減少した



■ 注意事項

パディング数が大きすぎると、データの連續性が損なわれ
ハードウェア prefetch が効かなくなることがあります。

Cache	L1I miss rate (/Effective instruction)	Load-store instruction	L1D miss	L1D miss rate (/Load-store instruction)	L1D miss demand rate (%) (/L1D miss)	L1D miss hardware prefetch rate (%) (/L1D miss)	L1D miss software prefetch rate (%) (/L1D miss)	L2 miss	L2 miss rate (/Load-store instruction)	L2 miss demand rate (%) (/L2 miss)	L2 miss hardware prefetch rate (%) (/L2 miss)	L2 miss software prefetch rate (%) (/L2 miss)
改善前	0.00	3.08E+09	1.30E+09	0.42	67.73%	32.27%	0.00%	1.50E+04	0.00	37.70%	72.22%	0.00%
改善後	0.00	2.62E+09	4.58E+08	0.17	8.20%	91.80%	0.00%	9.69E+03	0.00	46.19%	58.89%	0.00%

配列aのそれぞれのストリーム同士が16KB境界にあるためL1Dキャッシュラッシングが発生します。そのため、浮動小数点ロードL2アクセス待ちが多くなっています。

改善前ソース

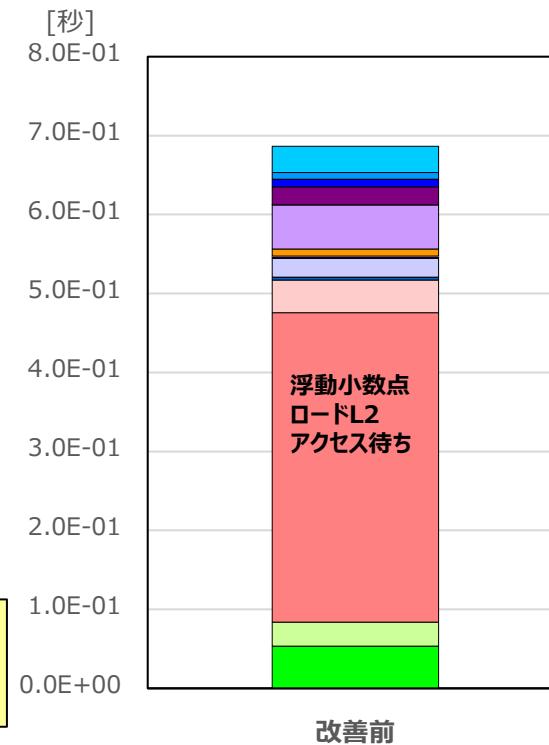
```

30     void sub(void){
31         int i, j;
32
33         #pragma omp parallel for collapse(2)
34         <<< Loop-information Start >>>
35         <<< [OPTIMIZATION]
36         <<< SIMD(VL: 8)
37         <<< SOFTWARE PIPELINING(IPC: 2.66, ITR: 104, MVE: 7, POL: S)
38         <<< PREFETCH(HARD) Expected by compiler :
39         <<< a
40         <<< Loop-information End >>>
41         p   v   for (j = 0; j < m; j++){
42             p   v   for (i = 0; i < n; i++){
43                 p   v       a[7][j][i] = a[0][j][i] + a[1][j][i] + a[2][j][i] + a[3][j][i] +
44                               a[4][j][i] + a[5][j][i] + a[6][j][i];
45             p   v   }
46             p   v   }
47         }

```

配列宣言
double a[8][256][256];
aの配列サイズ
256×256×8B=
32×16KB(16KB境界)

同一配列のストリーム



配列が連続アクセスであるにもかかわらずL1Dミスが高い
⇒ L1Dキャッシュラッシングが発生している

Cache	L1I miss rate (/Effective instruction)	Load-store instruction	L1D miss	L1D miss rate (/Load-store instruction)	L1D miss demand rate (%) (/L1D miss)	L1D miss hardware prefetch rate (%) (/L1D miss)	L1D miss software prefetch rate (%) (/L1D miss)	L2 miss	L2 miss rate (/Load-store instruction)	L2 miss demand rate (%) (/L2 miss)	L2 miss hardware prefetch rate (%) (/L2 miss)	L2 miss software prefetch rate (%) (/L2 miss)
改善前	0.00	3.08E+09	1.33E+09	0.43	68.32%	31.68%	0.00%	2.65E+04	0.00	49.48%	60.25%	0.00%

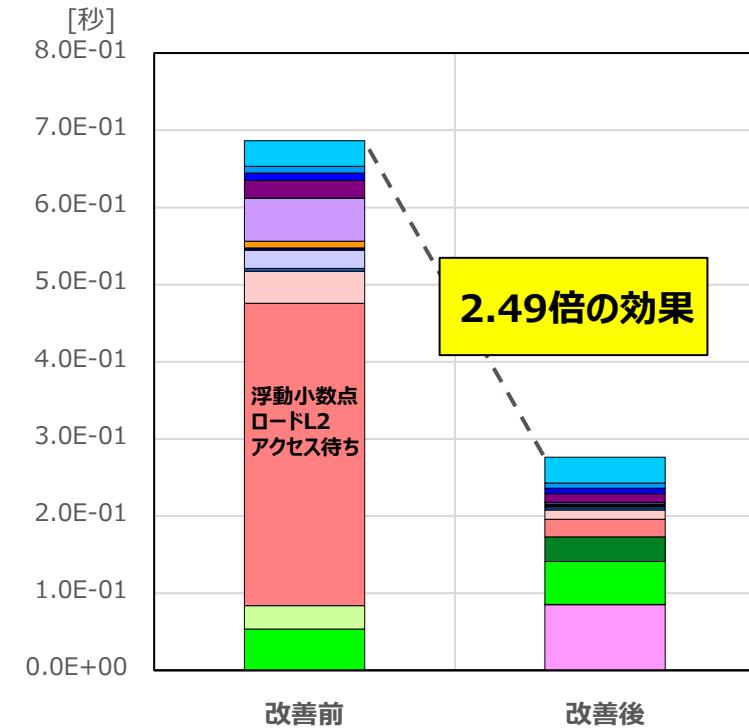
配列aのそれぞれのストリームの最終次元にパディング(+1)することで、L1Dキャッシュスラッシングを回避しました。その結果、浮動小数点ロードL2アクセス待ちが改善されました。

改善後ソース

```

30 void sub(void){
31     int i, j;
32
33     #pragma omp parallel for collapse(3)
34     <<< Loop-information Start >>>
35     <<< [OPTIMIZATION]
36     <<< SIMD(VL: 8)
37     <<< SOFTWARE PIPELINING(IPC: 2.66, ITR: 104, MVE: 7, POL: S)
38     <<< PREFETCH(HARD) Expected by compiler :
39     <<< a
40     <<< Loop-information End >>>
41     p    v    for (j = 0; j < m; j++){
42     p    v    for (i = 0; i < n; i++){
43     p    v        a[7][j][i] = a[0][j][i] + a[1][j][i] + a[2][j][i] +
44                           a[3][j][i] + a[4][j][i] + a[5][j][i] + a[6][j][i];
45     p    v    }
46     p    v    }
47 }
```

配列宣言
double a[8][256][257];
配列aの最終次元を+1することで16KB境界からずらす



L1Dミスが減少した

■ 注意事項
パディング数が大きすぎると、データの連續性が損なわれ
ハードウェア prefetch が効かなくなることがあります。

Cache	L1I miss rate (/Effective instruction)	Load-store instruction	L1D miss	L1D miss rate (/Load-store instruction)	L1D miss demand rate (%) (/L1D miss)	L1D miss hardware prefetch rate (%) (/L1D miss)	L1D miss software prefetch rate (%) (/L1D miss)	L2 miss	L2 miss rate (/Load-store instruction)	L2 miss demand rate (%) (/L2 miss)	L2 miss hardware prefetch rate (%) (/L2 miss)	L2 miss software prefetch rate (%) (/L2 miss)
改善前	0.00	3.08E+09	1.33E+09	0.43	68.32%	31.68%	0.00%	2.65E+04	0.00	49.48%	60.25%	0.00%
改善後	0.00	2.67E+09	4.60E+08	0.17	8.50%	91.50%	0.00%	2.26E+04	0.00	22.14%	83.91%	0.00%

以下の翻訳オプション(Fortran固有)を指定することで、ソースチューニングと同等の効果を得ることができます。

翻訳オプション	機能説明
-Karraypad_const[=N] ($1 \leq N \leq 2,147,483,647$)	1次元目が明示上下限であり、かつその上下限式が定数式の配列に対して N 要素のパディングを行います。 N が省略された場合、対象の配列ごとにコンパイラがパディングの量を決めます。パディングとは、配列の内部に隙間を作ることです。
-Karraypad_expr=N ($1 \leq N \leq 2,147,483,647$)	上下限式が定数式かどうかにかかわらず、1次元目が明示上下限の配列に対して N 要素のパディングを行います。

■ 使用例 (改善前ソース時)

```
$ frtpx -Kfast,parallel sample.f90 -Karraypad_expr=1
```

対象配列を自動選出しパディングを適用

■ 注意事項

- 対象配列を使うソース全てにオプション指定が必要です。
- パディングの効果はプログラムによって異なります。
- 正しくない使い方をした場合には計算結果が異なる場合があります。
- -Karraypad_const [=N]オプションと-Karray_expr=Nオプションは、同時に指定できません。

2次元目の配列要素数を増やすパディング

- 1次元目の配列要素数を増やすパディングで改善されないケース
- 2次元目の配列要素数を増やすパディングについて
- 2次元目の配列要素数を増やすパディング（改善前）
- 2次元目の配列要素数を増やすパディングの効果（ソースチューニング）

1次元目の配列要素数を増やすパディングで改善されないケース

FUJITSU

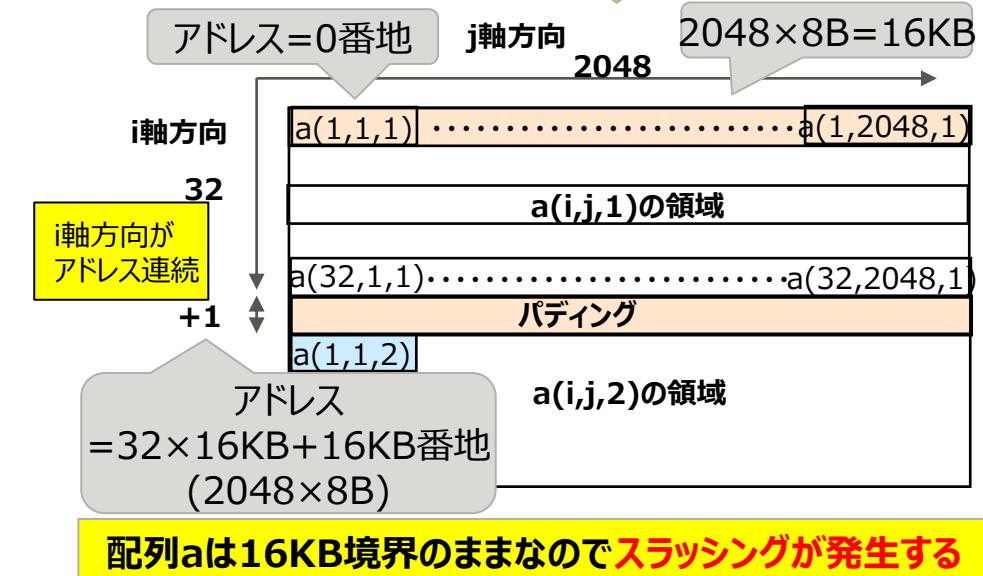
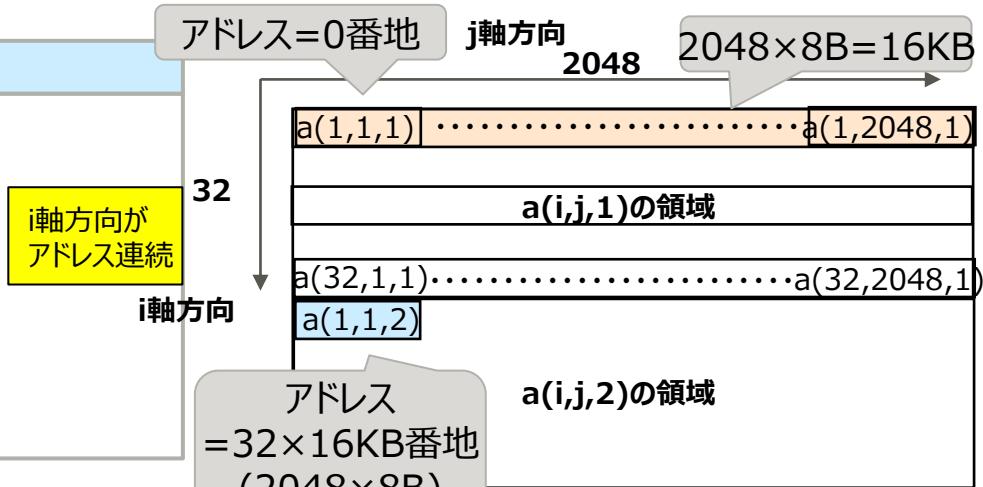
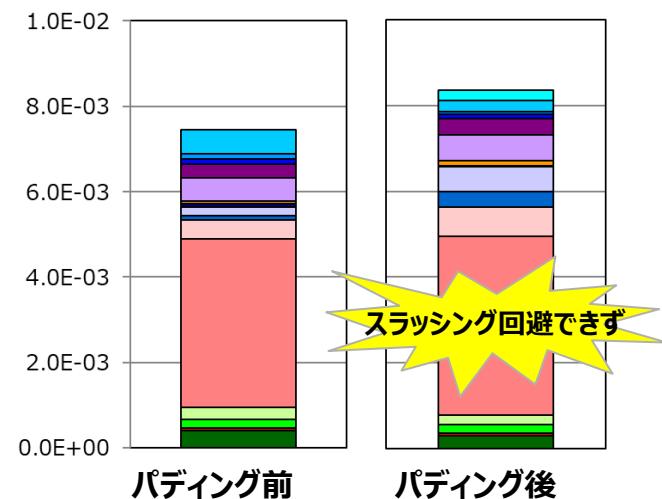
配列のサイズによっては、1次元目の配列要素にパディング(+1)しても改善されないことがあります。

改善前ソース

```
parameter(k=32,l=2048)
real*8 a(k, l, 8)
common /com/a
do j = 1 , l
  do i = 1 , k
    a(i, j, 8) = a(i, j, 1) + a(i, j, 2) + a(i, j, 3) + &
      a(i, j, 4) + a(i, j, 5) + a(i, j, 6) + a(i, j, 7)
  enddo
enddo
end
```

	L1I miss rate (/Effective instruction)	L1D miss	L1D miss rate (/Load-store instruction)	L1D miss demand rate (%) (/L1D miss)
改善前	0.00	1.04E+07	0.39	68.52
改善後	0.00	1.32E+07	0.52	75.20

[秒]



2次元目の配列要素数を増やすパディングについて

FUJITSU

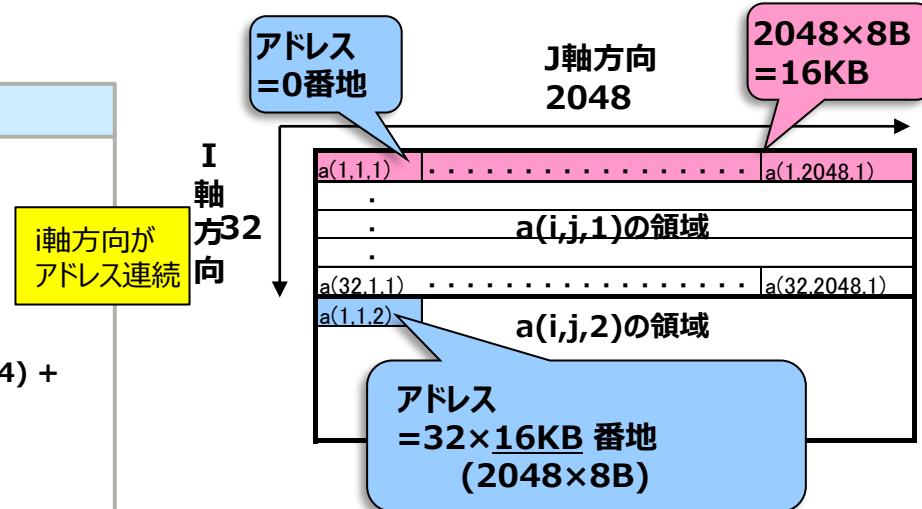
2次元目にパディング(+1)することで16KB境界がずれるため、L1Dキャッシュスラッシングが回避されます。

改善前ソース

```

33 parameter(k=32,l=2048)
34 real*8 a(k, l, 8)
35 common /com/a
36 do j = 1 , l
37  do i = 1 , k
38    a(i, j, 8) = a(i, j, 1) + a(i, j, 2) + a(i, j, 3) + a(i, j, 4) +
                  a(i, j, 5) + a(i, j, 6) + a(i, j, 7)
39 enddo
40 enddo
41 end

```



16KB境界になるのでスラッシングが発生してしまう

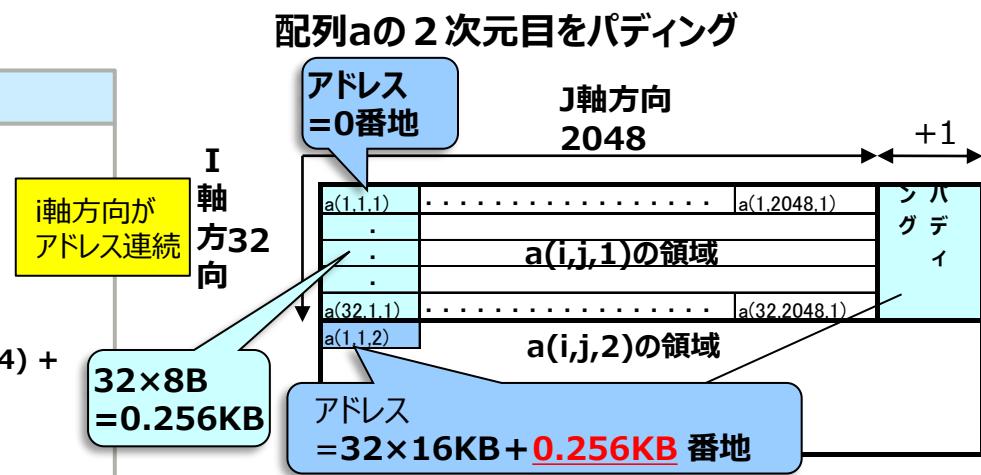
↓

改善後ソース

```

33 parameter(k=32,l=2048)
34 real*8 a(k, l+1, 8)
35 common /com/a
36 do j = 1 , l
37  do i = 1 , k
38    a(i, j, 8) = a(i, j, 1) + a(i, j, 2) + a(i, j, 3) + a(i, j, 4) +
                  a(i, j, 5) + a(i, j, 6) + a(i, j, 7)
39 enddo
40 enddo
41 end

```



16KB境界が崩れるのでスラッシングが回避される

配列aのそれぞれのストリーム同士が16KB境界にあるためL1Dキャッシュスラッシングが発生します。そのため、浮動小数点ロードL2アクセス待ちが多くなっています。

改善前ソース

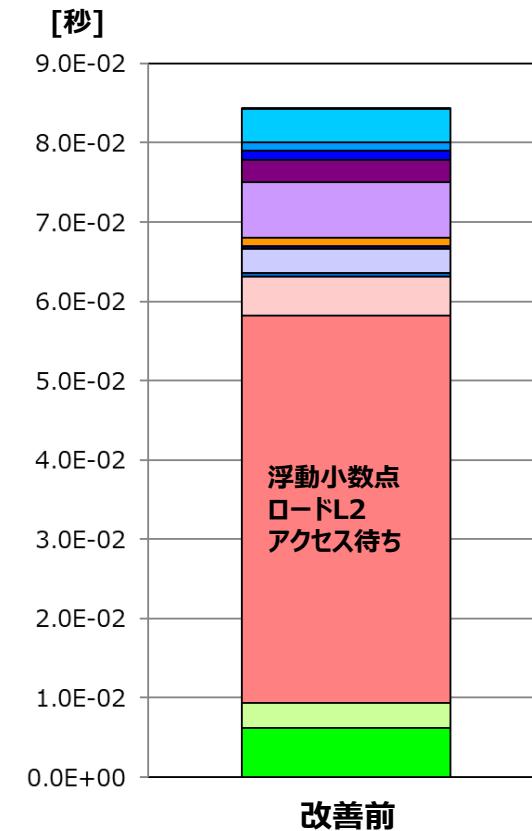
```

38      parameter(n=32,m=2048)
39      real*8 a(n, m, 8)
40      common /com/a
41      <<< Loop-information Start
42      1 pp v      do j = 1 , m
43      2 p   v      a(i, j, 8) = a(i, j, 1) + a(i, j, 2) + a(i, j, 3) + &
44      2           a(i, j, 4) + a(i, j, 5) + a(i, j, 6) + a(i, j, 7)
45      2 p   v      enddo
46      1 p   enddo

```

配列サイズ
32×2048×8B=
32×16KB
(16KB境界)

同一配列の
ストリーム



配列が連続アクセスであるにもかかわらずL1Dミスが高い
⇒ L1Dキャッシュスラッシングが発生している

Cache	L1 miss rate (/Effective instruction)	Load-store instruction	L1D miss	L1D miss rate (/Load-store instruction)	L1D miss demand rate (%) (/L1D miss)	L1D miss hardware prefetch rate (%) (/L1D miss)	software prefetch rate (%) (/L1D miss)	L2 miss	L2 miss rate (/Load-store instruction)	L2 miss demand rate (%) (/L2 miss)	hardware prefetch rate (%) (/L2 miss)	software prefetch rate (%) (/L2 miss)
改善前	0.00	3.88E+08	1.63E+08	0.42	67.76%	32.22%	0.01%	1.01E+04	0.00	71.17%	39.65%	0.00%

配列aのそれぞれのストリームの2次元目にパディング(+1)することで、L1Dキャッシュスラッシング回避しました。その結果、浮動小数点ロードL2アクセス待ちが改善されました。

改善後ソース

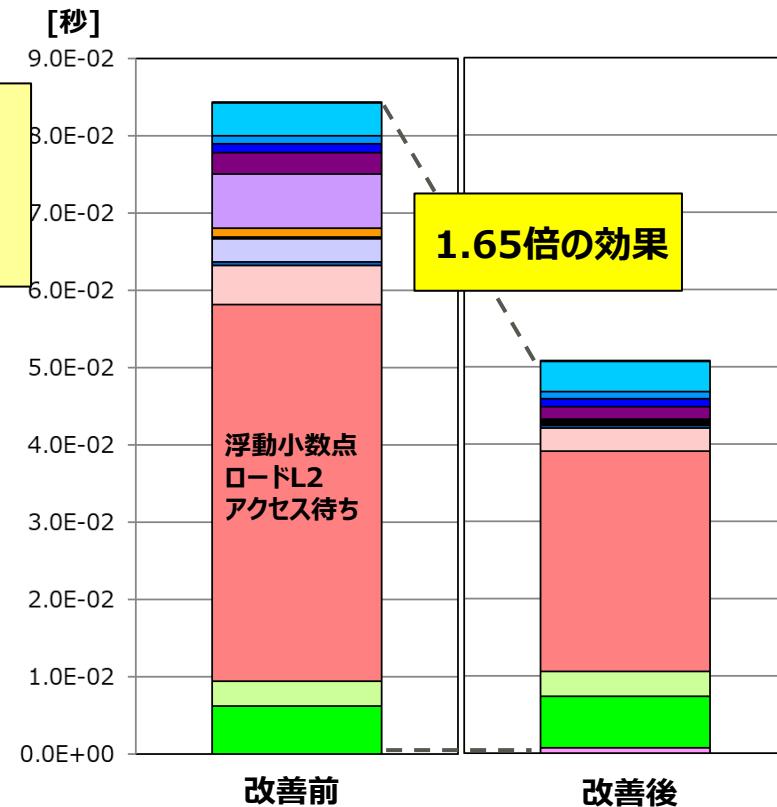
```

39      parameter(n=32,m=2048)
40      real*8 a(n, m+1, 8)
41      common /com/a
<<< Loop-information Start >>>
<<< [PARALLELIZATION]
<<< Standard iteration count
<<< [OPTIMIZATION]
<<< COLLAPSED
<<< SIMD(VL: 8)
<<< SOFTWARE PIPELINING(IPC: 2.66, ITR: 104,
                           MVE: 7, POL: S)
<<< PREFETCH(HARD) Expected by compiler :
<<<     a
<<< Loop-information End >>>
42 1 pp v    do j = 1 , m
<<< Loop-information Start >>>
<<< [OPTIMIZATION]
<<< COLLAPSED
<<< Loop-information End >>>
43 2 p      do i = 1 , n
44 2 p v      a(i, j, 8) = a(i, j, 1) + a(i, j, 2) + a(i, j, 3) + &
45 2           a(i, j, 4) + a(i, j, 5) + a(i, j, 6) + a(i, j, 7)
46 2 p v      enddo
47 1 p      enddo

```

**mを+1することで
16KB境界からずらす**

L1Dミスが減少した



Cache	L1 miss rate (/Effective instruction)	Load-store instruction	L1D miss	L1D miss rate (/Load-store instruction)	L1D miss demand rate (%) / L1D miss	L1D miss hardware prefetch rate (%) (/L1D miss)	L1D miss software prefetch rate (%) (/L1D miss)	L2 miss	L2 miss rate (/Load-store instruction)	L2 miss demand rate (%) / L2 miss	L2 miss hardware prefetch rate (%) (/L2 miss)	L2 miss software prefetch rate (%) (/L2 miss)
改善前	0.00	3.88E+08	1.63E+08	0.42	67.76%	32.22%	0.01%	1.01E+04	0.00	71.17%	39.65%	0.00%
改善後	0.00	3.26E+08	1.07E+08	0.33	51.02%	48.98%	0.00%	9.13E+03	0.00	72.77%	33.09%	0.00%

配列aのそれぞれのストリーム同士が16KB境界にあるためL1Dキャッシュスラッシングが発生します。そのため、浮動小数点ロードL2アクセス待ちが多くなっています。

改善前ソース

```

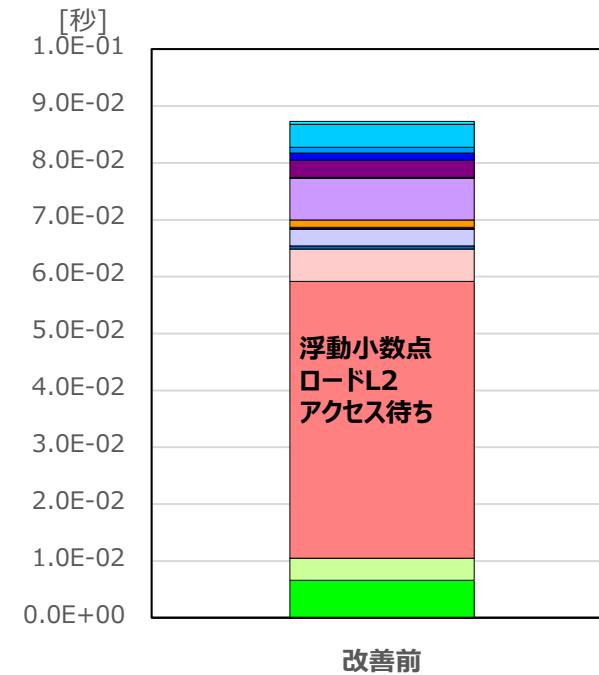
30 void sub(void){
31     int i, j;
32
33     #pragma omp parallel for collapse(2)
34     <<< Loop-information Start >>>
35     <<< [OPTIMIZATION]
36     <<< SIMD(VL: 8)
37     <<< SOFTWARE PIPELINING(IPC: 2.66, ITR: 104, MVE: 7, POL: S)
38     <<< PREFETCH(HARD) Expected by compiler
39
40     p v for (j = 0; j < m; j++) {
41         p v for (i = 0; i < n; i++) {
42             p v a[7][j][i] = a[0][j][i] + a[1][j][i] + a[2][j][i]
43                         + a[3][j][i] + a[4][j][i] + a[5][j][i] + a[6][j][i];
44         p v }
45     p v }
46 }
```

配列宣言
double a[8][2048][32];

aの配列サイズ
32×2048×8B=

32×16KB(16KB境界)

同一配列のストリーム



配列が連続アクセスであるにもかかわらずL1Dミスが高い
⇒ L1Dキャッシュスラッシングが発生している

Cache	L1I miss rate (/Effective instruction)	Load-store instruction	L1D miss	L1D miss rate (/Load-store instruction)	L1D miss demand rate (%) (/L1D miss)	L1D miss hardware prefetch rate (%) (/L1D miss)	L1D miss software prefetch rate (%) (/L1D miss)	L2 miss	L2 miss rate (/Load-store instruction)	L2 miss demand rate (%) (/L2 miss)	L2 miss hardware prefetch rate (%) (/L2 miss)	L2 miss software prefetch rate (%) (/L2 miss)
改善前	0.00	3.92E+08	1.66E+08	0.42	68.39%	31.61%	0.00%	2.04E+04	0.00	39.77%	72.33%	0.00%

配列aのそれぞれのストリームの2次元目にパディング(+1)することで、L1Dキャッシュラッシュを回避しました。その結果、浮動小数点ロードL2アクセス待ちが改善されました。

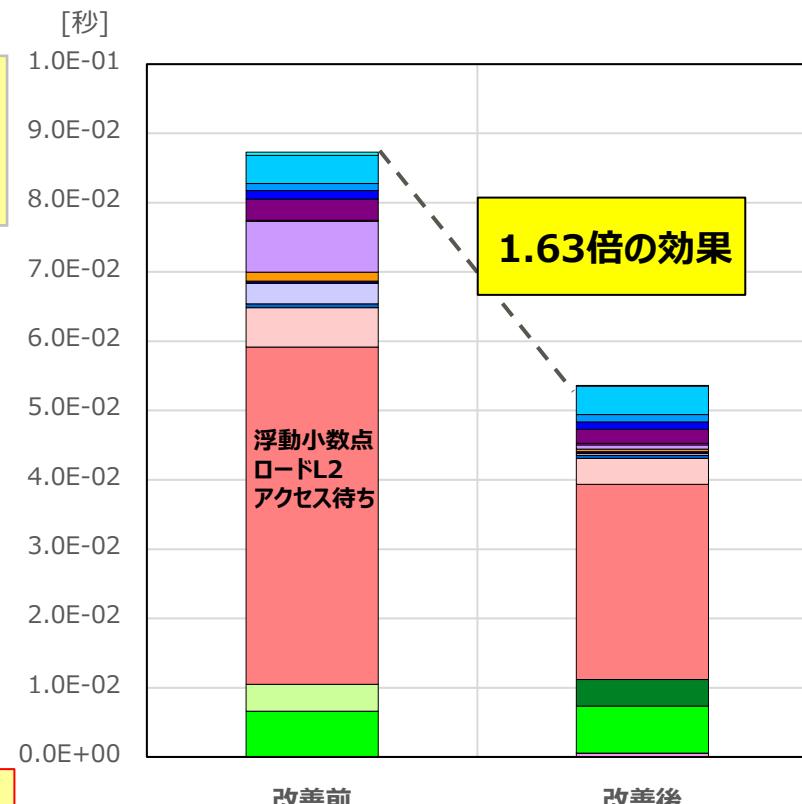
改善後ソース

```

30 void sub(void){
31     int i, j;
32
33     #pragma omp parallel for collapse(2)
<<< Loop-information Start >>>
<<< [OPTIMIZATION]
<<< SIMD(VL: 8)
<<< SOFTWARE PIPELINING(IPC: 2.66, ITR: 104, MVE: 7, POL: S)
<<< PREFETCH(HARD) Expected by compiler :
<<< a
<<< Loop-information End >>>
34 p   v   for (j = 0; j < m; j++){
35 p   v   for (i = 0; i < n; i++){
36 p   v       a[7][j][i] = a[0][j][i] + a[1][j][i] + a[2][j][i] + a[3][j][i]
                     + a[4][j][i] + a[5][j][i] + a[6][j][i];
37 p   v   }
38 p   v   }
39 }
```

配列宣言
double a[8][**2049**][32];
配列aの2次元目を+1することで16KB境界からずらす

L1Dミスが減少した



Cache	L1I miss rate (/Effective instruction)	Load-store instruction	L1D miss	L1D miss rate (/Load-store instruction)	L1D miss demand rate (%) (/L1D miss)	L1D miss hardware prefetch rate (%) (/L1D miss)	L1D miss software prefetch rate (%) (/L1D miss)	L2 miss	L2 miss rate (/Load-store instruction)	L2 miss demand rate (%) (/L2 miss)	L2 miss hardware prefetch rate (%) (/L2 miss)	L2 miss software prefetch rate (%) (/L2 miss)
改善前	0.00	3.92E+08	1.66E+08	0.42	68.39%	31.61%	0.00%	2.04E+04	0.00	39.77%	72.33%	0.00%
改善後	0.00	3.44E+08	1.07E+08	0.31	51.12%	48.88%	0.00%	2.21E+04	0.00	19.24%	84.94%	0.00%

ダミー配列によるパディング

- ・ ダミー配列によるパディング（改善前）
- ・ ダミー配列によるパディングの効果（ソースチューニング）
- ・ ダミー配列によるパディングの効果（翻訳オプションチューニング）

それぞれの配列が16KB境界にあるため、L1Dキャッシュストラッシングが発生します。
そのため、浮動小数点ロードL2アクセス待ちが多くなっています。

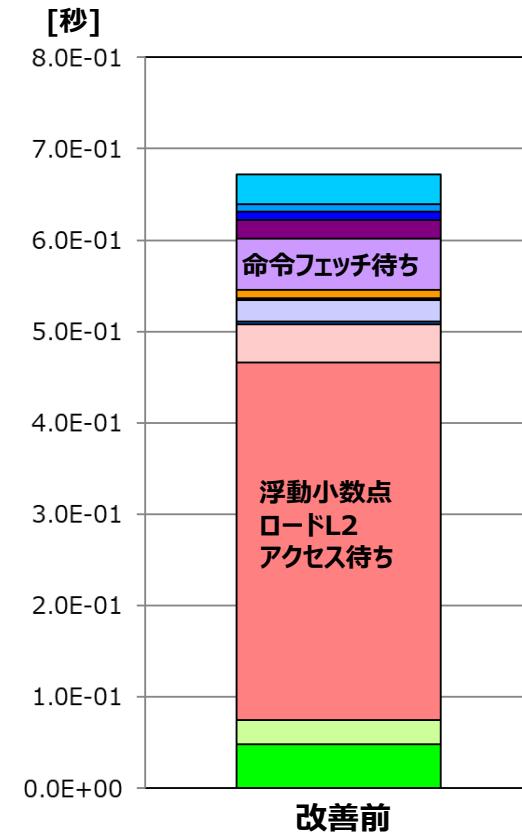
改善前ソース

```

1      parameter(n=256,m=256)
2      real*8 a(n, m),b(n,m),c(n,m),d(n,m),e(n,m),f(n,m),g(n,m),h(n,m)
3      character (1),parameter :: null0=z'00'
4      common /test/a,b,c,d,e,f,g,h
:
27    1 s   s     call sub()
:
34      subroutine sub()
35      parameter(n=256,m=256)
36      real*8 a(n, m),b(n,m),c(n,m),d(n,m),e(n,m),f(n,m),g(n,m),h(n,m)
37      common /test/a,b,c,d,e,f,g,h
<<< Loop-information Start >>>
<<< [PARALLELIZATION]
<<< Standard iteration count: 433
<<< [OPTIMIZATION]
<<< COLLAPSED
<<< SIMD(VL: 8)
<<< SOFTWARE PIPELINING(IPC: 2.66, ITR: 104, MVE: 7, POL: S)
<<< Loop-information End >>>
38    1 pp  v   do j = 1 , m
<<< Loop-information Start >>>
<<< [OPTIMIZATION]
<<< COLLAPSED
<<< Loop-information End >>>
39    2 p   do i = 1 , n
40    2 p   v   a(i, j) = b(i, j) + c(i, j) + d(i, j) + e(i, j) + f(i ,j) + g(i ,j) + h(i ,j)
41    2 p   v   enddo
42    1 p   enddo

```

**配列サイズ
256×256×8B=32×16KB
(16KB境界)**



配列が連続アクセスであるにもかかわらずL1Dミスが高い
⇒ L1Dキャッシュストラッシングが発生している

Cache

	L1 miss rate (/Effective instruction)	Load-store instruction	L1D miss	L1D miss rate (/Load-store instruction)	L1D miss demand rate (%) (/L1D miss)	L1D miss hardware prefetch rate (%) (/L1D miss)	L1D miss software prefetch rate (%) (/L1D miss)	L2 miss	L2 miss rate (/Load-store instruction)	L2 miss demand rate (%) (/L2 miss)	L2 miss hardware prefetch rate (%) (/L2 miss)	L2 miss software prefetch rate (%) (/L2 miss)
改善前	0.00	3.05E+09	1.30E+09	0.43	67.76%	32.24%	0.00%	8.35E+03	0.00	86.64%	19.76%	0.00%

配列間にダミー配列を追加して16KB境界からずらしたため、L1Dキャッシュストラッシングが回避されました。その結果、浮動小数点ロードL2アクセス待ちが改善されました。

改善後ソース

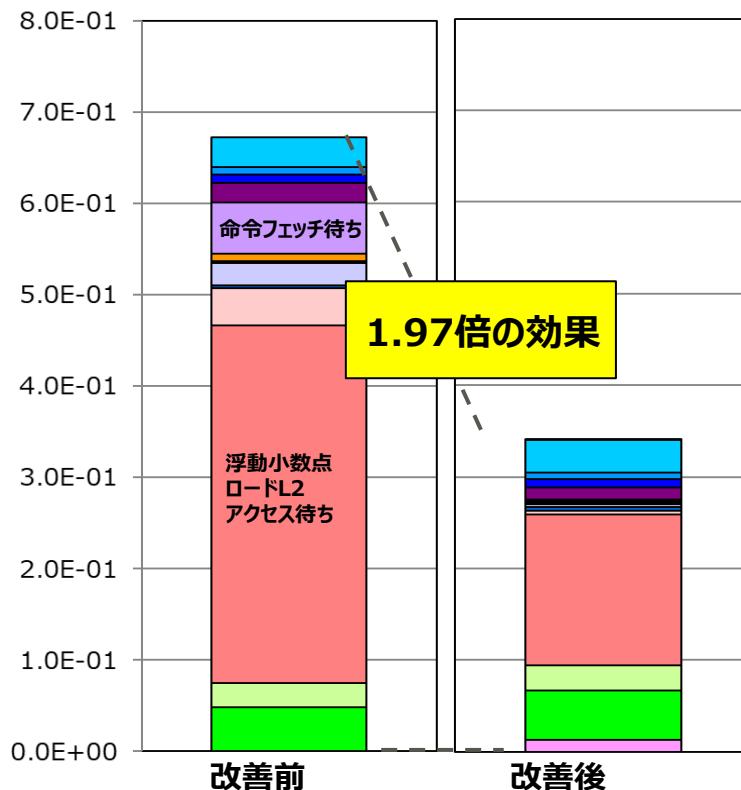
```

1      parameter(n=256,m=256)
2      real*8 a(n, m),dummy1(64),b(n,m),dummy2(64),&
3          c(n,m),dummy3(64),d(n,m),dummy4(64)
4      real*8 e(n, m),dummy5(64),f(n,m),dummy6(64),&
5          g(n,m),dummy7(64),h(n,m)
6      character (1),parameter :: null=' '
7
8      1 s s    call sub()
9
10     subroutine sub()
11
12
13     <<< Loop-information Start >>>
14     <<< [PARALLELIZATION]
15     <<< Standard iteration count: 433
16     <<< [OPTIMIZATION]
17     <<< COLLAPSED
18     <<< SIMD(VL: 8)
19     <<< SOFTWARE PIPELINING(IPC: 2.66, ITR: 104, MVE: 7)
20     <<< Loop-information End >>>
21
22      1 pp v    do j = 1 , m
23      <<< Loop-information Start >>>
24      <<< [OPTIMIZATION]
25      <<< COLLAPSED
26      <<< Loop-information End >>>
27
28      2 p      do i = 1 , n
29      2 p      v      a(i, j) = b(i, j) + c(i, j) + d(i, j) + e(i, j) + f(i,j) + g(i,j) + h(i,j)
30      2 p      v      enddo
31      1 p      enddo

```

配列間にダミー配列を追加することで16KB境界からずらす

[秒]



Cache

L1Dミス、L1Dミスdm率が改善した

	Load-store instruction	L1D miss	L1D miss rate (/Load-store instruction)	L1D miss demand rate (%) (/L1D miss)
改善前	3.05E+09	1.30E+09	0.43	67.76%
改善後	2.79E+09	6.06E+08	0.22	31.03%

それぞれの配列が16KB境界にあるため、L1Dキャッシュラッシングが発生します。
そのため、浮動小数点ロードL2アクセス待ちが多くなっています。

改善前ソース

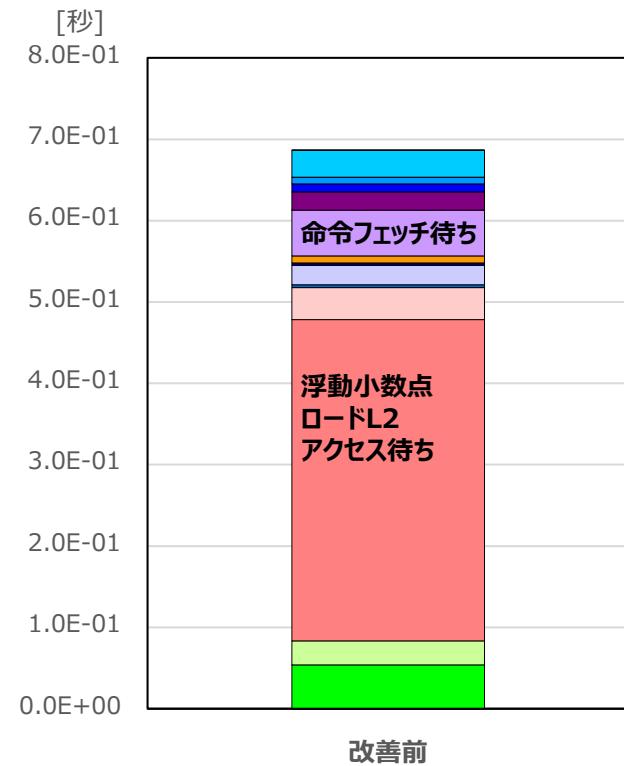
```

35 void sub(void){
36     int i, j;
37
38     #pragma omp parallel for collapse(2)
39     <<< Loop-information Start >>>
40     <<< [OPTIMIZATION]
41     <<< SIMD(VL: 8)
42     <<< SOFTWARE PIPELINING(IPC: 2.66, ITR: 104, MVE: 7, POL: S)
43     <<< PREFETCH(HARD) Expected by compiler :
44     <<< b, c, f, e, g, h, d, a
45     <<< Loop-information End >>>
46
47     p v for (j = 0; j < m; j++){
48         p v for (i = 0; i < n; i++){
49             p v a[j][i] = b[j][i] + c[j][i] + d[j][i] + e[j][i] + f[j][i] + g[j][i] + h[j][i];
50         p v }
51     p v }
52 }
```

配列宣言
`double a[256][256],
b[256][256], c[256][256],
d[256][256], e[256][256],
f[256][256], g[256][256],
h[256][256];`

配列のサイズ256×256×8B=
32×16KB(16KB境界)

配列が連続アクセスであるにもかかわらずL1Dミスが高い
⇒ L1Dキャッシュラッシングが発生している



Cache	L1 miss rate (/Effective instruction)	Load-store instruction	L1D miss	L1D miss rate (/Load-store instruction)	L1D miss demand rate (%) (/L1D miss)	L1D miss hardware prefetch rate (%) (/L1D miss)	L1D miss software prefetch rate (%) (/L1D miss)	L2 miss	L2 miss rate (/Load-store instruction)	L2 miss demand rate (%) (/L2 miss)	L2 miss hardware prefetch rate (%) (/L2 miss)	L2 miss software prefetch rate (%) (/L2 miss)
改善前	0.00	3.12E+09	1.33E+09	0.43	68.30%	31.70%	0.00%	2.34E+04	0.00	46.43%	66.98%	0.00%

配列間にダミー配列を追加して16KB境界からずらしたため、L1Dキャッシュストラッシングが回避されました。その結果、浮動小数点ロードL2アクセス待ちが改善されました。

改善後ソース

```

36     void sub(void){
37         int i, j;
38
39         #pragma omp parallel for collapse(:
40         <<< Loop-information Start >>>
41         <<< [OPTIMIZATION]
42         <<< SIMD(VL: 8)
43         <<< SOFTWARE PIPELINING(IPC:
44         <<< PREFETCH(HARD) Expected b
45         <<< b, c, f, e, g, h, d, a
46         <<< Loop-information End >>>
47         p   for (j = 0; j < m; j++){
48             p   for (i = 0; i < n; i++){
49                 a[j][i] = b[j][i] + c[j][i] + d[j][i] + e[j][i] + f[j][i] + g[j][i] + h[j][i];
50             }
51         }

```

配列a,b,c,d,e,f,g,hの宣言間にダミー配列double型の64要素の配列(例: dummy[64])宣言を追加することで16KB境界からずらす

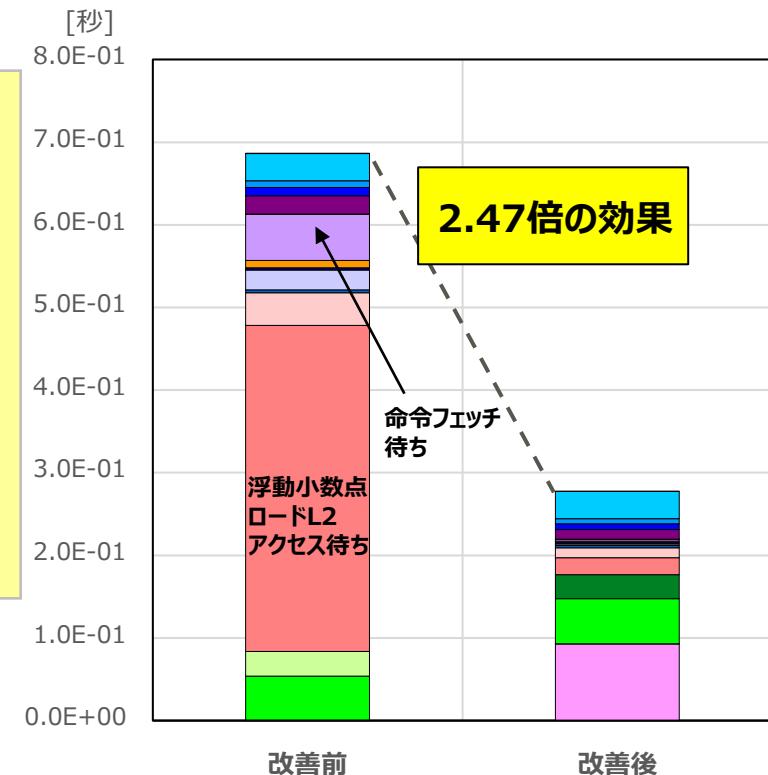
配列宣言

```

double a[256][256],
dummy1[64], b[256][256],
dummy2[64], c[256][256],
dummy3[64], d[256][256],
dummy4[64], e[256][256],
dummy5[64], f[256][256],
dummy6[64], g[256][256],
dummy7[64], h[256][256];

```

L1Dミス、L1Dミスdm率が改善した



Cache	Load-store instruction	L1D miss	L1D miss rate (/Load-store instruction)	L1D miss demand rate (%) (/L1D miss)
改善前	3.12E+09	1.33E+09	0.43	68.30%
改善後	2.67E+09	4.56E+08	0.17	8.37%

以下の翻訳オプション(Fortran固有)を指定することで、ソースチューニングと同等の効果を得ることができます。

翻訳オプション	機能説明
-Kcommonpad[=N] ($4 \leq N \leq 2,147,483,644$)	データキャッシュの使用効率をあげるために共通ブロック内の各変数の領域の間に、隙間をつくることを指定します。 N を省略した場合は、コンパイラが自動的に最良な値を決定します。

■ 使用例（改善前ソース時）

```
$ frtpx -Kfast,parallel sample.f90 -Kcommonpad=512
```

■ 対象配列を自動選出→パディングを適用

■ 注意事項

- 分割コンパイルする場合は、共通ブロックを含むファイルに対して翻訳オプション -Kcommonpadを指定した場合、同じ名前の共通ブロックを含む他のファイルに対しても、翻訳オプション-Kcommonpadを指定しなければなりません。
- 複数のファイルに対して翻訳オプション-Kcommonpad=Nを指定して翻訳する場合、Nの値は同じでなければなりません。
- 同じ共通ブロック名に対し、その要素を変えて使用している場合、翻訳オプション -Kcommonpadを指定するとプログラムが正しく実行されないことがあります。

ダミー配列によるパディング (サイズが異なる配列)

- サイズが異なる配列同士の競合について
- ダミー配列によるパディング（サイズが異なる配列：改善前）
- ダミー配列によるパディングの効果（サイズが異なる配列：ソースチューニング）

サイズが異なる配列同士の競合について (1/2)

一般的にサイズの異なる配列は、定常にキャッシュスラッシングは発生しません。

ソース例

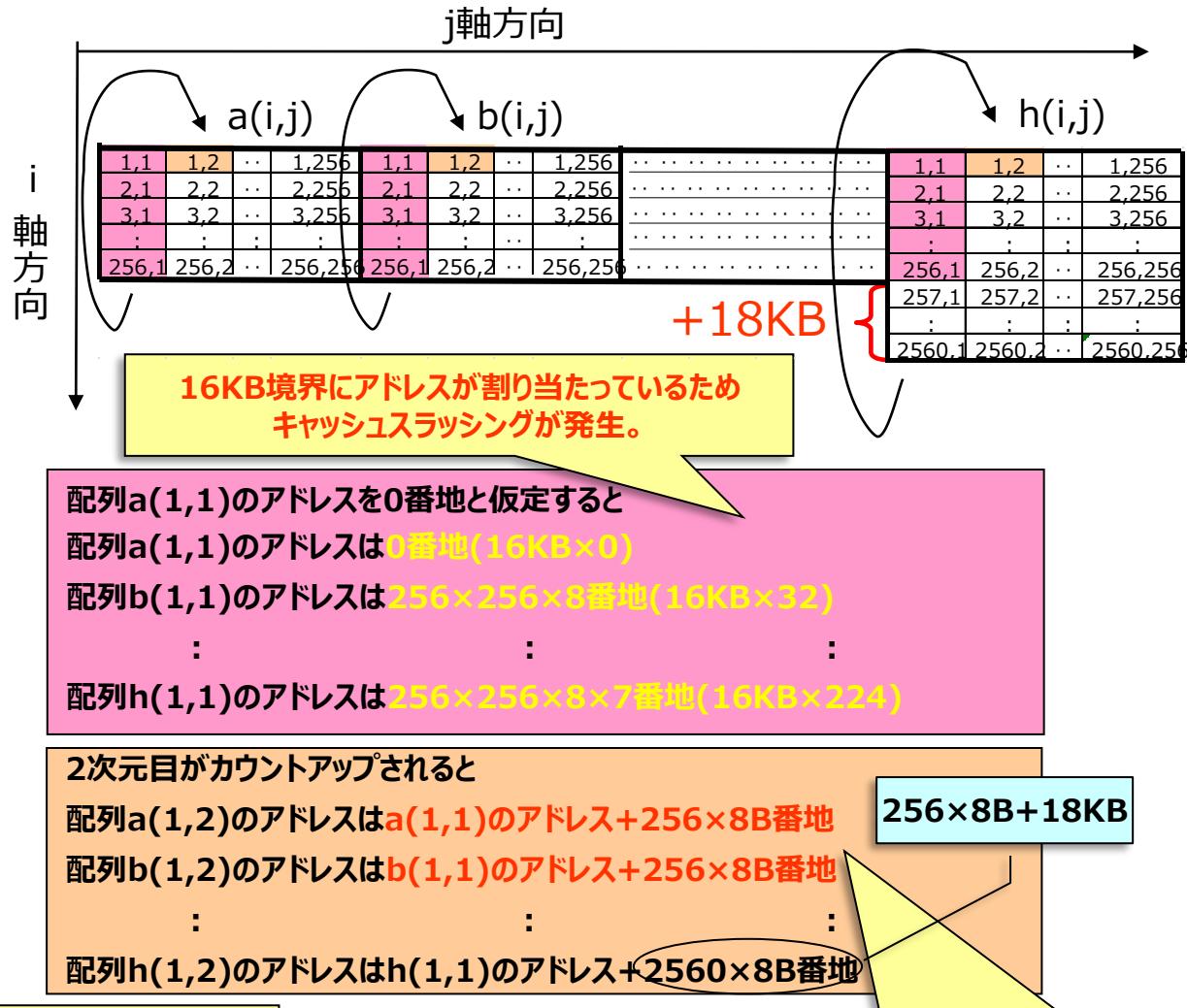
```

parameter(n=256,m=256)
parameter(k=2560,l=256)

real*8 a(n,m), b(n,m), c(n,m),
d(n,m), e(n,m), f(n,m),
g(n,m), h(k,l)

common /test/a,b,c,d,e,f,g,h
do j = 1 , m
  do i = 1 , n
    a(i, j) = b(i, j) + c(i, j) + d(i, j) +
      e(i, j) + f(i, j) + g(i, j) +
      h(i, j)
  enddo
enddo

```



一般的にサイズの異なる配列は、
定常にキャッシュスラッシングは発生しない

配列 a, b, c, d, e, f, g は $16KB$ 境界を維持するが、
配列 h のアドレスは $16KB$ 境界からずれる。

サイズが異なる配列同士の競合について (2/2)

サイズの異なる配列でも配列サイズによっては16KB境界を維持するため、定常にキヤッシュラッキングを発生します。

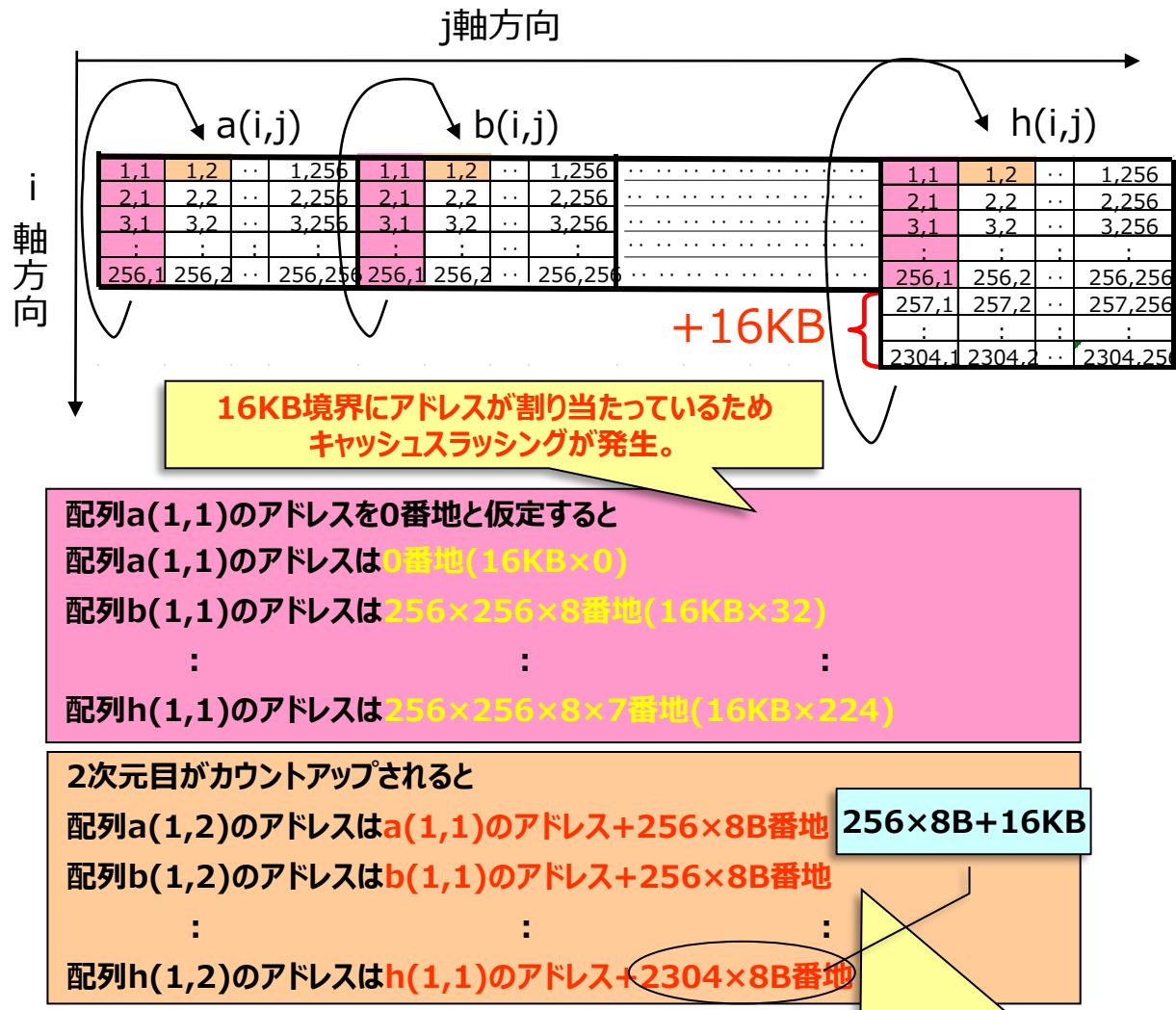
ソース例

```
parameter(n=256,m=256)
```

```
parameter(k=2304,l=256)
```

```
real*8 a(n,m), b(n,m), c(n,m),
d(n,m), e(n,m), f(n,m),
g(n,m), h(k,l)
```

```
common /test/a,b,c,d,e,f,g,h
do j = 1 , m
  do i = 1 , n
    a(i, j) = b(i, j) + c(i, j) + d(i, j) +
      e(i, j) + f(i, j) + g(i, j) +
      h(i, j)
  enddo
enddo
```



配列hを含む全てのサイズの配列に対して
スラッキング対策が必要

配列a,b,c,d,e,f,g,h全て16KB境界を維持する。

それぞれの配列が16KB境界にあるためL1Dキャッシュストラッシングが発生します。
そのため、浮動小数点ロードL2アクセス待ちが多くなっています。

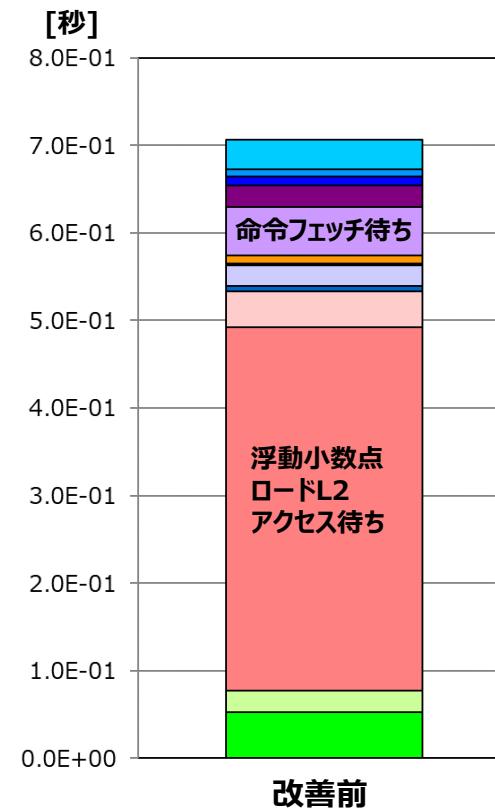
改善前ソース

```

52      integer k,l,n,m
53      parameter(n=256,m=256)
54      parameter(k=2304,l=256)
55
56      real*8 a(n,m), b(n,m), c(n,m), d(n,m), &
57          e(n,m), f(n,m), g(n,m), h(k,l)
58
59      common /test/a,b,c,d,e,
60      <<< Loop-information Start >>>
61      <<< [PARALLELIZATION] >>>
62      <<< Standard iteration control >>>
63      <<< Loop-information End >>>
64
65      1 pp      do j = 1 , m
66      <<< Loop-information Start >>>
67      <<< [OPTIMIZATION] >>>
68      <<< SIMD(VL: 8) >>>
69      <<< SOFTWARE PIPELINING(IPC: 2.83, ITR: 112,
70                      MVE: 13, POL: S) >>>
71      <<< Loop-information End >>>
72
73      2 p v      do i = 1 , n
74      2 p v          a(i, j) = b(i, j) + c(i, j) + d(i, j) + e(i, j) +
75                      f(i, j) + g(i, j) + h(i, j)
76
77      2 p v      enddo
78      1 p      enddo

```

**2次元目がカウントアップされても
各配列間は16KBを維持**



配列が連続アクセスであるにもかかわらずL1Dミスが高い
⇒ L1Dキャッシュストラッシングが発生している

	L1 miss rate (/Effective instruction)	Load-store instruction	L1D miss	L1D miss rate (/Load-store instruction)	L1D miss demand rate (%) (/L1D miss)	L1D miss hardware prefetch rate (%) (/L1D miss)	L1D miss software prefetch rate (%) (/L1D miss)	L2 miss	L2 miss rate (/Load-store instruction)	L2 miss demand rate (%) (/L2 miss)	L2 miss hardware prefetch rate (%) (/L2 miss)	L2 miss software prefetch rate (%) (/L2 miss)
改善前	0.00	3.00E+09	1.39E+09	0.46	70.14%	29.86%	0.00%	3.34E+04	0.00	34.30%	77.94%	0.00%

配列間にダミー配列を追加することで16KB境界からずれたため、L1Dキャッシュストラッシングが回避されました。その結果、浮動小数点ロードL2アクセス待ちが改善されました。

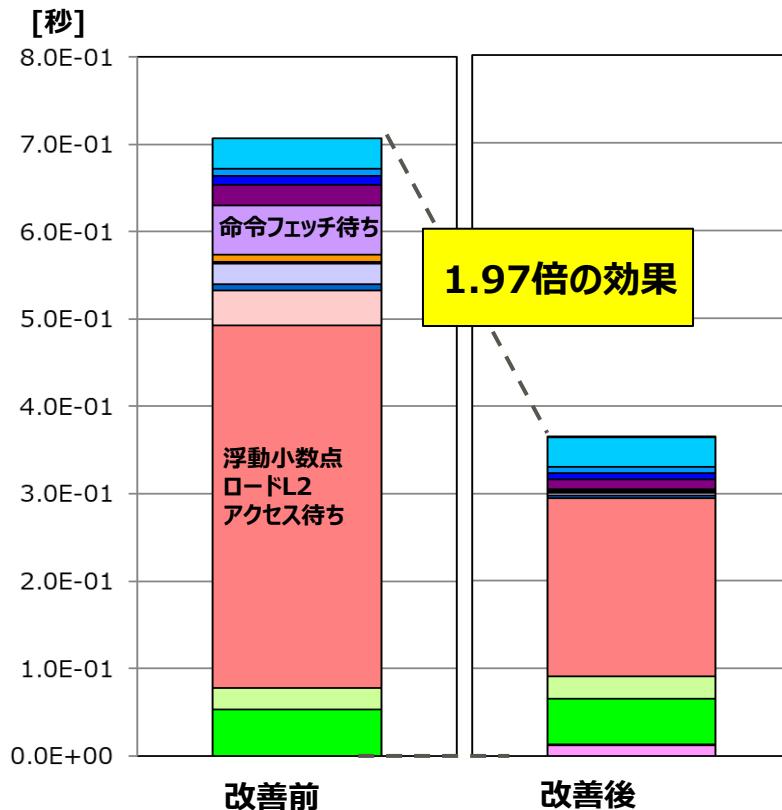
改善後ソース

```

52      integer k,l,n,m
53      parameter(n=256,m=256)
54      parameter(k=2304,l=256)
55
56      real*8 a(n,m),dummy1(64),b(n,m),dummy2(64),&
57          c(n,m),dummy3(64),d(n,m),dummy4(64),&
58          e(n,m),dummy5(64),f(n,m),dummy6(64),&
59          g(n,m),dummy7(64),h(k,l)
60      common /test/a,dummy1,b,dummy2,c,dummy3,d,dummy4,e,
61          dummy5,f,dummy6,g,dummy7,h
62      <<< Loop-information Start >>>
63      <<< [PARALLELIZATION]
64      <<< Standard iteration count: 2
65      <<< Loop-information End >>>
66      1 pp      do j = 1 , m
67      <<< Loop-information Start >>>
68      <<< [OPTIMIZATION]
69      <<< SIMD(VL: 8)
70      <<< SOFTWARE PIPELINING(IPC: 2.83, ITR: 112, MVE: 13, POL: S)
71      <<< Loop-information End >>>
72      2 p      v      do i = 1 , n
73      2 p      v      a(i, j) = b(i, j) + c(i, j) + d(i, j) + e(i, j) + f(i, j) + g(i, j) + h(i, j)
74      2 p      v      enddo
75      1 p      enddo

```

配列間にダミーの配列を挟む



L1Dミス、L1Dミスdm率が改善した。

	L1I miss rate (/Effective instruction)	Load-store instruction	L1D miss	L1D miss rate (/Load-store instruction)	L1D miss demand rate (%) (/L1D miss)	L1D miss hardware prefetch rate (%) (/L1D miss)	software prefetch rate (%) (/L1D miss)	L2 miss	L2 miss rate (/Load-store instruction)	demand rate (%) (/L2 miss)	hardware prefetch rate (%) (/L2 miss)	miss software prefetch rate (%) (/L2 miss)
改善前	0.00	3.00E+09	1.39E+09	0.46	70.14%	29.86%	0.00%	3.34E+04	0.00	34.30%	77.94%	0.00%
改善後	0.00	2.57E+09	6.90E+08	0.27	38.70%	61.31%	0.00%	2.54E+04	0.00	29.53%	80.62%	0.00%

それぞれの配列が16KB境界にあるためL1Dキャッシュストラッシングが発生します。
そのため、浮動小数点ロードL2アクセス待ちが多くなっています。

改善前ソース

```

24 void sub(void){
25     int i, j;
26
27     #pragma omp parallel for
28     <<< Loop-information Start >>>
29     <<< [OPTIMIZATION]
30     <<< PREFETCH(HARD) Expected by compiler :
31     <<< b, c, f, e, g, h, d, a
32     <<< Loop-information End >>>
33     for (j = 0; j < m; j++){
34         <<< Loop-information Start >>>
35         <<< [OPTIMIZATION]
36         <<< SIMD(VL: 8)
37         <<< SOFTWARE PIPELINING(IPC: 2.66, ITR: 104, MVE: 13, POL: S)
38         <<< PREFETCH(HARD) Expected by compiler :
39         <<< b, c, f, e, g, h, d, a
40         <<< Loop-information End >>>
41         p v for (i = 0; i < n; i++){
42             p v a[j][i] = b[j][i] + c[j][i] + d[j][i] + e[j][i] + f[j][i] + g[j][i] + h[j][i];
43         p v }
44         p v }
45     }

```

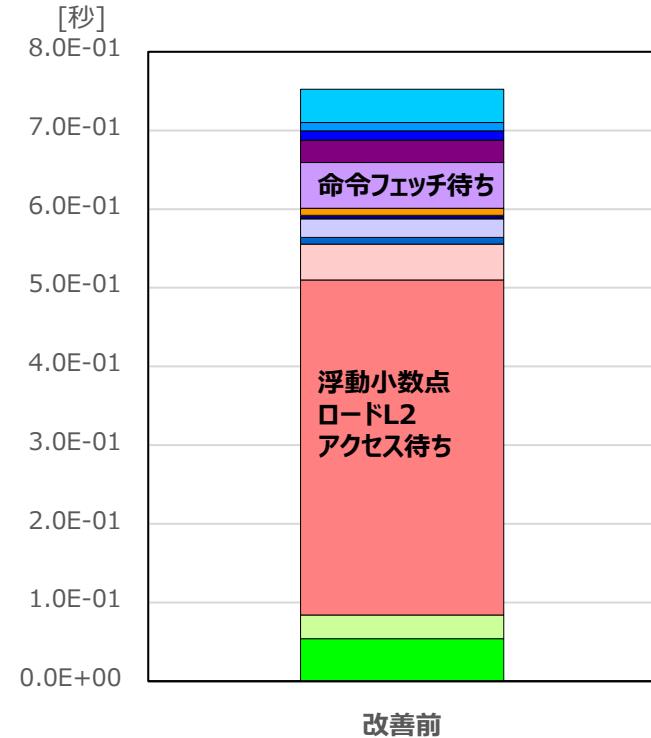
配列宣言

```
double a[256][256],
b[256][256], c[256][256],
d[256][256], e[256][256],
f[256][256], g[256][256],
h[256][2304];
```

配列aのサイズ256×256×8B=32×16KB(16KB境界)

2次元目がカウントアップされても各配列間は16KBを維持

配列が連続アクセスであるにもかかわらずL1Dミスが高い
⇒ L1Dキャッシュストラッシングが発生している



Cache	L1I miss rate (/Effective instruction)	Load-store instruction	L1D miss	L1D miss rate (/Load-store instruction)	L1D miss demand rate (%) (/L1D miss)	L1D miss hardware prefetch rate (%) (/L1D miss)	L1D miss software prefetch rate (%) (/L1D miss)	L2 miss	L2 miss rate (/Load-store instruction)	L2 miss demand rate (%) (/L2 miss)	L2 miss hardware prefetch rate (%) (/L2 miss)	L2 miss software prefetch rate (%) (/L2 miss)
改善前	0.00	3.62E+09	1.38E+09	0.38	71.21%	28.78%	0.01%	1.04E+05	0.00	91.27%	10.77%	0.00%

配列間にダミー配列を追加することで16KB境界からずれたため、L1Dキャッシュストラッシングが回避されました。その結果、浮動小数点ロードL2アクセス待ちが改善されました。

改善後ソース

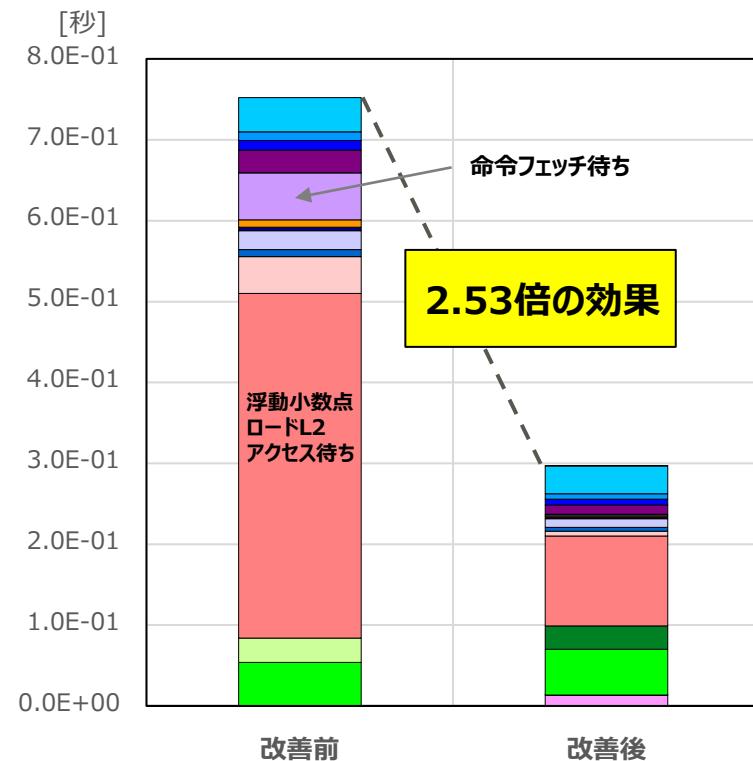
```

25 void sub(void){
26     int i, j;
27
28     #pragma omp parallel for
29     <<< Loop-information Start >>>
30     <<< [OPTIMIZATION]
31     <<< PREFETCH(HARD) Expected
32     <<< b, c, f, e, g, h, d, a
33     <<< Loop-information End >>>
34     p
35     for (j = 0; j < m; j++) {
36         <<< Loop-information Start >>>
37         <<< [OPTIMIZATION]
38         <<< SIMD(VL: 8)
39         <<< SOFTWARE PIPELINING(IPC)
40         <<< PREFETCH(HARD) Expected
41         <<< b, c, f, e, g, h, d, a
42         <<< Loop-information End >>>
43         p
44         v
45         for (i = 0; i < n; i++) {
46             p
47             v
48             a[j][i] = b[j][i] + c[j][i] + d[j][i] + e[j][i] + f[j][i] + g[j][i] + h[j][i];
49         }
50     }
51 }
```

配列a,b,c,d,e,f,g,hの宣言間にダミー配列double型の64要素の配列(例: dummy[64])宣言を追加することで16KB境界からずらす

配列宣言

```
double a[256][256],
dummy1[64], b[256][256],
dummy2[64], c[256][256],
dummy3[64], d[256][256],
dummy4[64], e[256][256],
dummy5[64], f[256][256],
dummy6[64], g[256][256],
dummy7[64], h[256][2304];
```



L1Dミス、L1Dミスdm率が改善した。

Cache	L1 miss rate (/Effective instruction)	Load-store instruction	L1D miss	L1D miss rate (/Load-store instruction)	L1D miss demand rate (%) (/L1D miss)	L1D miss hardware prefetch rate (%) (/L1D miss)	L1D miss software prefetch rate (%) (/L1D miss)	L2 miss	L2 miss rate (/Load-store instruction)	L2 miss demand rate (%) (/L2 miss)	L2 miss hardware prefetch rate (%) (/L2 miss)	L2 miss software prefetch rate (%) (/L2 miss)
改善前	0.00	3.62E+09	1.38E+09	0.38	71.21%	28.78%	0.01%	1.04E+05	0.00	91.27%	10.77%	0.00%
改善後	0.00	2.69E+09	4.81E+08	0.18	17.67%	82.33%	0.00%	1.20E+05	0.00	54.33%	62.57%	0.00%

配列マージ（スラッシングの改善）

- 配列マージとは
- 配列マージ（スラッシングの改善）（改善前）
- 配列マージ（スラッシングの改善）（ソースチューニング）
- 配列マージ（翻訳オプションチューニング）

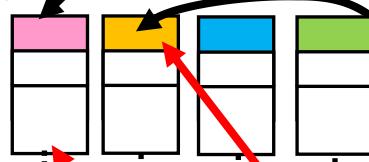
配列マージとは

配列マージとは、複数の配列を1つの配列とするチューニングです。

次の例のように配列数を削減することで、データを同一キャッシュライン上に載せることができます。

改善前

(L1Dキャッシュ)



L1Dキャッシュスラッシング発生

(メモリ上のデータ配置)

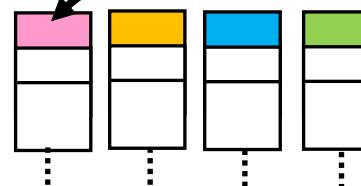
a(1, 1)	256×256×8B =32×16KB 離れたアクセス
a(2, 1)	
...	
a(256, 256)	
b(1, 1)	256×256×8B =32×16KB 離れたアクセス
...	
c(1, 1)	256×256×8B =32×16KB 離れたアクセス
...	
d(1, 1)	256×256×8B =32×16KB 離れたアクセス
...	
e(1, 1)	256×256×8B =32×16KB 離れたアクセス
...	
f(1, 1)	256×256×8B =32×16KB 離れたアクセス
...	
g(1, 1)	256×256×8B =32×16KB 離れたアクセス
...	
h(1, 1)	256×256×8B =32×16KB 離れたアクセス
...	
h(256, 256)	

ソース例

```
subroutine sub()
parameter(n=256,m=256)
real*8 a(n,m),b(n,m),c(n,m),d(n,m),e(n,m),
f(n,m),g(n,m),h(n,m)
common /test/a,b,c,d,e,f,g,h
do j = 1 , m
  do i = 1 , n
    a(i, j) = b(i,j)+c(i,j)+d(i,j)+ &
    e(i,j)+f(i,j)+g(i,j)+h(i,j)
  enddo
enddo
End
```

改善後

(L1Dキャッシュ)



8つとも同一キャッシュライン上のためデータを有効利用できる

(メモリ上のデータ配置)

abcd(1, 1, 1)
abcd(2, 1, 1)
...
abcd(1, 2, 1)
...
abcd(1, 4, 1)
...
efgh(1, 1, 1)
...
efgh(1, 4, 1)
...
abcd(1, 2, 1)
...
abcd(1,256, 1)
abcd(2,256, 1)
abcd(3,256, 1)
abcd(4,256, 1)
...
abcd(4,256,256)
...

ソース例

```
subroutine sub()
parameter(n=256,m=256)
real*8 abcd(n,4,m),efgh(n,4,m)
common /test/abcd,efgh
do j = 1 , m
  do i = 1 , n
    abcd(i,1,j)=abcd(i,2,j)+abcd(i,3,j)+ &
    abcd(i,4,j)+ &
    efgh(i,1,j)+ efgh(i,2,j)+ &
    efgh(i,3,j)+ efgh(i,4,j)
  enddo
enddo
End
```



キャッシュへの格納



キャッシュへの格納（競合）



メモリへのアクセス順番

それぞれの配列が16KB境界にあるためL1Dキャッシュスラッシングが発生します。
そのため、浮動小数点ロードL2アクセス待ちが多くなっています。

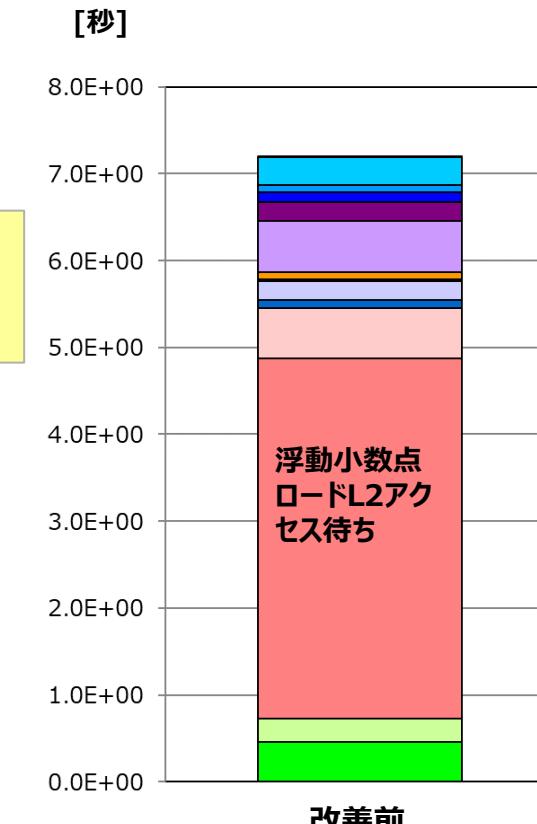
改善前ソース

```

40      parameter(n=256,m=256)
41      real*8 a(n, m),b(n,m),c(n,m),d(n,m),&
42          e(n,m),f(n,m),g(n,m),h(n,m)
43      common /test/a,b,c,d,e,f,g,h
44      <<< Loop-information Start >>>
45      <<< [PARALLELIZATION]
46      <<< Standard iteration count: 433
47      <<< [OPTIMIZATION]
48      <<< COLLAPSED
49      <<< SIMD(VL: 8)
50      <<< SOFTWARE PIPELINING(IPC: 2.66, ITR: 104,
51                      MVE: 7, POL: S)
52      <<< PREFETCH(HARD) Expected by compiler :
53          b, c, f, e, g, h, d, a
54      <<< Loop-information End >>>
55      1 pp   v    do j = 1 , m
56      <<< Loop-information Start >>>
57      <<< [OPTIMIZATION]
58      <<< COLLAPSED
59      <<< Loop-information End >>>
60      2 p     do i = 1 , n
61      2 p   v    a(i,j)=b(i,j)+c(i,j)+d(i,j)+e(i,j)+f(i,j)+g(i,j)+h(i,j)
62      2 p   v    enddo
63      1 p     enddo

```

配列サイズ
256×256×8B=
32×16KB
(16KB境界)



改善前

配列が連続アクセスであるにもかかわらずL1Dミスが高い
⇒ L1Dキャッシュスラッシングが発生している

Cache	L1I miss rate (/Effective instruction)	Load-store instruction	L1D miss	L1D miss rate (/Load-store instruction)	L1D miss demand rate (%) (/L1D miss)	L1D miss hardware prefetch rate (%) (/L1D miss)	L1D miss software prefetch rate (%) (/L1D miss)	L2 miss	L2 miss rate (/Load-store instruction)	L2 miss demand rate (%) (/L2 miss)	L2 miss hardware prefetch rate (%) (/L2 miss)	L2 miss software prefetch rate (%) (/L2 miss)
改善前	0.00	3.40E+10	1.33E+10	0.39	68.21%	31.79%	0.00%	1.12E+04	0.00	72.56%	37.92%	0.00%

配列マージを行うことでストリーム数が8から2となり、L1Dキャッシュスラッシングを回避しました。
その結果、浮動小数点ロードL2アクセス待ちが改善されました。

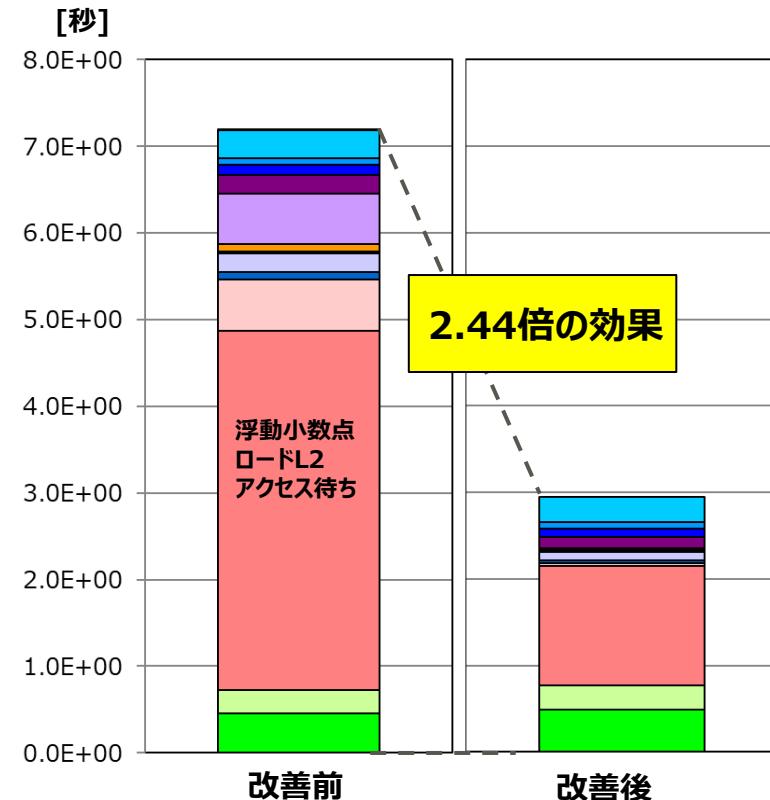
改善後ソース (ソースチューニング)

```

44      parameter(n=256,m=256)
45      real*8 abcd(n,4,m),efgh(n,4,m)
46      common /test/abcd,efgh
47      <<< Loop-information Start >>>
48      1 pp      do j = 1,m
49      2 p v      abcd(i,1,j)=abcd(i,2,j)+abcd(i,3,j)+abcd(i,4,j)+&
50      2           efgh(i,1,j)+efgh(i,2,j)+efgh(i,3,j)+efgh(i,4,j)
51      2 p v      enddo
52      1 p       enddo

```

8つの配列を
4つずつマージする



L1Dミスが減少した

	L1I miss rate (/Effective instruction)	Load-store instruction	L1D miss	L1D miss rate (/Load-store instruction)	L1D miss demand rate (%) (/L1D miss)	L1D miss hardware prefetch rate (%) (/L1D miss)	L1D miss software prefetch rate (%) (/L1D miss)	L2 miss	L2 miss rate (/Load-store instruction)	L2 miss demand rate (%) (/L2 miss)	L2 miss hardware prefetch rate (%) (/L2 miss)	L2 miss software prefetch rate (%) (/L2 miss)
改善前	0.00	3.40E+10	1.33E+10	0.39	68.21%	31.79%	0.00%	1.12E+04	0.00	72.56%	37.92%	0.00%
改善後	0.00	2.54E+10	4.40E+09	0.17	84.98%	15.03%	0.00%	1.70E+04	0.00	69.93%	46.61%	0.00%

それぞれの配列が16KB境界にあるためL1Dキャッシュスラッシングが発生します。
そのため、浮動小数点ロードL2アクセス待ちが多くなっています。

改善前ソース

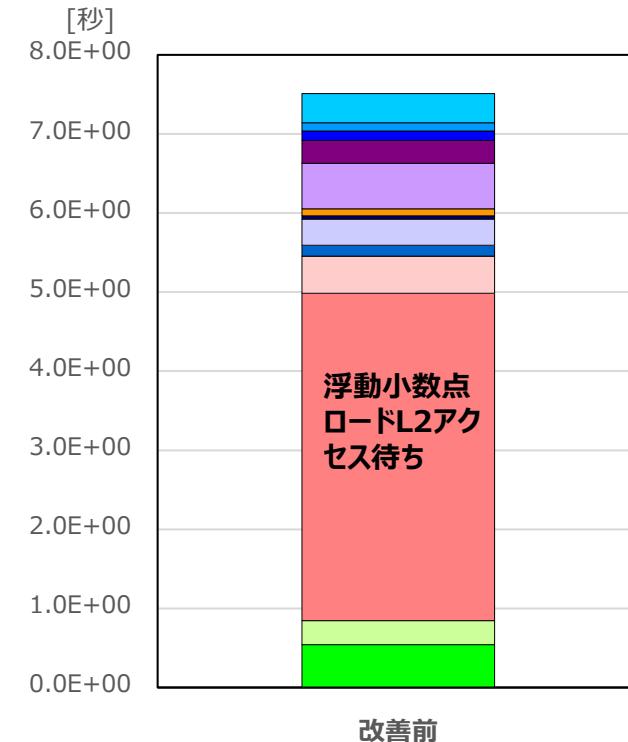
```

31     void sub(void)
32     {
33         int i,j;
34         #pragma omp parallel for
35         <<< Loop-information Start >>>
36         <<< [OPTIMIZATION]
37         <<< PREFETCH(HARD) Expected by compiler :
38         <<< b, c, f, e, g, h, d, a
39         <<< Loop-information End >>>
40         p     for(j=0;j<m;j++){
41             <<< Loop-information Start >>>
42             <<< [OPTIMIZATION]
43             <<< SIMD(VL: 8)
44             <<< SOFTWARE PIPELINING(IPC: 2.66, ITR: 104, MVE: 7, POL: S)
45             <<< PREFETCH(HARD) Expected by compiler :
46             <<< b, c, f, e, g, h, d, a
47             <<< Loop-information End >>>
48         p     v     for(i=0;i<n;i++){
49             p     v     a[j][i]=b[j][i]+c[j][i]+d[j][i]+e[j][i]+f[j][i]+g[j][i]+h[j][i];
50         p     v     }
51     }

```

配列宣言
 double a[256][256],
 b[256][256], c[256][256],
 d[256][256], e[256][256],
 f[256][256], g[256][256],
 h[256][256];
 配列aのサイズ256×256×8B=
32×16KB(16KB境界)

配列が連続アクセスであるにもかかわらずL1Dミスが高い
 ⇒ L1Dキャッシュスラッシングが発生している



改善前

Cache	L1 miss rate (/Effective instruction)	Load-store instruction	L1D miss	L1D miss rate (/Load-store instruction)	L1D miss demand rate (%) (/L1D miss)	L1D miss hardware prefetch rate (%) (/L1D miss)	L1D miss software prefetch rate (%) (/L1D miss)	L2 miss	L2 miss rate (/Load-store instruction)	L2 miss demand rate (%) (/L2 miss)	L2 miss hardware prefetch rate (%) (/L2 miss)	L2 miss software prefetch rate (%) (/L2 miss)
改善前	0.00	3.32E+10	1.36E+10	0.41	69.48%	30.52%	0.00%	1.16E+05	0.00	86.19%	26.53%	0.00%

配列マージを行うことでストリーム数が8から2となり、L1Dキャッシュスラッシングを回避しました。その結果、浮動小数点ロードL2アクセス待ちが改善されました。

改善後ソース（ソースチューニング）

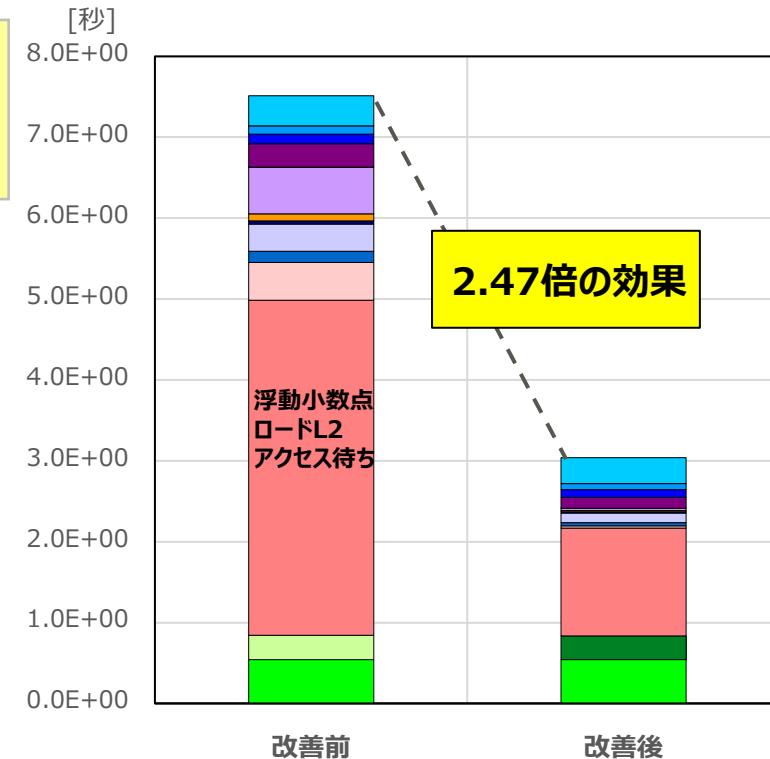
```

31 void sub(void)
32 {
33     int i,j;
34     #pragma omp parallel for
35     <<< Loop-information Start >>>
36     <<< [OPTIMIZATION]
37     <<< PREFETCH(HARD) Expected by compiler :
38     <<< egh, abcd
39     <<< Loop-information End >>>
40     p   for(j=0;j<m;j++){
41         <<< Loop-information Start >>>
42         <<< [OPTIMIZATION]
43         <<< SIMD(VL: 8)
44         <<< SOFTWARE PIPELINING(IPC: 2.28, ITR: 112, MVE: 8, POL: S)
45         <<< PREFETCH(HARD) Expected by compiler :
46         <<< egh, abcd
47         <<< Loop-information End >>>
48     p   v   for(i=0;i<n;i++){
49         p   v   abcd[j][0][i]=abcd[j][1][i]+abcd[j][2][i]+abcd[j][3][i]+
50             egh[j][0][i]+egh[j][1][i]+egh[j][2][i]+egh[j][3][i];
51     p   v   }
52     p   }
53 }
```

配列宣言
double abcd[256][4][256],
efgh[256][4][256];

8つの配列を4つずつマージする

L1Dミスが減少した



Cache	L1I miss rate (/Effective instruction)	Load-store instruction	L1D miss	L1D miss rate (/Load-store instruction)	L1D miss demand rate (%) (/L1D miss)	L1D miss hardware prefetch rate (%) (/L1D miss)	L1D miss software prefetch rate (%) (/L1D miss)	L2 miss	L2 miss rate (/Load-store instruction)	L2 miss demand rate (%) (/L2 miss)	L2 miss hardware prefetch rate (%) (/L2 miss)	L2 miss software prefetch rate (%) (/L2 miss)
改善前	0.00	3.32E+10	1.36E+10	0.41	69.48%	30.52%	0.00%	1.16E+05	0.00	86.19%	26.53%	0.00%
改善後	0.00	2.63E+10	4.40E+09	0.17	90.03%	9.98%	-0.01%	4.45E+04	0.00	80.52%	35.47%	0.00%

以下の翻訳オプション（Fortran固有）を指定することで、ソースチューニングと同等の効果を得ることができます。

翻訳オプション	機能説明
-Karray_merge_common [=name]	共通ブロック内の複数配列をマージすることを指示します。nameには共通ブロック名を指定できます。Nameを省略した場合、すべて名前付き共通ブロック内の配列が対象となります。
-Karray_merge_local	ローカル配列の複数配列をマージすることを指示します。 -Karray_merge_local_size=1000000も同時に有効になります。
-Karray_merge_local_size=N ($2 \leq N \leq 2,147,483,647$)	マージ対象とするローカル配列の大きさが、Nバイト以上であることを指示します。-Karray_merge_localオプションが有効な場合に意味があります。
-Karray_merge	本オプションは、-Karray_merge_localおよび-Karray_merge_commonオプションが指定された場合と等価です。

■ 使用例（改善前ソース時）

```
$frtpx -Kfast,parallel sample.f90 -Karray_merge_common
```

■ 注意事項

- 対象配列を使うソースすべてにオプション指定が必要です。
- マージの効果はプログラムによって異なります。
- 正しくない使い方をした場合には計算結果が異なることがあります。
- デバッグオプション(-gおよび-H)と併用はできません。

ループ分割（スラッシングの改善）

- ループ分割とは
- ループ分割（スラッシングの改善）（改善前）
- ループ分割の効果（スラッシングの改善）（ソースチューニング）
- ループ分割の効果（最適化制御行チューニング）

ループ分割とは

- ループ分割とは、主に以下の目的でループ[°]を複数の小さなループに分割する手法です。
- ソフトウェアパイプラインの促進
- キャッシュメモリ利用効率の改善
- レジスタ不足の解消

ループ分割を行うことによりループ内でアクセスする配列数を削減し、ソフトウェアパイプラインの促進や、キャッシュスラッシング回避できる場合があります。

ただし、分割の仕方によってはキャッシュ上のデータを有効利用できなくなる場合があるため注意が必要です。

改善前ソース例

```
parameter(n=65536)
real*8 a(n),b(n),c(n),d(n),e(n),f(n),
g(n),h(n)
common /com/a,b,c,d,e,f,g,h
do i=1,n
  a(i) = s / b(i)
  c(i) = s / d(i)
  e(i) = s / f(i)
  g(i) = s / h(i)
enddo
```

スラッシングが発生



改善後ソース例

```
parameter(n=65536)
real*8 a(n),b(n),c(n),d(n),e(n),f(n),
g(n),h(n)
common /com/a,b,c,d,e,f,g,h
!OCL LOOP_NOFUSION
do i=1,n
  a(i) = s / b(i)
  c(i) = s / d(i)
enddo
do i=1,n
  e(i) = s / f(i)
  g(i) = s / h(i)
enddo
```

ループ融合を抑止

ループを分割
スラッシングを抑止

それぞれの配列が16KB境界にあるためL1Dキャッシュスラッシング発生します。
そのため、浮動小数点ロードアクセス待ちが多くなっています。

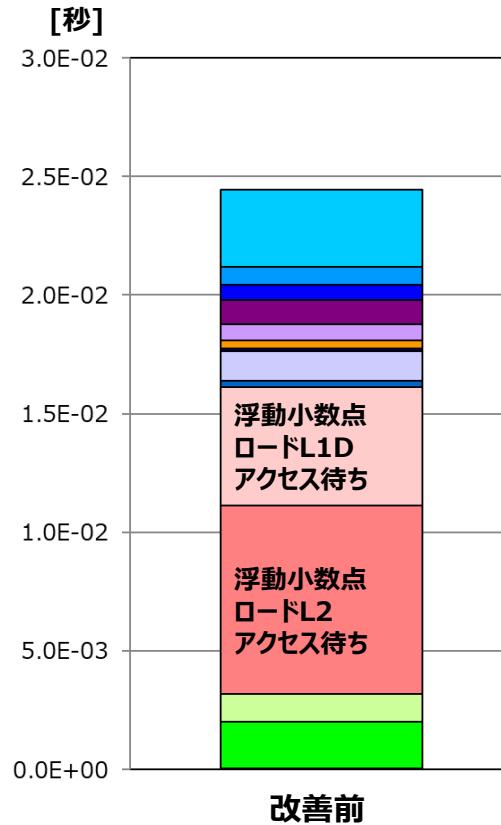
改善前ソース

```

46      parameter(n=65536)
47      real*8 a(n),b(n),c(n),d(n),e(n),f(n),g(n),h(n)
48      common /com/a,b,c,d,e,f,g,h
49
<<< Loop-information Start >>
<<< [PARALLELIZATION]
<<< Standard iteration count:
<<< [OPTIMIZATION]
<<< SIMD(VL: 8)
<<< SOFTWARE PIPELINING(IPC: 1.90, ITR: 56,
                           MVE: 4, POL: S)
<<< PREFETCH(HARD) Expected by compiler :
<<<   f, d, b, h, g, e, c, a
<<< Loop-information End >>
50      1 pp v    do i=1,n
51      1 p  v      a(i) = s / b(i)
52      1 p  v      c(i) = s / d(i)
53      1 p  v      e(i) = s / f(i)
54      1 p  v      g(i) = s / h(i)
55      1 p  v      enddo

```

配列サイズ
65536×8B=
32×16KB
(16KB境界)



Cache

	L1I miss rate (/Effective instruction)	Load-store instruction	L1D miss	L1D miss rate (/Load-store instruction)	L1D miss demand rate (%) (/L1D miss)	L1D miss hardware prefetch rate (%) (/L1D miss)	prefetch rate (%) (/L1D miss)	L2 miss	(/Load-store instruction)	demand rate (%) (/L2 miss)	prefetch rate (%) (/L2 miss)	prefetch rate (%) (/L2 miss)
改善前	0.00	1.25E+08	3.85E+07	0.31	57.48%	42.52%	0.01%	1.84E+04	0.00	31.49%	72.74%	0.00%

配列が連続アクセスであるにも関わらず
L1Dミス率が高く、L1ミスdm率が高い
⇒ L1Dキャッシュスラッシングが発生している

ループ分割を行うことでストリーム数が8から4となり、L1Dキャッシュスラッシングを回避しました。
その結果、浮動小数点ロードL2アクセス待ちが改善されました。

改善後ソース

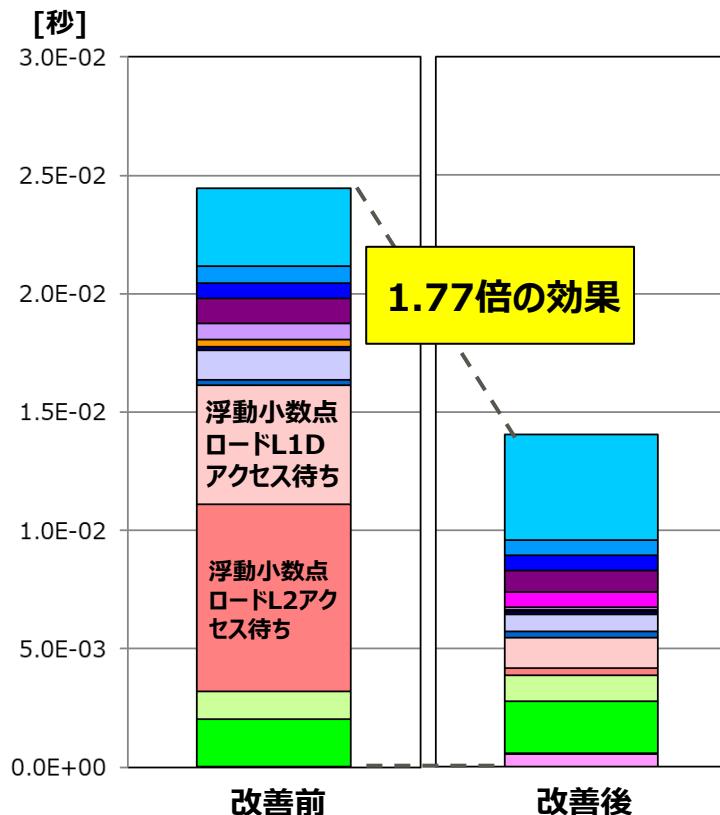
```

46      parameter(n=65536)
47      real*8 a(n),b(n),c(n),d(n),e(n),f(n),g(n),h(n)
48      common /com/a,b,c,d,e,f,g,h
49
50      !OCL LOOP_NOFUSION
51      <<< Loop-information Start >>>
52      <<< [PARALLELIZATION]
53      <<< Standard iteration count: 411
54      <<< [OPTIMIZATION]
55      <<< SIMD(VL: 8)
56      <<< SOFTWARE PIPELINING(IPC: 2.36, ITR: 80,
57                                MVE: 2, POL: S)
58
59      <<< Loop-information End >>>
60      1 pp 2v    do i=1,n
61      1 p 2v      a(i) = s / b(i)
62      1 p 2v      c(i) = s / d(i)
63      1 p 2v    enddo
64      <<< Loop-information Start >>>
65      <<< [PARALLELIZATION]
66      <<< Standard iteration count: 411
67      <<< [OPTIMIZATION]
68      <<< SIMD(VL: 8)
69      <<< SOFTWARE PIPELINING(IPC: 2.45, ITR: 80,
70                                MVE: 2, POL: S)
71
72      <<< Loop-information End >>>
73      1 pp 2v    do i=1,n
74      1 p 2v      e(i) = s / f(i)
75      1 p 2v      g(i) = s / h(i)
76      1 p 2v    enddo

```

ループ融合を抑止

ループ分割



L1Dミス、L1Dミスdm率が減少した。

	L1I miss rate (/Effective instruction)	Load-store instruction	L1D miss	L1D miss rate (/Load-store instruction)	L1D miss demand rate (%) (/L1D miss)	L1D miss hardware prefetch rate (%) (/L1D miss)	L1D miss software prefetch rate (%) (/L1D miss)
改善前	0.00	1.25E+08	3.85E+07	0.31	57.48%	42.52%	0.01%
改善後	0.00	1.10E+08	1.78E+07	0.16	8.37%	91.63%	0.00%

それぞれの配列が16KB境界にあるためL1Dキャッシュスラッシング発生します。
そのため、浮動小数点ロードアクセス待ちが多くなっています。

改善前ソース

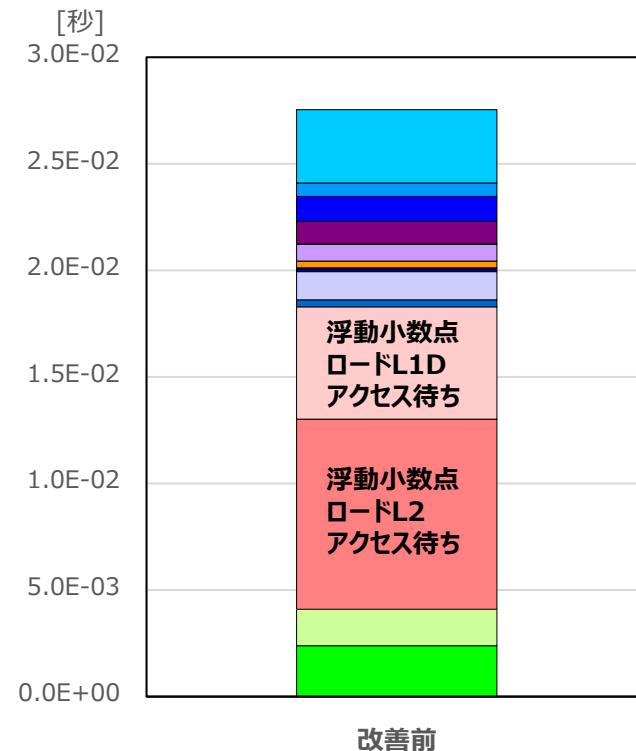
```

45     void sub(double s)
46     {
47         int i;
48
49         #pragma omp parallel for
50         <<< Loop-information Start >>>
51         <<< [OPTIMIZATION]
52         <<< SIMD(VL: 8)
53         <<< SOFTWARE PIPELINING(IPC: 1.90, ITR: 56,
54             MVE: 4, POL: S)
55         <<< PREFETCH(HARD) Expected by compiler :
56         <<< f, d, b, h, g, e, c, a
57         <<< Loop-information End >>>
58         p v for(i=0;i<n; i++)
59         p v {
60             p v     a[i] = s / b[i];
61             p v     c[i] = s / d[i];
62             p v     e[i] = s / f[i];
63             p v     g[i] = s / h[i];
64             p v }
65
66         return;
67     }

```

配列宣言
`double a[65536],
b[65536], c[65536],
d[65536], e[65536],
f[65536], g[65536],
h[65536];`

配列a,b,c,d,e,f,g,h
のサイズ65536 x 8B=
32×16KB(16KB境界)



配列が連続アクセスであるにも関わらず
L1Dミス率が高く、L1ミスdm率が高い
⇒ **L1Dキャッシュスラッシングが発生している**

Cache	L1 miss rate (/Effective instruction)	Load-store instruction	L1D miss	L1D miss rate (/Load-store instruction)	L1D miss demand rate (%) (/L1D miss)	L1D miss hardware prefetch rate (%) (/L1D miss)	prefetch rate (%) (/L1D miss)	L2 miss	(/Load-store instruction)	demand rate (%) (/L2 miss)	prefetch rate (%) (/L2 miss)	software prefetch rate (%) (/L2 miss)
改善前	0.00	1.17E+08	3.93E+07	0.34	58.44%	41.56%	0.00%	1.94E+04	0.00	16.93%	87.92%	0.00%

ループ分割を行うことでストリーム数が8から4となり、L1Dキャッシュスラッシングを回避しました。
その結果、浮動小数点ロードL2アクセス待ちが改善されました。

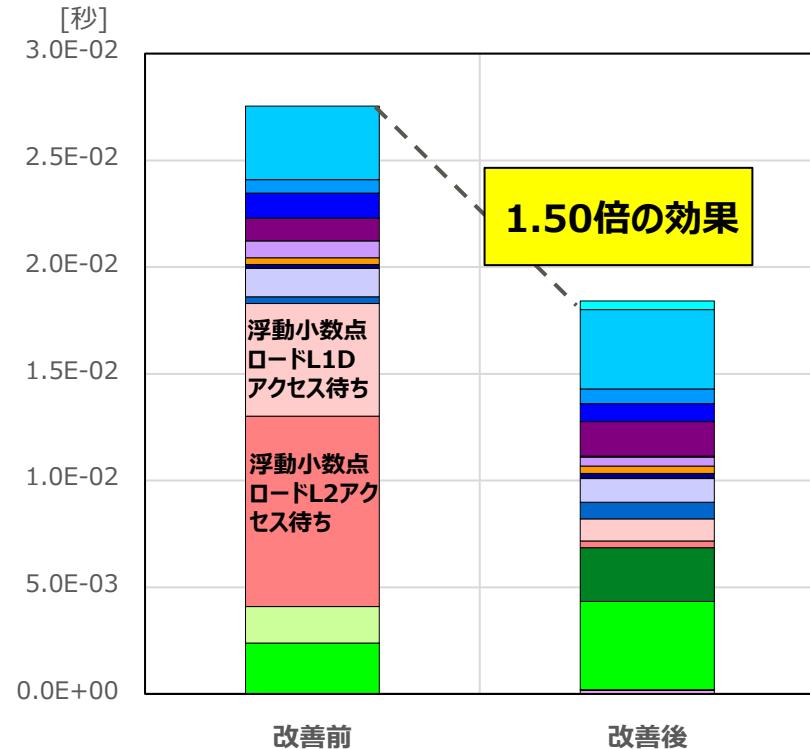
改善後ソース

```

45 void sub(double s)
46 {
47     int i;
48
49     #pragma omp parallel for
<<< Loop-information Start >>>
<<< [OPTIMIZATION]
<<< SIMD(VL: 8)
<<< SOFTWARE PIPELINING(IPC: 2.09, ITR: 64,
    MVE: 2, POL: S)
<<< PREFETCH(HARD) Expected by compiler :
<<< d, b, c, a
<<< Loop-information End >>>
50 p 2v for(i=0;i<n; i++)
51 p 2v {
52 p 2v   a[i] = s / b[i];
53 p 2v   c[i] = s / d[i];
54 p 2v }
55
56     #pragma omp parallel for
<<< Loop-information Start >>>
<<< [OPTIMIZATION]
<<< SIMD(VL: 8)
<<< SOFTWARE PIPELINING(IPC: 2.09, ITR: 64,
    MVE: 2, POL: S)
<<< PREFETCH(HARD) Expected by compiler :
<<< h, f, g, e
<<< Loop-information End >>>
57 p 2v for(i=0;i<n; i++)
58 p 2v {
59 p 2v   e[i] = s / f[i];
60 p 2v   g[i] = s / h[i];
61 p 2v }
62     return;
63 }
```

配列宣言
double a[65536],
b[65536], c[65536],
d[65536], e[65536],
f[65536], g[65536],
h[65536];

ループ分割



L1Dミス、L1Dミスdm率が減少した。

Cache	L1I miss rate (/Effective instruction)	Load-store instruction	L1D miss	L1D miss rate (/Load-store instruction)	L1D miss demand rate (%) (/L1D miss)	L1D miss hardware prefetch rate (%) (/L1D miss)	L1D miss software prefetch rate (%) (/L1D miss)
改善前	0.00	1.17E+08	3.93E+07	0.34	58.44%	41.56%	0.00%
改善後	0.00	1.51E+08	1.86E+07	0.12	12.30%	87.72%	-0.02%

ループ分割の効果（最適化制御行チューニング）

FUJITSU

以下の最適化制御行を指定することで、ソースチューニングと同等の効果を得ることができます。

最適化指示子 (Fortran)	意味	指定可能な最適化行			
		プログラム単位	DOループ単位	文単位	配列代入文単位
FISSION_POINT[(n1)] (n1は1~6の10進数)	ループ内の指定された位置でループ分割することを指示します。最内ループから数えてn1重ネストされた多重ループを対象にループします。	×	×	○	×

最適化指示子 (C/C++のtradモードのみ)	意味	指定可能な最適化行			
		global行	procedure行	loop行	statement行
fission_point[(n1)] (n1は1~6の10進数)	ループ内の指定された位置でループ分割することを指示します。最内ループから数えてn1重ネストされた多重ループを対象にループします。	×	×	×	○

改善後ソース（最適化制御行チューニング）（Fortranの例）

```

46      parameter(n=65536)
47      real*8 a(n),b(n),c(n),d(n),e(n),f(n),g(n),h(n)
48      common /com/a,b,c,d,e,f,g,h
49
      <<< Loop-information Start >>>
      <<< [PARALLELIZATION]
      <<< Standard iteration count: 411
      <<< [OPTIMIZATION]
      <<< FISSION(num: 2)
      <<< SIMD(VL: 8)
      <<< SOFTWARE PIPELINING(IPC: 2.45, ITR: 80, MVE: 2, POL: S)
      <<< PREFETCH(HARD) Expected by compiler :
      <<<     b, d, c, a, f, h, g, e
      <<< Loop-information End >>>
50  1 pp 2v    do i=1,n
51  1 p 2v      a(i) = s / b(i)
52  1 p 2v      c(i) = s / d(i)
53  1          !ocl fission_point(1)
54  1 p 2v      e(i) = s / f(i)
55  1 p 2v      g(i) = s / h(i)
56  1 p 2v      enddo
57      end subroutine sub
58
Diagnostic messages: program name(sub)
jwd8212o-i "a.f90", line 50: ループを2分割しました。

```

ラージページ環境変数によるパディング

- XOS_MMM_L_FORCE_MMAP_THRESHOLDについて
- ラージページ環境変数によるパディング(改善前)
- ラージページ環境変数によるパディング(改善後)

環境変数名	指定値 (_付はdefault)	説明
XOS_MMM_L_FORCE_MMAP_THRESHOLD	0 1	<p>MALLOC_MMAP_THRESHOLD_(デフォルトは128MiB)以上のサイズのメモリ獲得時にmmap(2)を優先するかどうかの設定です。</p> <p>「0」の場合、mmap(2) は優先しません。まずヒープ領域の空きを検索し、空きがあればヒープ領域の空きメモリを返します。ヒープ領域の空きが見つからないときにのみ mmap(2) でメモリを獲得します。「1」の場合、mmap(2) を優先します。</p> <p>ヒープ領域の空きは検索せず、(例え空きがあっても)mmap(2) でメモリを獲得します。</p>

動的配列の各ストリームが16KB境界にあるためL1Dキャッシュスラッシングが発生します。そのため、浮動小数点ロードL2アクセス待ちが多くなっています。

改善前ソース

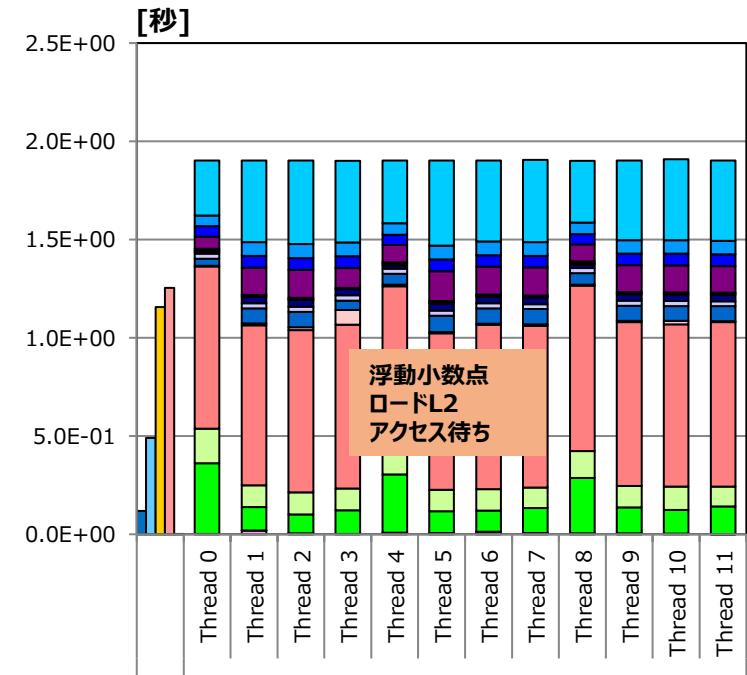
```

3      parameter(n=256,m=256)
4      real*8,allocatable ::a(:, :, b(:, :, c(:, :, d(:, :, e(:, :, 
5      f(:, :, g(:, :, h(:, :, 
6      allocate(a(n,m))
7      allocate(b(n,m))
8      allocate(c(n,m))
9      allocate(d(n,m))
10     allocate(e(n,m))
11     allocate(f(n,m))
12     allocate(g(n,m))
13     allocate(h(n,m))
14     ..... 
36    1 p      do j = 1 , m
17      <<< SIMD(VL: 8)
18      <<< SOFTWARE PIPELINE
19          MVE: 2, POL: S
20      <<< PREFETCH(HARD)
21      <<< c, b, d, e, f, g, h, a
22      <<< Loop-information End >>>
37    2 p v      do i = 1 , n
38    2 p v          a(i, j) = b(i, j) + c(i, j) + d(i, j) + e(i, j)
39          + f(i, j) + g(i, j) + h(i, j)
40    2 p v      enddo
41    1 p      enddo

```

配列サイズ
256×256×8B=
32×16KB
(16KB境界)

同一サイズのストリーム

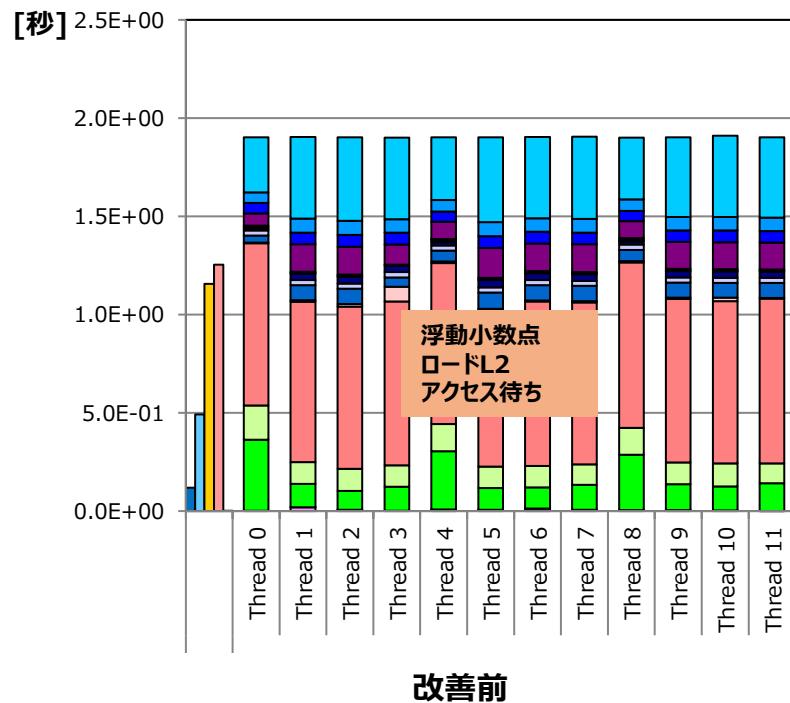


改善前

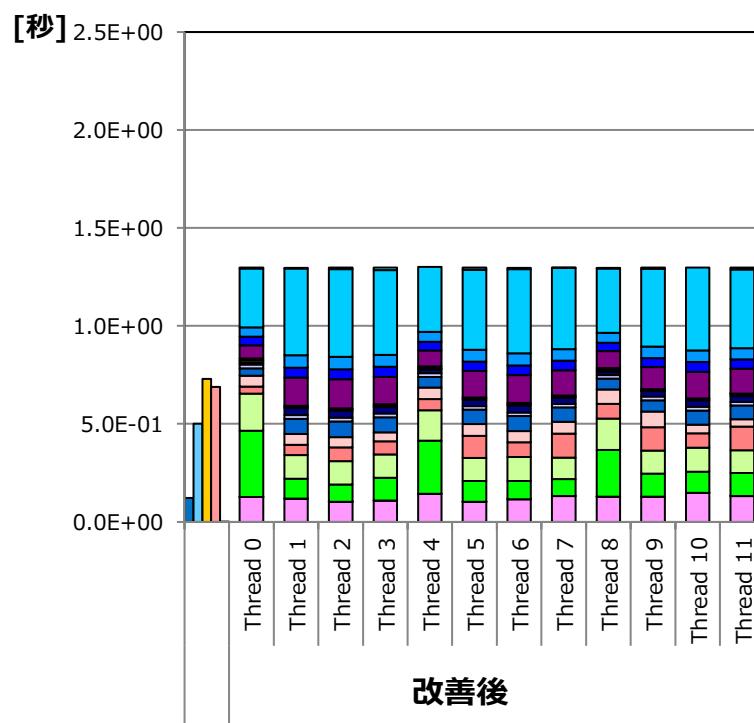
L1Dミスが高い
⇒ L1Dキャッシュスラッシングが発生している

Cache	L1I miss rate (/Effective instruction)	Load-store instruction	L1D miss	L1D miss rate (/Load-store instruction)	L1D miss demand rate (%) (/L1D miss)	L1D miss hardware prefetch rate (%) (/L1D miss)	L1D miss software prefetch rate (%) (/L1D miss)	L2 miss	L2 miss rate (/Load-store instruction)	L2 miss demand rate (%) (/L2 miss)	L2 miss hardware prefetch rate (%) (/L2 miss)	L2 miss software prefetch rate (%) (/L2 miss)
改善前	0.00	1.34.E+10	3.44.E+09	0.26	51.52%	48.28%	0.20%	1.61.E+03	0.00	73.10%	53.55%	0.00%

環境変数 **MALLOC_MMAP_THRESHOLD_ = 204800** を指定することで、各配列のアドレス配置が変わり、L1Dキャッシュラッシングを回避しました。その結果、浮動小数点ロードL2アクセス待ちが改善されました。



改善前



改善後

Cache

L1Dミス、L1Dミスdm率が改善した

	L1I miss rate (/Effective instruction)	Load-store instruction	L1D miss	L1D miss rate (/Load-store instruction)	L1D miss demand rate (%) (/L1D miss)	L1D miss hardware prefetch rate (%) (/L1D miss)	L1D miss software prefetch rate (%) (/L1D miss)	L2 miss	L2 miss rate (/Load-store instruction)	L2 miss demand rate (%) (/L2 miss)	L2 miss hardware prefetch rate (%) (/L2 miss)	L2 miss software prefetch rate (%) (/L2 miss)
改善前	0.00	1.34.E+10	3.44.E+09	0.26	51.52%	48.28%	0.20%	1.61.E+03	0.00	73.10%	53.55%	0.00%
改善後	0.00	1.35.E+10	1.81.E+09	0.13	9.47%	90.51%	0.02%	2.43.E+03	0.00	68.31%	58.01%	0.00%

動的配列の各ストリームが16KB境界にあるためL1Dキャッシュスラッシングが発生します。そのため、浮動小数点ロードL2アクセス待ちが多くなっています。

改善前ソース

```

62 void sub(int n, int m, double (* restrict a)[n], double (* restrict b)[n],
       double (* restrict c)[n], double (* restrict d)[n], double (* restrict e)[n],
       double (* restrict f)[n], double (* restrict g)[n], double (* restrict h)[n])
63 {
64     int i,j;
65     #pragma omp parallel for private(i)
66     <<< Loop-information Start >>>
67     <<< [OPTIMIZATION]
68     <<< PREFETCH(HARD) Expected by
69     <<< (unknown)
70     <<< Loop-information End >>>
71     p     for(j = 0; j<m; j++)
72     p     {
73         <<< Loop-information Start >>>
74         <<< [OPTIMIZATION]
75         <<< SIMD(VL: 8)
76         <<< SOFTWARE PIPELINING(IPC: 2.66, ITR: 104 MVE: 7 DOL: 6)
77         <<< PREFETCH(HARD) Expected by compiler :
78         <<< (unknown)
79         <<< Loop-information End >>>
80         p     v     for(i = 0; i<n; i++)
81         p     v     {
82             a[j][i] = b[j][i] + c[j][i] + d[j][i] + e[j][i] + f[j][i] + g[j][i] + h[j][i];
83         p     v     }
84         return;
85     }

```

配列宣言

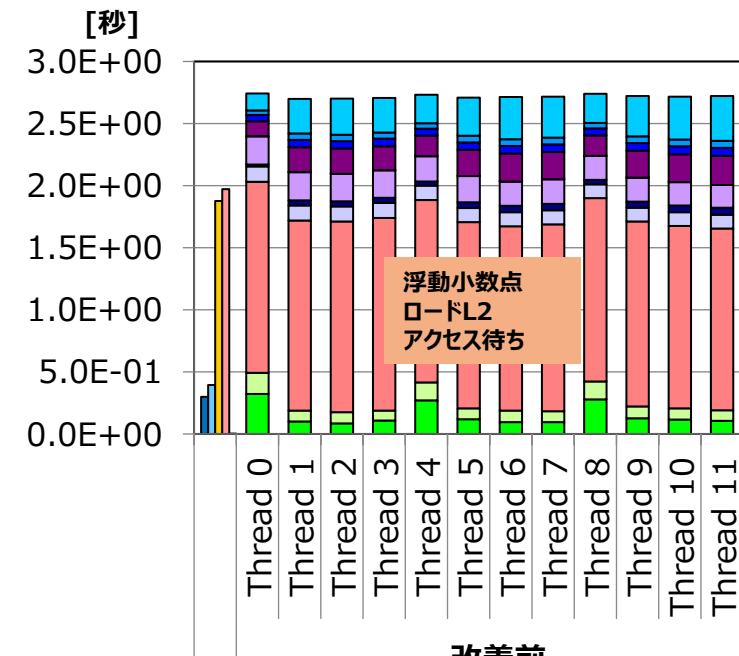
```

double a[256][256],
b[256][256], c[256][256],
d[256][256], e[256][256],
f[256][256], g[256][256],
h[256][256];

```

配列aのサイズ256×256×8B=
32×16KB(16KB境界)

同一サイズのストリーム

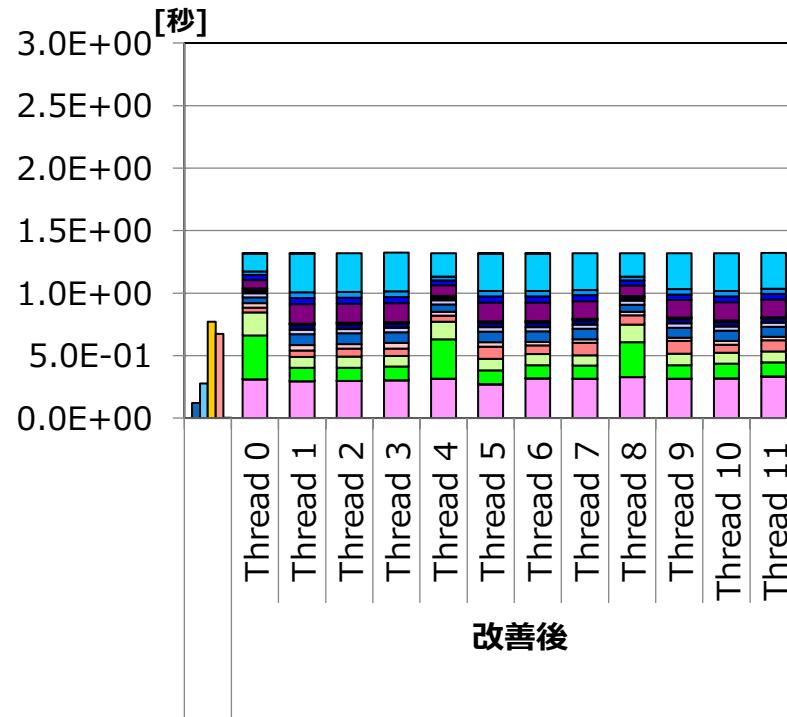
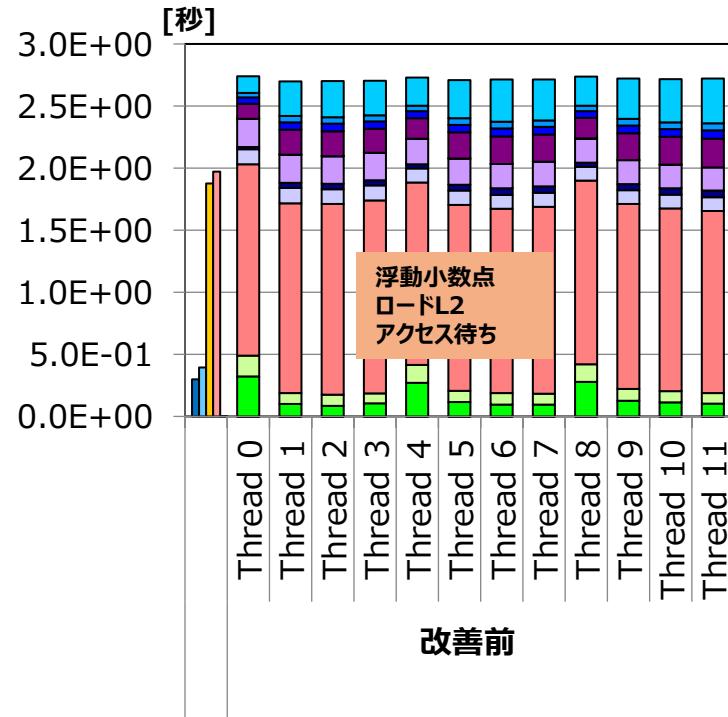


改善前

L1Dミスが高い
⇒ L1Dキャッシュスラッシングが発生している

Cache	L1I miss rate (/Effective instruction)	Load-store instruction	L1D miss	L1D miss rate (/Load-store instruction)	L1D miss demand rate (%) (/L1D miss)	L1D miss hardware prefetch rate (%) (/L1D miss)	L1D miss software prefetch rate (%) (/L1D miss)	L2 miss	L2 miss rate (/Load-store instruction)	L2 miss demand rate (%) (/L2 miss)	L2 miss hardware prefetch rate (%) (/L2 miss)	L2 miss software prefetch rate (%) (/L2 miss)
改善前	0.00	1.39E+10	5.24E+09	0.38	69.38%	30.62%	0.01%	7.49E+04	0.00	68.31%	44.20%	0.00%

環境変数 `MALLOC_MMAP_THRESHOLD_ = 204800` を指定することで、各配列のアドレス配置が変わり、L1Dキャッシュラッシングを回避しました。その結果、浮動小数点ロードL2アクセス待ちが改善されました。



L1Dミス、L1Dミスdm率が改善した

Cache	L1I miss rate (/Effective instruction)	Load-store instruction	L1D miss	L1D miss rate (/Load-store instruction)	L1D miss demand rate (%) (/L1D miss)	L1D miss hardware prefetch rate (%) (/L1D miss)	L1D miss software prefetch rate (%) (/L1D miss)	L2 miss	L2 miss rate (/Load-store instruction)	L2 miss demand rate (%) (/L2 miss)	L2 miss hardware prefetch rate (%) (/L2 miss)	L2 miss software prefetch rate (%) (/L2 miss)
改善前	0.00	1.39E+10	5.24E+09	0.38	69.38%	30.62%	0.01%	7.49E+04	0.00	68.31%	44.20%	0.00%
改善後	0.00	1.25E+10	1.80E+09	0.14	9.31%	90.69%	0.00%	2.98E+04	0.00	49.45%	57.41%	0.00%

演算待ち (SIMD化の促進)

- ループピーリング
- 定義関係が明らかでないループ
- ポインタ変数が含まれるループ

ループピーリング

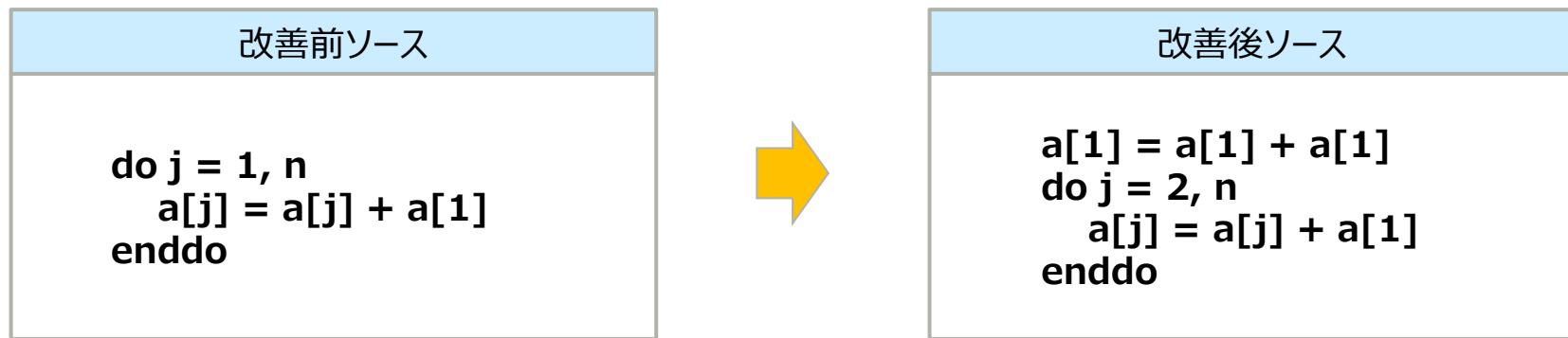
- ループピーリングとは
- ループピーリング（改善前）
- ループピーリング（ソースチューニング）

ループピーリングとは

ループピーリングとはループ内の一 部の処理をループから切り離す処理です。

ループ内でデータの依存関係がある場合、SIMD化や効果的なソフトウェアパイプラインングが行われない場合があります。

以下のような、ループ内に依存関係がある場合は、その部分だけループから切り離すループピーリングを行うことで対処することができます。



左の例では、 $j=1$ の時に定義された $a[1]$ が $2 \leq j \leq n$ で使用されており、部分的な依存関係があります。右の例のように $j=1$ の場合だけピーリングすることでループ内の依存関係を取り除くことができます。

配列aについて、 $i = 1$ の時に定義したものを $i = 2$ 以降で参照するデータ依存関係があるため、SIMD化が効果的に行われていません。

改善前ソース

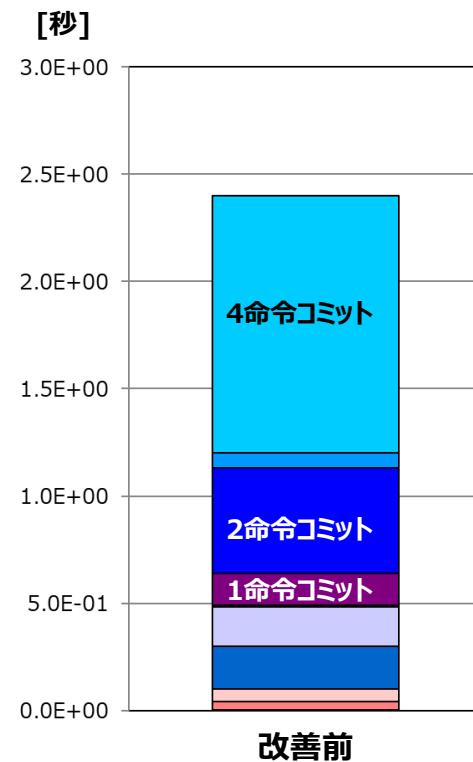
```

<<< Loop-information Start >>>
<<< [OPTIMIZATION]
<<< PREFETCH(HARD) Exp...
<<< b, (unknown)
<<< Loop-information End >...
8   2 p      do j = 1, m
<<< Loop-information Start >...
<<< [OPTIMIZATION]
<<< SIMD(VL: 8)
<<< SOFTWARE PIPELINING(IPC: 1.38, ITR: 144,
                           MVE: 5, POL: S)
<<< PREFETCH(HARD) Expected by compiler :
<<< a, b, (unknown)
<<< Loop-information End >>>
9   3 p 2m    do i = 1, n
10  3 p 2m      a(i,j) = c0 + a(1,j)*(c1 + b(i,j)*(c2 + b(i,j)*(c3 + b(i,j)*
11  3           & (c4 + b(I,j)*(c5 + b(I,j)*(c6 + b(I,j)*(c7 + b(I,j)*
12  3           & (c8 + b(i,j)*c9)))))))
13  3 p  v      end do
14  2 p        end do

```

一部しかSIMD化が行われていない

i=1の時にロード側配列aとストア側配列aで依存関係がある。



SIMD

	Effective instruction	SIMD instruction rate (%) (/Effective instruction)
改善前	1.49E+11	14.61%

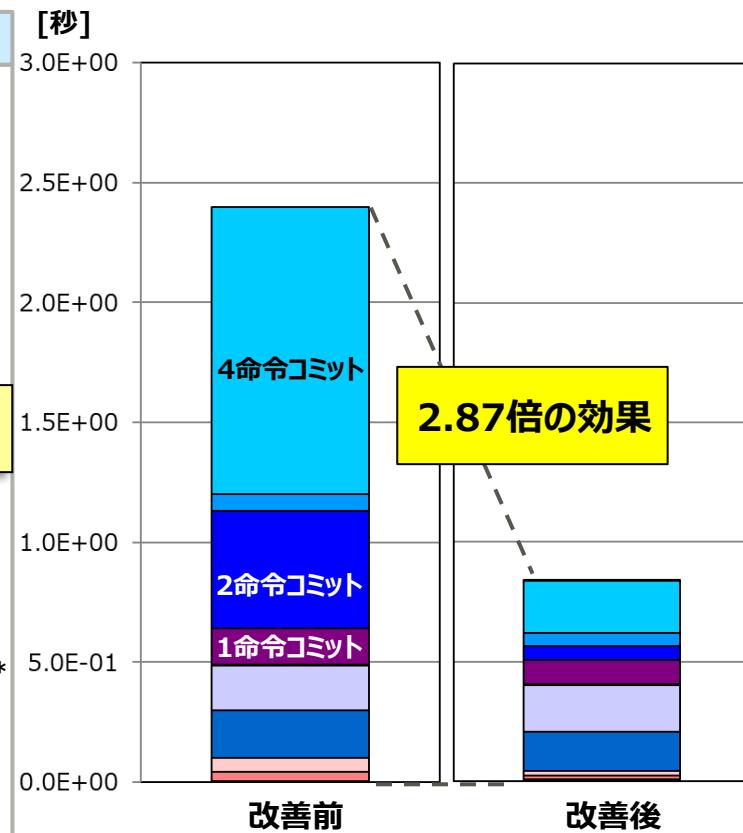
1回分のループピーリングを行うことで、SIMD化が促進されました。その結果、有効総命令数が削減され命令コミットが減少し性能が向上しました。

改善後ソース（ソースチューニング）

```

<<< Loop-information Start >>>
<<< [OPTIMIZATION]
<<< PREFETCH(HARD) Expected by compiler :
<<< b, a
<<< Loop-information End >>>
 8   2 p      do j = 1, m
 9   2 p      i = 1
10  2 p      a(i,j) = c0 + a(1,j)*(c1 + b(i,j)*(c2 + b(i,j)*(c3 + b(i,j)*
11  2          & (c4 + b(i,j)*(c5 + b(i,j)*(c6 + b(i,j)*(c7 + b(i,j)*
12  2          & (c8 + b(i,j)*c9)))))))
<<< Loop-information Start >>> 依存関係のある箇所を
<<< [OPTIMIZATION]          ピーリング
<<< SIMD(VL: 8)           依存関係がなくなったことで
<<< SOFTWARE PIPELINING(IPC: 1.23, ITR: 128,
                           MVE: 4, POL: S)
<<< PREFETCH(HARD) Expected by compiler :
<<< b, a
<<< Loop-information End >>>
13  3 p  2v      do i = 2, n
14  3 p  2v      a(i,j) = c0 + a(1,j)*(c1 + b(i,j)*(c2 + b(i,j)*(c3 + b(i,j)*
15  3          & (c4 + b(i,j)*(c5 + b(i,j)*(c6 + b(i,j)*(c7 + b(i,j)*
16  3          & (c8 + b(i,j)*c9)))))))
17  3 p  2v      end do
18  2 p          end do

```



SIMD

	Effective instruction	SIMD instruction rate (%) (/Effective instruction)
改善前	1.49E+11	14.61%
改善後	2.96E+10	81.26%

配列aについて、 $i = 0$ の時に定義したものを $i = 1$ 以降で参照するデータ依存関係があるため、SIMD化が効果的に行われていません。

改善前ソース

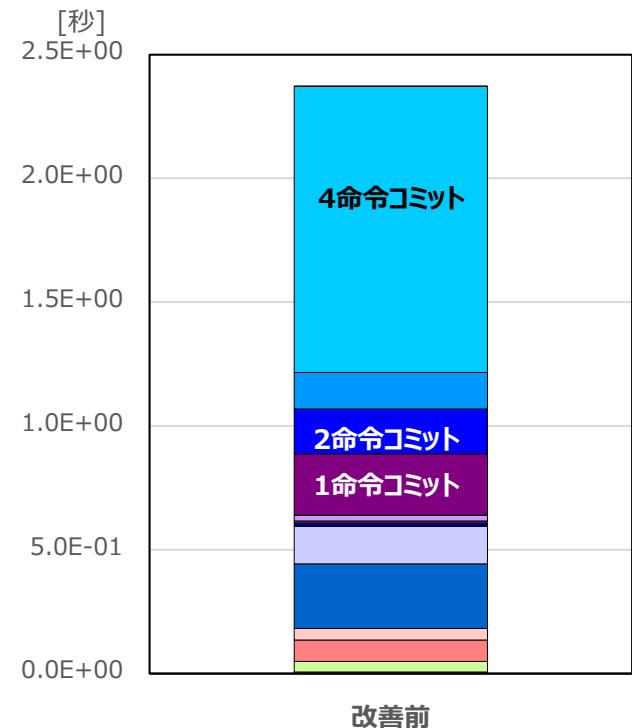
```

40     #pragma omp parallel
41     for(k=0;k<1000000;k++){
42         #pragma omp for nowait
43         <<< Loop-information Start >>>
44         <<< [OPTIMIZATION]
45         <<< PREFETCH(HARD) Expected by compiler :
46             (unknown)
47         <<< Loop-information End >>>
48         p   for(j=0;j<m;j++){
49             <<< Loop-information Start >>>
50             <<< [OPTIMIZATION]
51             <<< SIMD(VL: 8)
52             <<< SOFTWARE PIPELINING(IPC: 1.27, ITR: 144, MVE: 5, POL: S)
53             <<< PREFETCH(HARD) Expected by compiler :
54                 (unknown)
55             <<< Loop-information End >>>
56             p   2m    for(i=0;i<n;i++){
57                 p   2m    a[j][i] = c0 + a[j][0]*(c1 + b[j][i]*(c2 + b[j][i]*(c3 + b[j][i]*
58                     (c4 + b[j][i]*(c5 + b[j][i]*(c6 + b[j][i]*(c7 + b[j][i]*
59                         (c8 + b[j][i]*c9))))));
60             p   v   }
61             p   }
62         }

```

一部しかSIMD化が
行われていない

$i=0$ の時にロード側配列aとストア側配列aで依存関係がある。



Statistics	Effective instruction	SIMD instruction rate (%) (/Effective instruction)
改善前	1.33E+11	16.32%

1回分のループピーリングを行うことで、SIMD化が促進されました。その結果、有効総命令数が削減され命令コミットが減少し性能が向上しました。

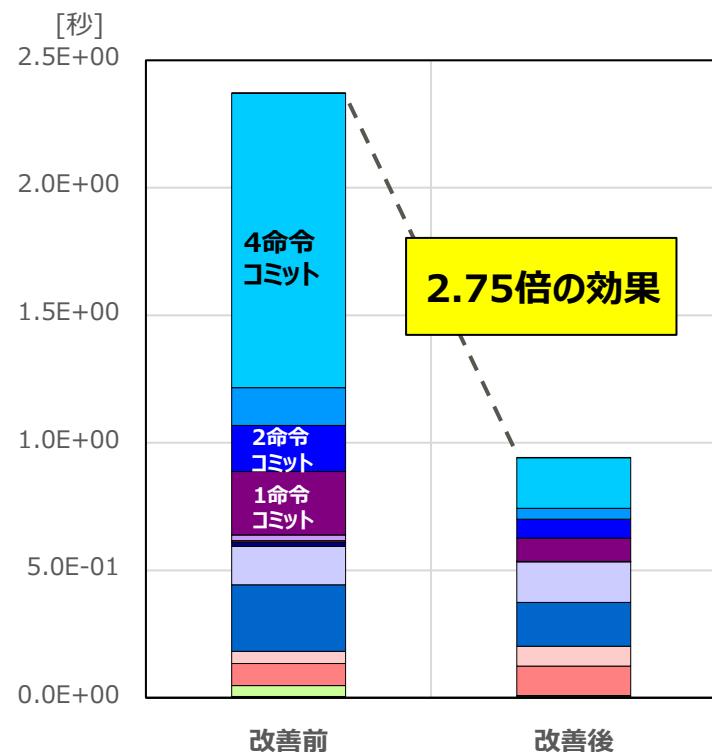
改善後ソース（ソースチューニング）

```

40     #pragma omp parallel
41     for(k=0;k<10000000;k++){
42         #pragma omp for nowait
43         <<< Loop-information Start >>>
44         <<< [OPTIMIZATION]
45         <<< PREFETCH(HARD) Expected by compiler :
46             (unknown)
47             <<< Loop-information End >>>
48         p     for(j=0;j<m;j++){
49             p         i=0;
50             p             a[j][i] = c0 + a[j][0]*(c1 + b[j][i]*(c2 + b[j][i]*(c3 + b[j][i]*
51                 (c4 + b[j][i]*(c5 + b[j][i]*(c6 + b[j][i]*(c7 + b[j][i]*
52                 (c8 + b[j][i]*c9)))))));
53         }
54     }
55 }
```

依存関係のある箇所を
ピーリング

依存関係がなくなったことで
SIMD化が促進された



Statistics	Effective instruction	SIMD instruction rate (%) (/Effective instruction)
改善前	1.33E+11	16.32%
改善後	2.79E+10	86.10%

定義関係が明らかでないループ[®]

- 定義関係が明らかでないループ（改善前）
- 定義関係が明らかでないループ（最適化制御行チューニング）
- 定義関係が明らかでないループ（最適化制御行）

以下の最適化制御行を指定します。

最適化指示子 (Fortran)	意味	指定可能な最適化制御行			
		プログ ラム 単位	DO ループ 単位	文単位	配列代 入文 単位
NORECURRENCE [(array1[,array2]...)]	DOループ [°] 内の演算対象となる配列の要素が、回転を跨いで定義引用されないことを本処理系に指示します。 (ループスライス可能な配列を指示します。) <i>array1</i> 、 <i>array2</i> 、…は配列名です。	○	○	×	○

最適化指示子 (C/C++)	意味	指定可能な最適化制御行			
		global 行	proce dure 行	loop行	state ment 行
norecurrence [(array1[,array2]...)]	ループ [°] の繰返しを跨いで定義引用されない配列を指示します。 (ループスライス可能な配列を指示します。) <i>array1</i> 、 <i>array2</i> 、…は配列名です。	○	○	×	○

配列aについて、データ依存関係が不明なため、SIMD化や効果的なソフトウェアパイプラインングが行われていません。そのため、整数演算待ちが多くなっています。

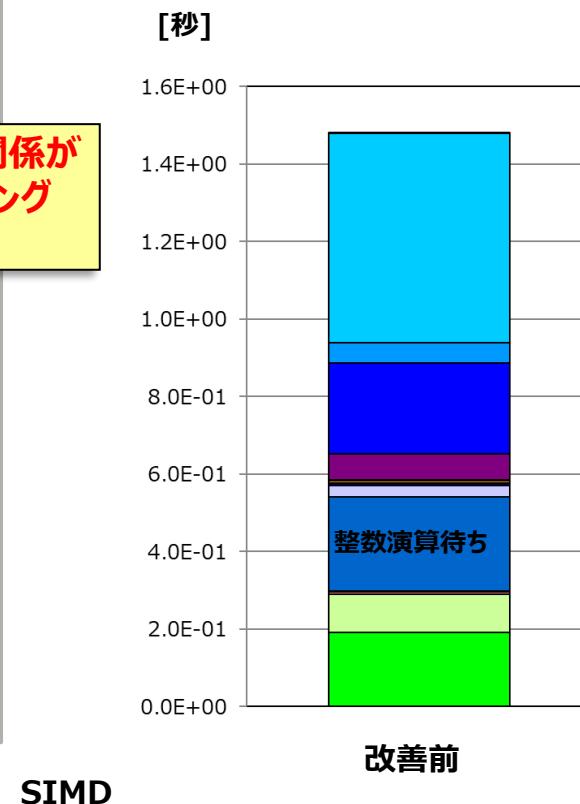
	改善前ソース
--	--------

```

53   1 pp      <<< Loop-information Start >>>
      <<< [PARALLELIZATION]
      <<< 配列aについてイタレーション間のデータ依存関係が
      <<< 不明のため、効果的なソフトウェアパイプラインング
      <<< が行われていない。
      <<< Loop-information End >>>
          do j=1,n1
<<< Loop-information Start >>>
<<< [OPTIMIZATION]
<<< SOFTWARE PIPELINING(IPC: 0.56, ITR: 3,
                           MVE: 2, POL: S)
<<< PREFETCH(HARD) Expected by compiler :
<<<   l, b, x
<<< Loop-information End >>>
      s           do i=1,n2
      2 p         a(l(i),j)=a(x(i),j)/b(i,j)
      2 p         end do
      2 p         v
      1 p         end do

```

ロード側の配列aとストア側の配列aの依存関係が不明。



	Effective instruction	SIMD instruction rate (%) (/Effective instruction)
改善前	7.10E+10	0.00%

NORECURRENCE指示子によりデータ依存関係がないことを明示化するので、SIMD化およびソフトウェアパイプラインングが促進されました。その結果、整数演算待ちが大幅に改善されました。

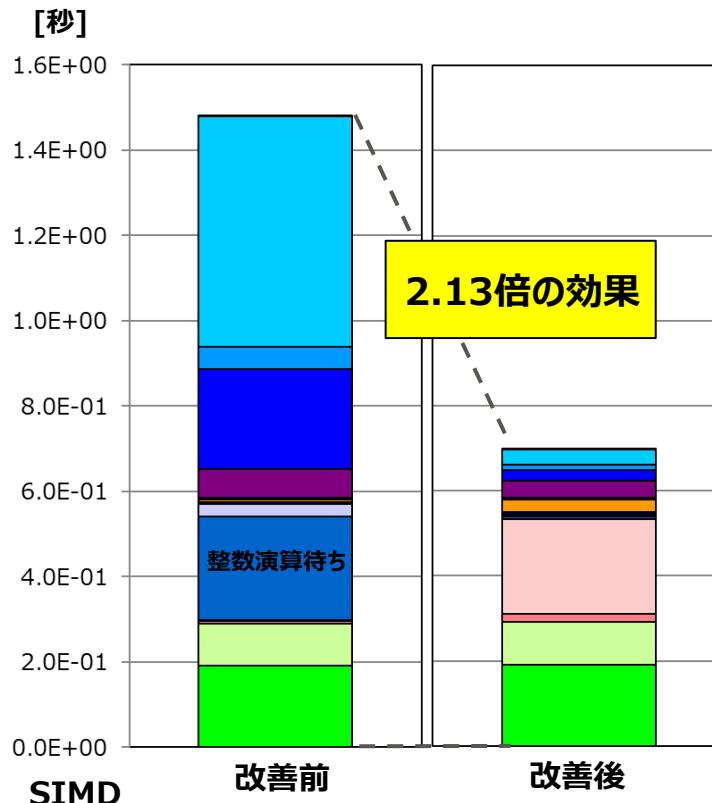
改善後ソース（最適化制御行チューニング）

```

53      !ocl norecurrence(a)
      <<< Loop-information
      <<< [PARALLELIZATION]
      <<< Standard iteration
      <<< [OPTIMIZATION]
      <<< PREFETCH(HARD) Expected by compiler :
      <<< x, b, l
      <<< Loop-information End >>>
54 1 pp          do j=1,n1
      <<< Loop-information Start >>>
      <<< [OPTIMIZATION]
      <<< SIMD(VL: 16)
      <<< SOFTWARE PIPELINING(IPC: 2.72, ITR: 288,
                                MVE: 3, POL: S)
      <<< PREFETCH(HARD) Expected by compiler :
      <<< x, b, l
      <<< Loop-information End >>>
55 2 p 2v          do i=1,n2
56 2 p 2v          a(l(i),j)=a(x(i),j)/b(i,j)
57 2 p 2v          end do
58 1 p          end do

```

SIMD化や効果的なソフトウェアパイプラインングが行われた



配列aについて、データ依存関係が不明なため、SIMD化や効果的なソフトウェアパイプラインニングが行われていません。そのため、整数演算待ちが多くなっています。

改善前ソース

```

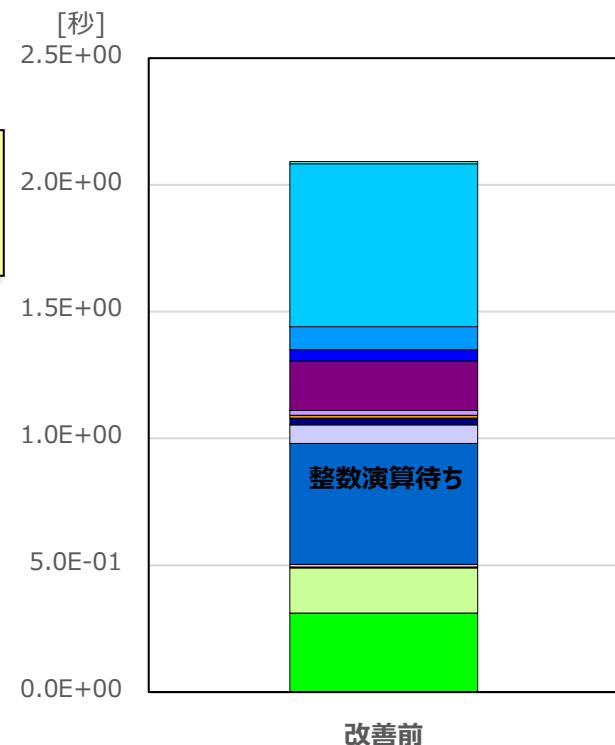
48     #pragma omp parallel
49     {
50         #pragma omp for nowait
51         <<< Loop-information Start >>>
52         <<< [OPTIMIZATION]
53         <<< SOFTWARE PIPELINING(IPC: 0.39, ITR: 6,
54             MVE: 2, POL: S)
55         <<< PREFETCH(HARD) Expected by compiler :
56             (unknown)
57         <<< Loop-information End >>>
58         p 2s   for(i=0; i<n2; i++)
59         p 2v   {
60             p 2m   a[j][l[i]] = a[j][x[i]] / b[j][i];
61             p 2v   }
62         }
63     }
64     return;
65 }

```

配列aについてイタレーション間のデータ依存性が不明のため、効果的なソフトウェアパイプラインングが行われていない。

ロード側の配列
ストア側の配列
依存関係が不明

ロード側の配列aと
ストア側の配列aの
依存関係が不明。



Statistics	Effective instruction	SIMD instruction rate (%) (/Effective instruction)
改善前	8.49E+10	0.00%

norecurrence指示子によりデータ依存関係がないことを明示化するので、SIMD化およびソフトウェアパイプラインングが促進されました。その結果、整数演算待ちが大幅に改善されました。

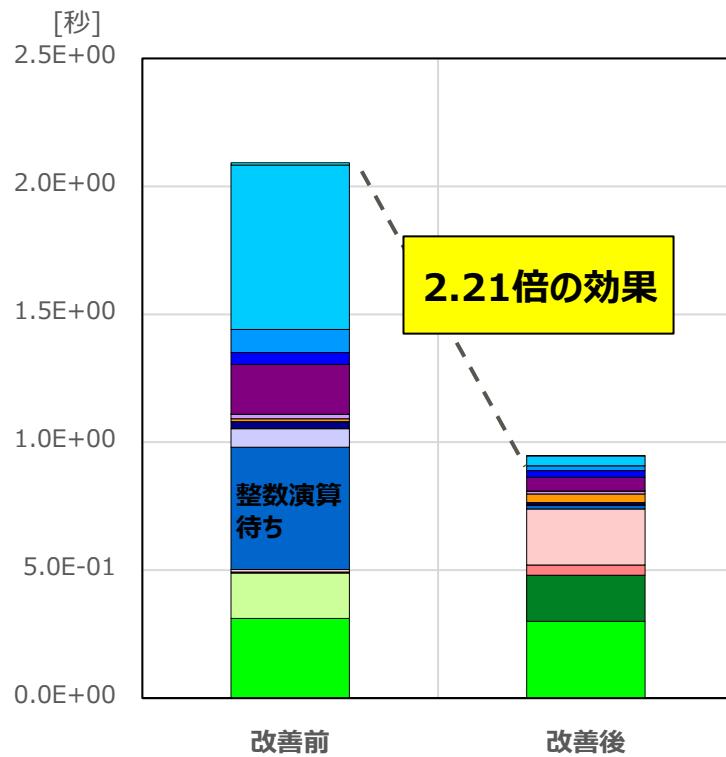
改善後ソース（最適化制御行チューニング）

```

48 #pragma omp parallel
49 {
50     #pragma omp for nowait
51     #pragma loop norecurrence
52     <<< Loop-information Start >>>
53     <<< [OPTIMIZATION]
54     <<< PREFETCH(HARD)
55     <<< (unknown)
56     <<< Loop-information
57     for(j=1; j<n1; j+
58     {
59         <<< Loop-information Start >>>
60         <<< [OPTIMIZATION]
61         <<< SIMD(VL: 16)
62         <<< SOFTWARE PIPELINING(IPC: 2.27, ITR: 256,
63                                     MVE: 3, POL: S)
64         <<< PREFETCH(HARD) Expected by compiler :
65         <<< (unknown)
66         <<< Loop-information End >>>
67         2v         for(i=0; i<n2; i++)
68         2v         {
69             2v             a[j][l[i]]=a[j][x[i]]/b[j][i];
70             2v         }
71     }
72     return
    }
```

**a[j, l[i]]とa[j, x[i]] に
データ依存関係がないこと
をコンパイラに知らせる**

**SIMD化や効果的なソフトウェア
パイプラインングが行われた**



Statistics	Effective instruction rate (%) (/Effective instruction)
改善前	8.49E+10
改善後	2.74E+10

ポインタ変数が含まれるループ

- ポインタ変数が含まれるループとは
- ポインタ変数が含まれるループ（改善前）
- ポインタ変数が含まれるループ（最適化制御行チューニング）
- ポインタ変数が含まれるループ（contiguous属性指示）

ポインタ変数が含まれるループとは

FUJITSU

ポインタ変数が記憶領域のどこを占めるかは実行時に決まるため、ポインタ変数を含むループでは最適化が促進されない場合があります。

ソースコード例

```
real,dimension(100),target::x  
real,dimension(:),pointer::a,b  
a=>x(1:10)  
b=>x(11:20)  
do i=1,10000  
    a(i)=2.0/b(i)+1.0  
end do
```



ソースコード例

```
real,dimension(100),target::x  
real,dimension(:),pointer::a,b  
!ocl noalias  
a=>x(1:10)  
b=>x(11:20)  
do i=1,10000  
    a(i)=2.0/b(i)+1.0  
end do
```

NOALIAS指示子を指定することで、異なるポインタ変数が同一の記憶領域を指さないことを翻訳時に判断することができます。これによりポインタ変数に関する最適化が促進されます。

ただし、ポインタ変数の結合状態がループ内で変更される場合、最適化指示子を指定しても最適化が促進されないことがあります。

最適化指示子 (Fortran)	意味	指定可能な最適化制御行			
		プログラム単位	DOループ単位	文単位	配列代入文単位
NOALIAS	ポインタ変数が他の変数と記憶域を共有しないことを指示します。	○	○	×	○

最適化指示子 (C/C++)	意味	指定可能な最適化制御行			
		global行	procedure行	loop行	statement行
noalias	ポインタ変数が他の変数と記憶域を共有しないことを指示します。	○	○	○	×

以下の翻訳オプションを指定することで、最適化制御行チューニングと同等の効果を得ることができます。

翻訳オプション	機能説明
-Knoalias [=spec]	<p>ポインタ変数またはポインタ成分が他の変数と記憶域を結合しないものとして、最適化を行うことを指示します。specには、sを指定することができます。Sを指定した場合、コンパイラはFortran規格のポインタ属性を持つ実体が、他の変数を結合していないものとして最適化します。以下の文脈では、ポインタ属性をもつ実体が他の変数と結合する可能性があります。</p> <ul style="list-style-type: none">- ポインタ代入文- ポインタ成分をもつ派生型の代入文- ポインタ成分をもつ派生型がSOURCE=指定子に現れるALLOCATE文- ポインタ属性またはポインタ成分をもつ仮引数- ポインタ属性またはポインタ成分をもつ変数への初期設定

■ 使用例（改善前ソース時）

■ \$ frtpx -Kfast,parallel sample.f90 **-Knoalias**

ポインタ変数a、bが異なる記憶領域を指すことが不明なためSIMD化が促進されません。そのため、整数演算待ちが多くなっています。

改善前ソース

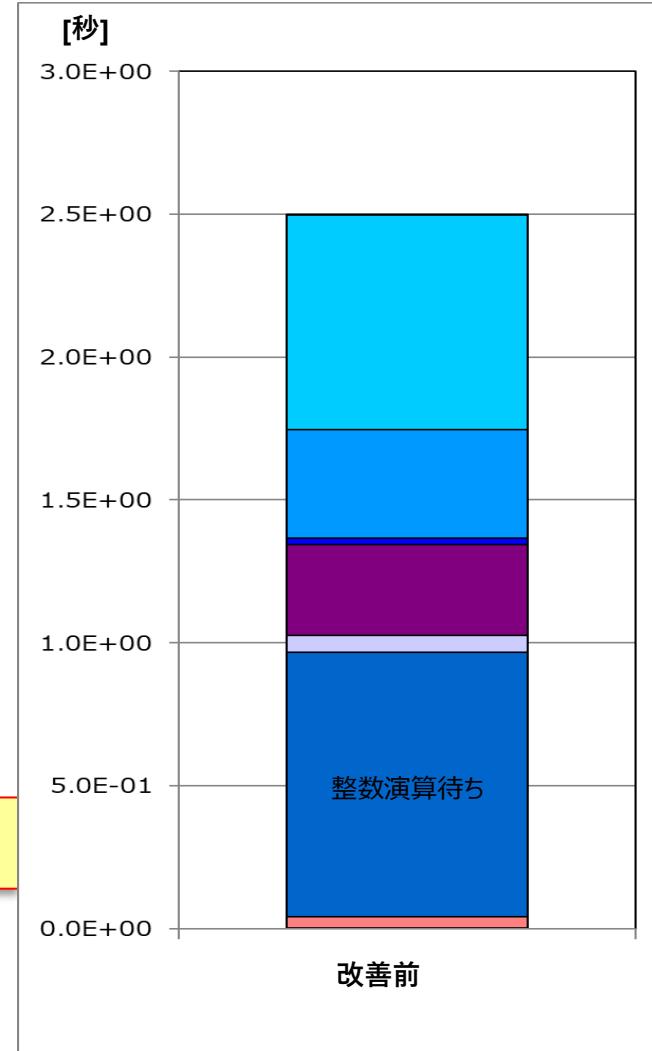
```

3      real,dimension(100000),target::x
4      integer :: kmax
5      real,dimension(:),pointer::a,b
:
9      a=>x(1:10000)
10     b=>x(10001:20000)
11     kmax = 1000000
12 !$omp parallel
13   1     do k=1,kmax
14   1       !$omp do
<<< Loop-information Start >>>
<<< [OPTIMIZATION]
<<< SOFTWARE PIPELINING(IPC: 0.19, ITR: 8, MVE: 2, POL: L)
<<< Loop-information End >>>
15   2 p 2s     do i=1,10000
16   2 p 2s       a(i)=2.0/b(i)+1.0
17   2 p 2s       end do
18   1           !$omp enddo nowait
19   1           end do
20           !$omp end parallel

```

SIMD化されていない

Statistics	Effective instruction	SIMD instruction rate (%) (/Effective instruction)
改善前	1.10E+11	0.00%



NOALIAS指示子によりデータ依存関係がないことを明示化することで、SIMD化およびソフトウェアパイプラインングが促進されました。その結果、整数演算待ちなどが大幅に改善されました。

改善後ソース（最適化制御行チューニング）

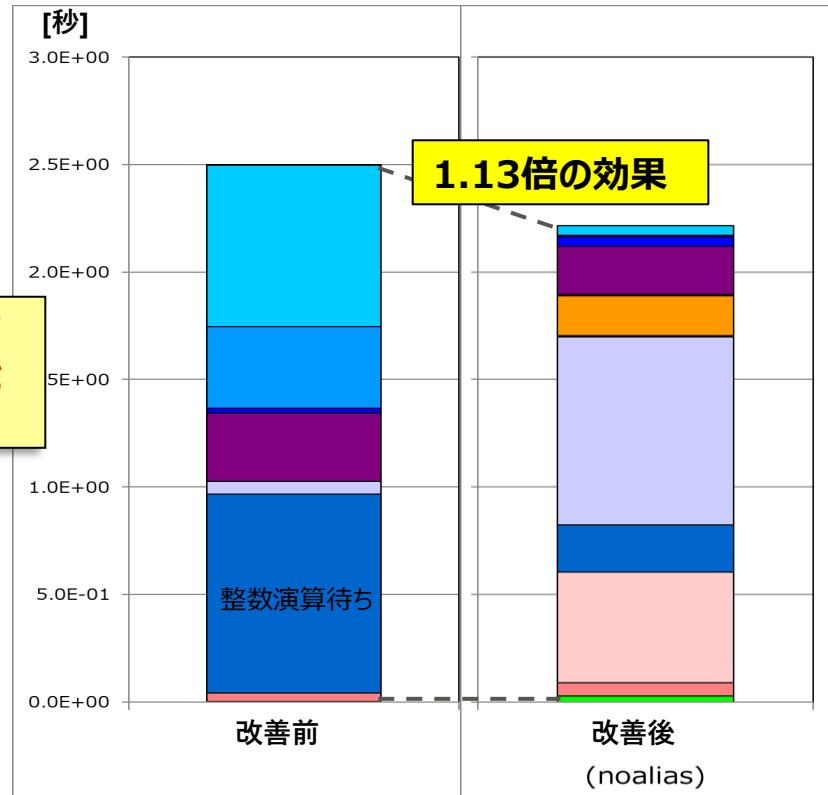
```

3      real,dimension(100000),target::x
4      integer :: kmax
5      real,dimension(:),pointer::a,b
6
7      a=>x(1:10000)
8      b=>x(10001:20000)
9      kmax = 1000000
10
11     !$omp parallel
12     do k=1,kmax
13     !$omp do
14     !ocl noalias
15
<<< Loop-information Start >>>
<<< [OPTIMIZATION]
<<< SIMD(VL: 16)
<<< SOFTWARE PIPELINING(IPC: 0.19, ITR: 96,
                           MVE: 2, POL: L)
<<< Loop-information End >>>
16    2 p 2v      do i=1,10000
17    2 p 2v          a(i)=2.0/b(i)+1.0
18    2 p 2v          end do
19    1           !$omp enddo nowait
20    1           end do
21    !$omp end parallel

```

a(i) と b (i) にデータ
依存関係がないことを
コンパイラに知らせる

SIMD化により有効総命令数が
削減された



Statistics	Effective instruction	SIMD instruction rate (%) (/Effective instruction)
改善前	1.10E+11	0.00%
改善後	1.26E+10	40.83%

contiguous属性の指示によりデータの連續性を明示化することで、**非連續命令が連續命令になりました**。その結果、最適化が促進され、浮動小数点ロードL1Dアクセス待ちや浮動小数点演算待ちなどが大幅に改善されました。

改善後ソース (contiguous属性指示)

```

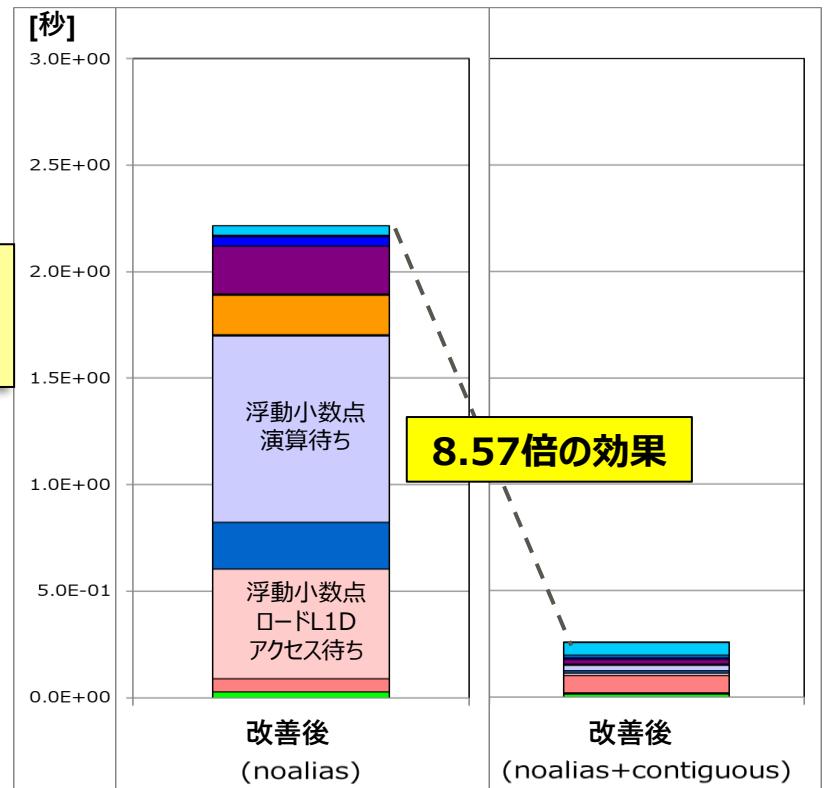
3      real,dimension(100000),target::x
4      integer :: kmax
5      real,dimension(:),pointer,contiguous::a,b
:
9      a=>x(1:10000)
10     b=>x(10001:20000)
11     kmax = 1000000
12     !$omp parallel
<<< Loop-information Start >>>
<<< [OPTIMIZATION]
<<< PREFETCH(HARD) Expected by compiler :
<<< (unknown)
<<< Loop-information End
13    1   do k=1,kmax
14    1     !$omp do
15    1       !ocl noalias
<<< Loop-information Start >>>
<<< [OPTIMIZATION]
<<< SIMD(VL: 16)
<<< SOFTWARE PIPELINING(IPC: 2.62, ITR: 352,
                           MVE: 3, POL: S)
<<< PREFETCH(HARD) Expected by compiler :
<<< (unknown)
<<< Loop-information End >>>
16    2   p 2v      do i=1,10000
17    2   p 2v      a(i)=2.0/b(i)+1.0
18    2   p 2v      end do
19    1     !$omp enddo nowait
20    1       end do
21      !$omp end parallel

```

ポイントア配列a、bがそれぞれ連続領域であることをコンパイラに知らせる

a(i)とb(i)にデータ依存関係がないことをコンパイラに知らせる

非連續命令が連続命令になった



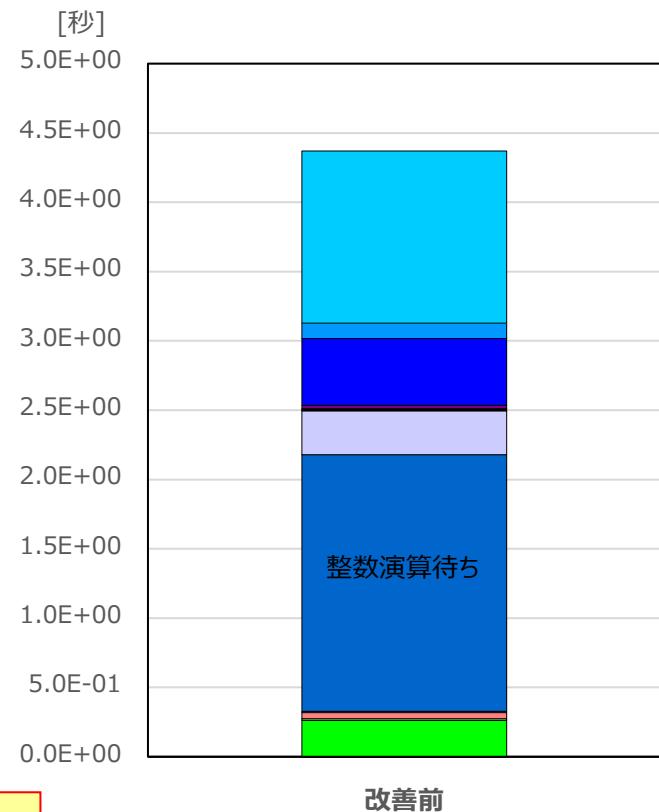
Instruction	Load-store instruction					
	Load instruction		Store instruction			
	SIMD	Non-SIMD	Non-SIMD	SIMD	Non-SIMD	Non-SIMD
Single vector contiguous load instruction	1.10E+01	6.36E+08	1.08E+09	1.56E+02	6.36E+08	1.33E+08
改善後 (noalias)	1.10E+01	6.36E+08	1.08E+09	1.56E+02	6.36E+08	1.33E+08
改善後 (noalias+contiguous)	6.36E+08	0.00E+00	5.00E+08	6.36E+08	0.00E+00	3.70E+07

ポインタ変数a、bが異なる記憶領域を指すことが不明なためSIMD化が促進されません。
そのため、整数演算待ちが多くなっています。

改善前ソース

```

17      #pragma omp parallel
18      {
19          <<< Loop-information Start >>>
20          <<< [OPTIMIZATION]
21          <<< PREFETCH(HARD) Expected by compiler :
22          <<< (unknown)
23          <<< Loop-information End >>>
24          for(k=0; k<kmax; k++)
25          {
26              #pragma omp for nowait
27              <<< Loop-information Start >>>
28              <<< [OPTIMIZATION]
29              <<< SOFTWARE PIPELINING(IPC: 0.21, ITR: 8, MVE: 2, POL: L)
30              <<< PREFETCH(HARD) Expected by compiler :
31              <<< (unknown)
32              <<< Loop-information End >>>
33      p  2s    for(i=0; i<10000; i++)
34      p  2s    {
35      p  2s    a[i]=2.0/b[i]+1.0;
36      p  2s    }
37      }
```



SIMD化されていない

Statistics	Effective instruction	SIMD instruction rate (%) (/Effective instruction)
改善前	1.50E+11	0.00%

NOALIAS指示子によりデータ依存関係がないことを明示化することで、SIMD化およびソフトウェアパイプラインングが促進されました。その結果、整数演算待ちなどが大幅に改善されました。

改善後ソース（最適化制御行チューニング）

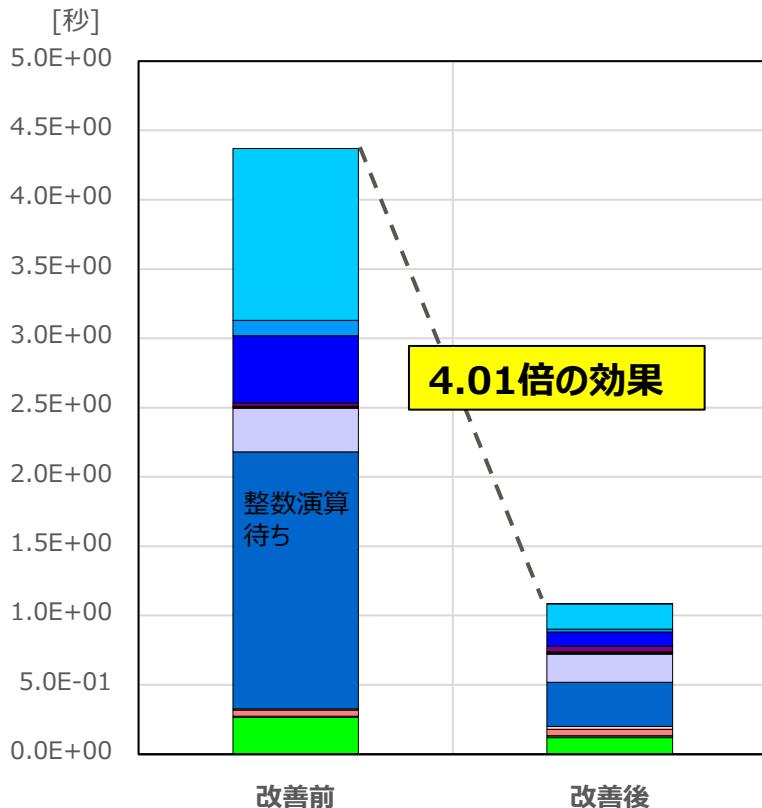
```

17     #pragma omp parallel
18     {
19         <<< Loop-information Start >>>
20         <<< [OPTIMIZATION]
21         <<< PREFETCH(HARD) Expected by compiler :
22             (unknown)
23             <<< Loop-information
24                 for(k=0; k<kmax;
25                     {
26                         #pragma omp for nowait
27                         #pragma loop noalias
28                         <<< Loop-information Start >>>
29                         <<< [OPTIMIZATION]
30                         <<< SIMD(VL: 8)
31                         <<< SOFTWARE PIPELINING(IPC: 0.18, ITR: 64,
32                                         MVE: 2, POL: L)
33                         <<< PREFETCH(HARD) Expected by compiler :
34                             (unknown)
35                         <<< Loop-information End >>>
36                         p 2v   for(i=0; i<10000; i++)
37                         p 2v   {
38                         p 2v       a[i]=2.0/b[i]+1.0;
39                         p 2v   }
40                     }
41                 }
42             }
43         }
44     }

```

a[i] と b[i] にデータ依存関係がないことをコンパイラに知らせる

SIMD化により有効総命令数が削減された



Statistics	Effective instruction	SIMD instruction rate (%) (/Effective instruction)
改善前	1.50E+11	0.00%
改善後	2.27E+10	72.14%

演算待ち（レイテンシの隠蔽）

- ・ ループ分割（ソフトウェアパイプラインング促進）
- ・ 適切なアンローリング展開数の指定およびソフトウェアパイプラインングの抑止
- ・ ストライピング（インターリーブ）展開数の指定および
ソフトウェアパイプラインングの抑止
- ・ 外側ループでのソフトウェアパイプラインング
- ・ リローリング
- ・ ループアンスイッチング

ループ分割 (ソフトウェアパイプラインング促進)

- ・ ループ分割（ソフトウェアパイプラインング促進）（改善前）
- ・ ループ分割（ソフトウェアパイプラインング促進）（ソースチューニング）

演算の連鎖が長く、多くのレジスタを必要とします。

そのため、SWPLなどのスケジューリングの最適化ができず、浮動小数点演算待ちが多くなっています。

改善前ソース

```

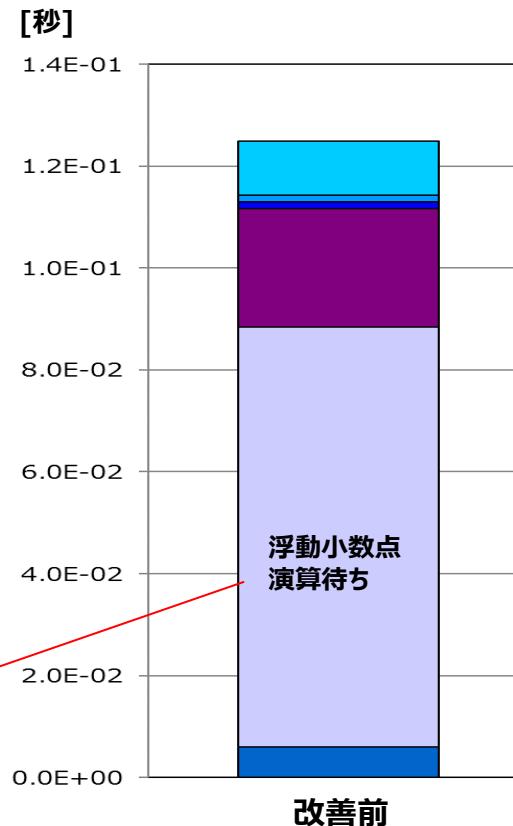
<<< Loop-information Start >>>
<<< [OPTIMIZATION]
<<< SIMD(VL: 8)
<<< PREFETCH(HARD) Expected by compiler :
<<<   x1, y
<<< Loop-information End >>>
18  2    2v      do i = 1, n
19  2    2v      y(i) = c0 / x1(i) / c1 / x1(i) &
20  2          / c2 / x1(i) / c3 / x1(i) / c4 / x1(i)
21  2    2v      end do

```

※コンパイルオプション
-Knoevalを指定

演算連鎖が長く、多くのレジスタを使用するため、
ソフトウェアパイプラインングを適用できない

スケジューリングが最適化できていないため、
演算待ちが大きくみえる



Cache

	L1I miss rate (/Effective instruction)	Load-store instruction	L1D miss	L1D miss rate (/Load-store instruction)	L1D miss demand rate (%) (/L1D miss)	L1D miss hardware prefetch rate (%) (/L1D miss)	L1D miss software prefetch rate (%) (/L1D miss)	L2 miss	L2 miss rate (/Load-store instruction)	L2 miss demand rate (%) (/L2 miss)	L2 miss hardware prefetch rate (%) (/L2 miss)	L2 miss software prefetch rate (%) (/L2 miss)
改善前	0.00	7.69E+06	1.61E+03	0.00	71.62%	27.38%	1.00%	1.32E+02	0.00	48.92%	59.71%	0.00

ループ分割を行うことで演算の連鎖が短くなり、1ループに使用するレジスタが減少しました。
その結果、SWPLなどのスケジューリングの最適化が行われ、演算待ちが改善されました。

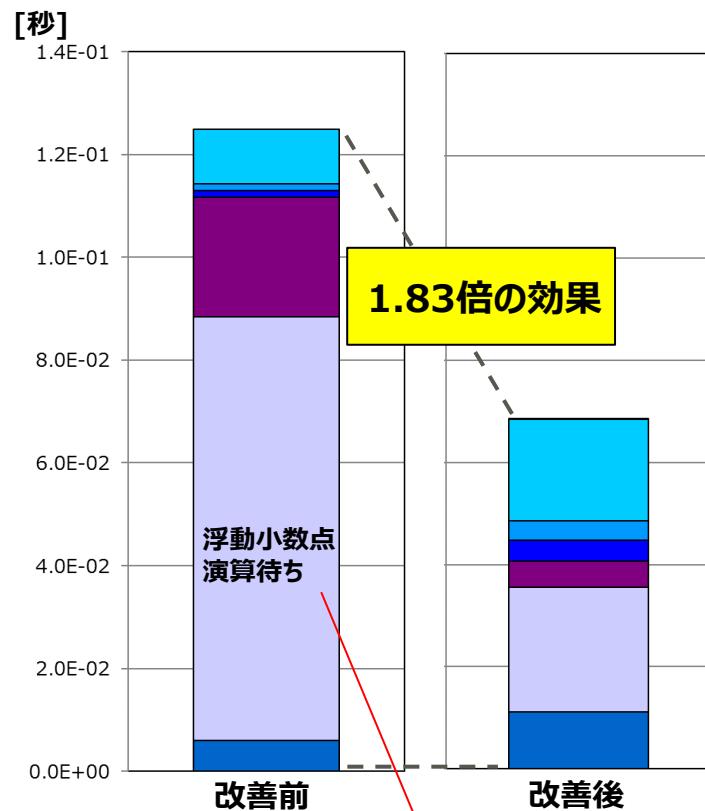
改善後ソース

```

20   1      !ocl loop_nofusion
20   1      <<< Loop-information Start >>>
20   1      <<< [OPTIMIZATION]
20   1      <<< SIMD(VL: 8)
20   1      <<< SOFTWARE PIPELINING(IPC: 1.02, ...
20   1      <<< PREFETCH(HARD) Expected by compiler :
20   1          y, x1
20   1      <<< SPILLS :
20   1          GENERAL : SPILL 0 FILL 0
20   1          SIMD&FP : SPILL 0 FILL 0
20   1          SCALABLE : SPILL 0 FILL 27
20   1          PREDICATE : SPILL 0 FILL 0
20   1      <<< Loop-information End >>>
21   2      2v      do i = 1, n
22   2      2v      y(i) = c2 / x1(i) / c3 / x1(i) / c4 / x1(i)
23   2      2v      end do
23   2      <<< Loop-information Start >>>
23   2      <<< [OPTIMIZATION]
23   2      <<< SIMD(VL: 8)
23   2      <<< SOFTWARE PIPELINING(IPC: 1.08, ...
23   2      <<< PREFETCH(HARD) Expected by compiler :
23   2          x1, y
23   2      <<< Loop-information End >>>
24   2      2v      do i = 1, n
25   2      2v      y(i) = c0 / x1(i) / c1 / x1(i) / y(i)
26   2      2v      end do

```

ループ融合を抑止
ループ分割



演算待ちが減少した

※コンパイルオプション-Knoevalを指定

	L1I miss rate (/Effective instruction)	Load-store instruction	L1D miss	L1D miss rate (/Load-store instruction)	L1D miss demand rate (%) (/L1D miss)	L1D miss hardware prefetch rate (%) (/L1D miss)	L1D miss software prefetch rate (%) (/L1D miss)
改善前	0.00	7.69E+06	1.61E+03	0.00	71.62%	27.38%	1.00%
改善後	0.00	3.65E+07	1.73E+03	0.00	70.79%	28.98%	0.23%

演算の連鎖が長く、多くのレジスタを必要とします。

そのため、SWPLなどのスケジューリングの最適化ができず、浮動小数点演算待ちが多くなっています。

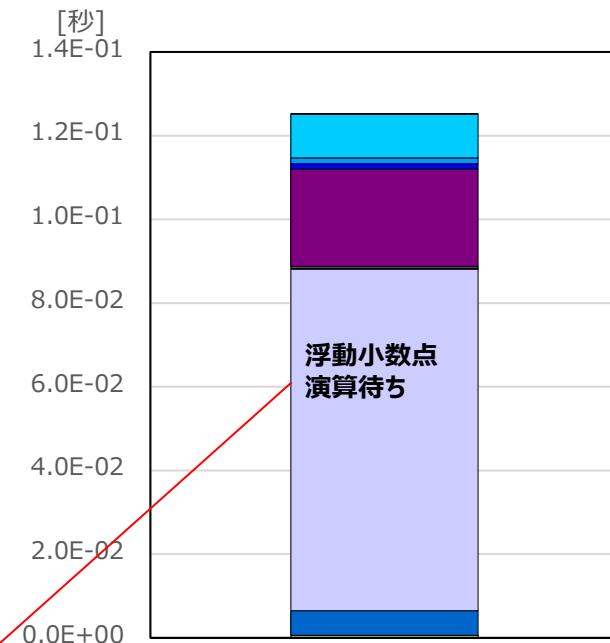
改善前ソース

```

18   for (iter = 0; iter < 10000; iter++) {
<<< Loop-information Start >>>
<<< [OPTIMIZATION]
<<< SIMD(VL: 8)
<<< PREFETCH(HARD) Expected by compiler :
<<< (unknown)
<<< Loop-information End >>>
19   2v   for (i = 0; i < n; i++) {
20     2v   y[i] = c0 / x1[i] / c1 / x1[i] / c2 / x1[i] / c3 / x1[i] / c4 / x1[i];
21   2v   }
22 }
```

※コンパイルオプション
-Knoevalを指定

演算連鎖が長く、多くのレジスタを使用するため、
ソフトウェアパイプラインングを適用できない



改善前

スケジューリングが最適化できていないため、
演算待ちが大きくみえる

Cache	L1I miss rate (/Effective instruction)	Load-store instruction	L1D miss	L1D miss rate (/Load-store instruction)	L1D miss demand rate (%) (/L1D miss)	L1D miss hardware prefetch rate (%) (/L1D miss)	L1D miss software prefetch rate (%) (/L1D miss)	L2 miss	L2 miss rate (/Load-store instruction)	L2 miss demand rate (%) (/L2 miss)	L2 miss hardware prefetch rate (%) (/L2 miss)	L2 miss software prefetch rate (%) (/L2 miss)
改善前	0.00	1.11E+09	1.72E+06	0.00	99.92%	0.08%	-0.01%	4.01E+03	0.00	86.58%	26.24%	0.00%

ループ分割を行うことで演算の連鎖が短くなり、1ループに使用するレジスタが減少しました。
その結果、SWPLなどのスケジューリングの最適化が行われ、演算待ちが改善されました。

改善後ソース

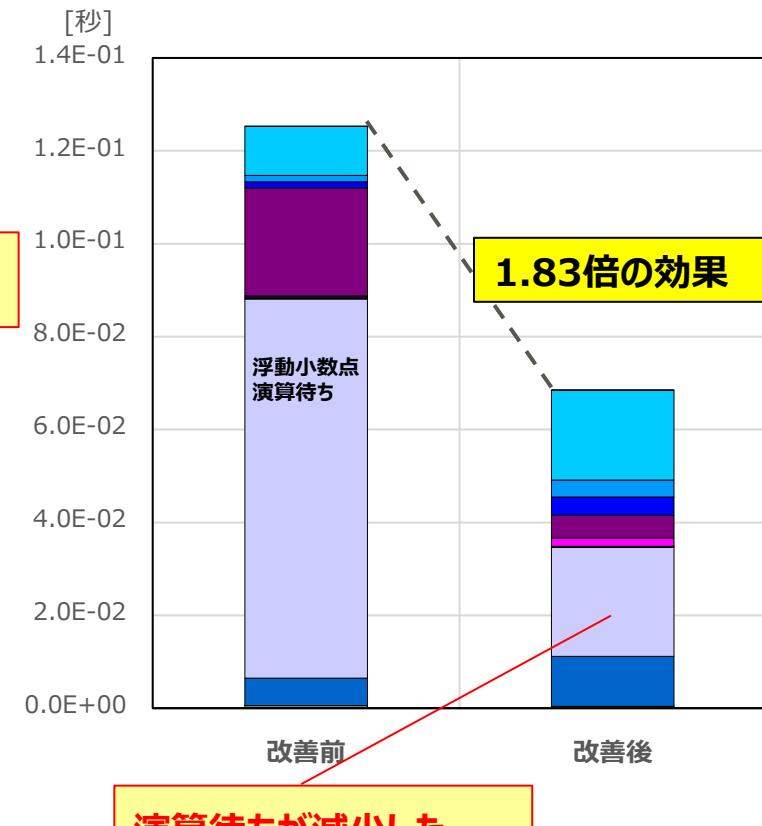
```

18   for (iter = 0; iter < 10000; iter++) {
19     #pragma loop loop_nofusion
20     <<< Loop-information Start >>>
21     <<< [OPTIMIZATION]
22     <<< SIMD(VL: 8)
23     <<< SOFTWARE PIPELINING(IPC: 1.02, ITR: 112, MVE: 3, POL: L)
24     <<< PREFETCH(HARD) Expected by compiler :
25     <<< (unknown)
26     <<< SPILLS :
27       GENERAL : SPILL 0 FILL 0
28       SIMD&FP : SPILL 0 FILL 0
29       SCALABLE : SPILL 0 FILL 27
30       PREDICATE : SPILL 0 FILL 0
31     <<< Loop-information End >>>
32     for (i = 0; i < n; i++){
33       y[i] = c2 / x1[i] / c3 / x1[i] / c4 / x1[i];
34     }
35     <<< Loop-information Start >>>
36     <<< [OPTIMIZATION]
37     <<< SIMD(VL: 8)
38     <<< SOFTWARE PIPELINING(IPC: 1.08, ITR: 64, MVE: 2, POL: S)
39     <<< PREFETCH(HARD) Expected by compiler :
40     <<< (unknown)
41     <<< Loop-information End >>>
42     for (i = 0; i < n; i++){
43       y[i] = c0 / x1[i] / c1 / x1[i] / y[i];
44     }
45   }

```

ループ融合を抑止

ループ分割



※コンパイルオプション-Knoevalを指定

Cache	L1I miss rate (/Effective instruction)	Load-store instruction	L1D miss	L1D miss rate (/Load-store instruction)	L1D miss demand rate (%) (/L1D miss)	L1D miss hardware prefetch rate (%) (/L1D miss)	L1D miss software prefetch rate (%) (/L1D miss)
改善前	0.00	1.11E+09	1.72E+06	0.00	99.92%	0.08%	-0.01%
改善後	0.00	6.46E+08	9.82E+05	0.00	99.88%	0.11%	0.00%

演算の連鎖が長く、多くのレジスタを必要とします。

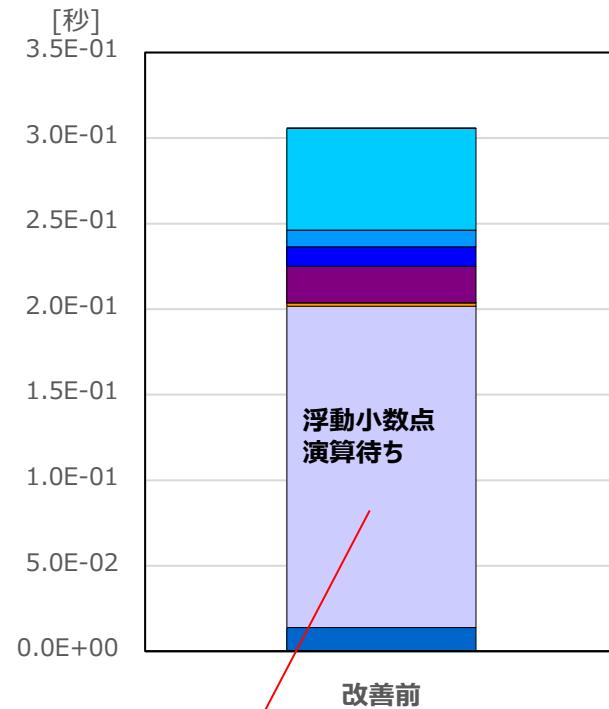
そのため、SWPLなどのスケジューリングの最適化ができず、浮動小数点演算待ちが多くなっています。

改善前ソース

```

19   for (iter = 0; iter < 10000; iter++) {
<<< Loop-information Start >>>
<<< [OPTIMIZATION]
<<< SIMD(VL: 8 Interleave: 1)
<<< SPILLS :
<<< GENERAL : SPILL 0 FILL 0
<<< SIMD&FP : SPILL 0 FILL 8
<<< SCALABLE : SPILL 4 FILL 6
<<< PREDICATE : SPILL 0 FILL 0
<<< Loop-information End >>>
20   v   for (i = 0; i < n; i++){
21     y[i] = sin(x1[i]*c0)
22       +cos(x1[i]*c1)
23       +atan(x1[i]*c2)
24       +log(x1[i]*c3);
25   }
26 }
```

演算連鎖が長く、多くのレジスタを使用するため、
ソフトウェアパイプラインングを適用できない



※コンパイラの特性を考慮して数学関数を例にした

スケジューリングが最適化できていないため、
演算待ちが大きくみえる

Cache	L1I miss rate (/Effective instruction)	Load-store instruction	L1D miss	L1D miss rate (/Load-store instruction)	L1D miss demand rate (%) (/L1D miss)	L1D miss hardware prefetch rate (%) (/L1D miss)	L1D miss software prefetch rate (%) (/L1D miss)	L2 miss	L2 miss rate (/Load-store instruction)	L2 miss demand rate (%) (/L2 miss)	L2 miss hardware prefetch rate (%) (/L2 miss)	L2 miss software prefetch rate (%) (/L2 miss)
改善前	0.00	1.76E+09	2.04E+06	0.00	99.85%	0.14%	0.01%	9.95E+03	0.00	88.46%	18.51%	0.00%

ループ分割を行うことで演算の連鎖が短くなり、1ループに使用するレジスタが減少しました。
その結果、SWPLなどのスケジューリングの最適化が行われ、演算待ちが改善されました。

改善後ソース

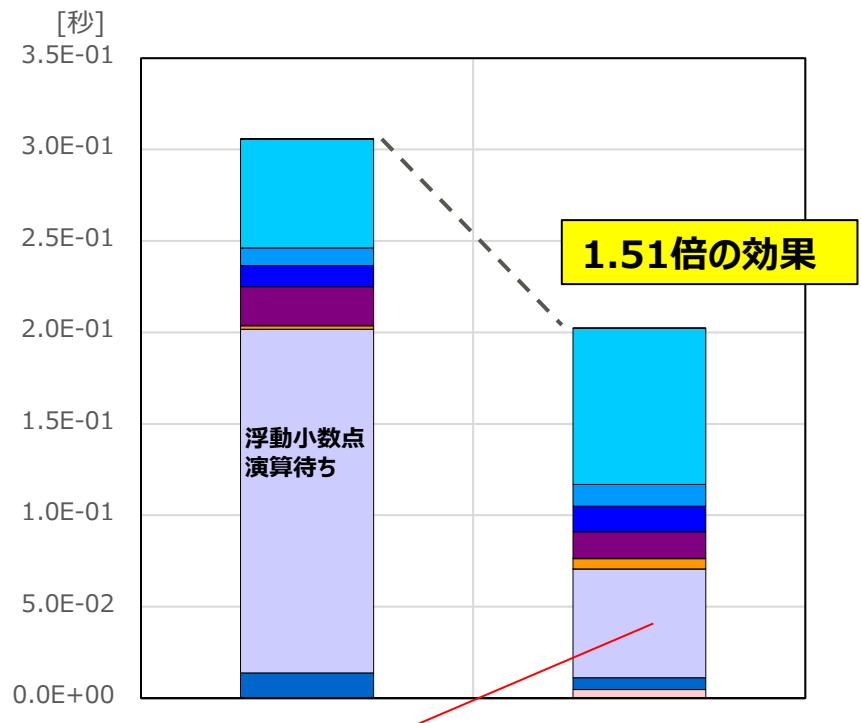
```

19   for (iter = 0; iter < 10000; iter++) {
20     #pragma loop loop_nofusion
21     <<< Loop-information Start >>>
22     <<< [OPTIMIZATION]
23     <<< SIMD(VL: 8 Interleave: 1)
24     <<< SOFTWARE PIPELINING
25     <<< SPILLS :
26       <<< GENERAL : SPILL 0 FILL 0
27       <<< SIMD&FP : SPILL 0 FILL 0
28       <<< SCALABLE : SPILL 8 FILL 17
29       <<< PREDICATE : SPILL 0 FILL 0
30     <<< Loop-information End >>>
31     v   for (i = 0; i < n; i++) {
32       y[i] = sin(x1[i]*c0);
33     }
34     <<< Loop-information Start >>>
35     <<< [OPTIMIZATION]
36     (略)
37     <<< Loop-information End >>>
38     v   for (i = 0; i < n; i++) {
39       y[i] += cos(x1[i]*c1);
40     }
41     <<< Loop-information Start >>>
42     <<< [OPTIMIZATION]
43     (略)
44     <<< Loop-information End >>>
45     v   for (i = 0; i < n; i++) {
46       y[i] += atan(x1[i]*c2);
47     }
48     <<< Loop-information Start >>>
49     <<< [OPTIMIZATION]
50     (略)
51     <<< Loop-information End >>>
52     v   for (i = 0; i < n; i++) {
53       y[i] += log(x1[i]*c3);
54     }
55   }

```

ループ融合を抑止

ループ分割



1.51倍の効果

※コンパイラの特性を考慮して数学関数を例にした

Cache	L1I miss rate (/Effective instruction)	Load-store instruction	L1D miss	L1D miss rate (/Load-store instruction)	L1D miss demand rate (%) (/L1D miss)	L1D miss hardware prefetch rate (%) (/L1D miss)	L1D miss software prefetch rate (%) (/L1D miss)
改善前	0.00	1.76E+09	2.04E+06	0.00	99.85%	0.14%	0.01%
改善後	0.00	2.67E+09	2.39E+06	0.00	99.91%	0.10%	-0.01%

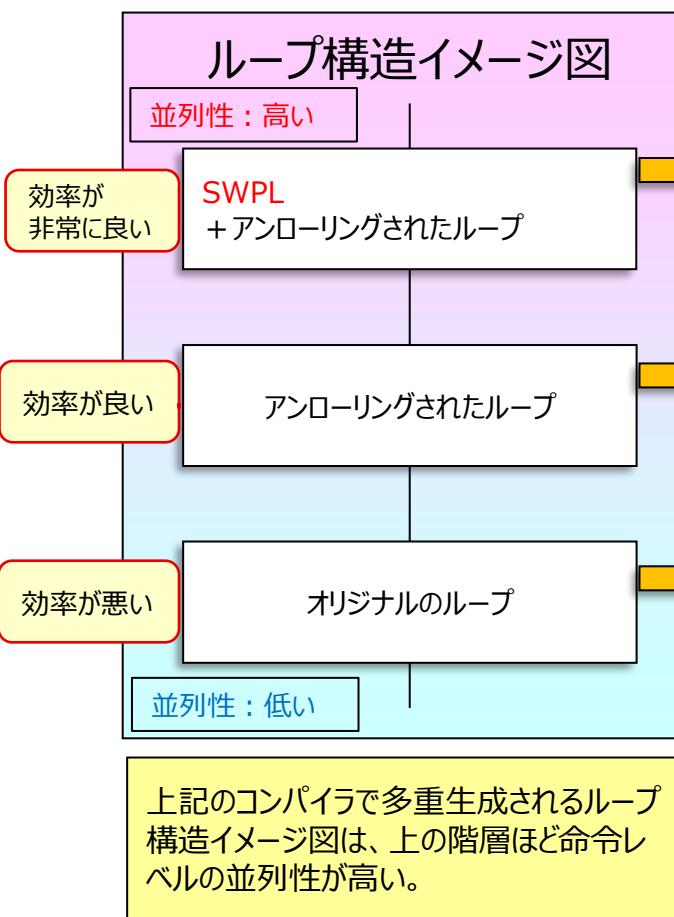
適切なアンローリング展開数の指定および ソフトウェアパイプライニングの抑止

- ・ ソフトウェアパイライン後のループの実行動作
- ・ 適切なアンローリング展開数の指定およびソフトウェアパイプライニングの抑止（改善前）
- ・ 適切なアンローリング展開数の指定およびソフトウェアパイプライニングの抑止（最適化制御行チューニング）
- ・ 適切なアンローリング展開数の指定およびソフトウェアパイプライニングの抑止（最適化制御行）

ソフトウェアパイプライン後のループの実行動作

FUJITSU

- ソフトウェアパイプラインではループの繰り返し数により、以下の処理ルートが決定されます。



例

<<< Loop-information Start >>>

<<< [OPTIMIZATION]

<<< SIMD(VL: 8)

<<< SOFTWARE PIPELINING(IPC: 2.50, ITR: 144, MVE: 4, POL: S)

<<< Loop-information End

17 1 pp 2v do i=1,N
18 1 p 2v b(i)=a(i)+c
19 1 p 2v enddo

例題の想定
ループ繰り返し数n=312
アンローリング展開数=2
SWPL
8SIMD

if(繰り返し数がSWPLの提示の値(144回転以上)の処理) then

ループの繰り返し数が144回以上の場合、ソフトウェアパイプラインを適用したループが実行時に選択されます。

例題ループの場合、i=1～288までが当ループで実行されます。

if(上記の残り繰り返し数の内、繰り返し数が16回転の倍数分の処理) then

アンローリング展開数2×8(SIMD)=16

例題ループの場合、i=289～304までが当ループで実行されます。

if(上記の残り繰り返し数の内、繰り返し数が8回転の倍数分の処理) then

8(SIMD)

例題ループの場合、残りのi=305～312までが当ループで実行されます。

上記のようにループ繰り返し数によって処理ルートが決まるため、
ループの繰り返し数が少ない場合は、命令レベルの並列性が高いループが
実行されないので命令スケジューリングの効果が小さくなる。

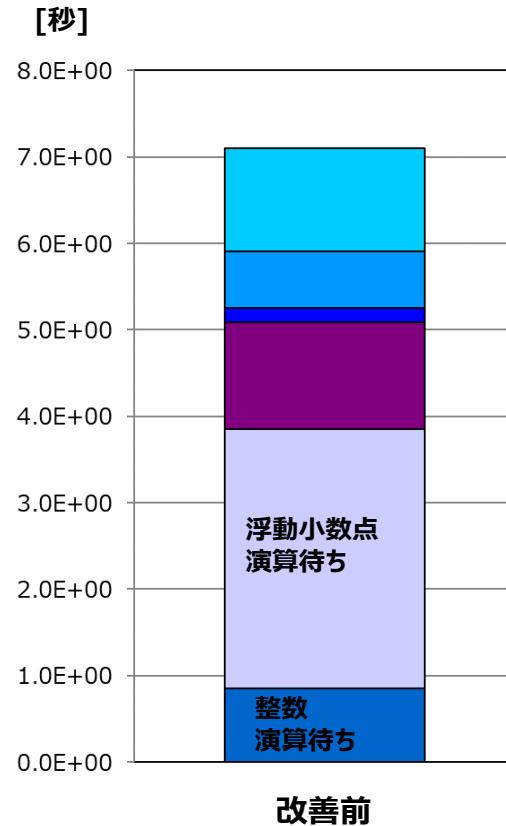
繰り返し数が少なく、アンローリングやソフトウェアパイプライニングが効果的に機能していません。
そのため、浮動小数点演算待ちや整数演算待ちが多くなっています。

改善前ソース

```
<<< Loop-information Start >>>
<<< [PARALLELIZATION]
<<< Standard iteration count: ループ繰り返し数n=40 に対し
<<< [OPTIMIZATION] アンロール展開数2
<<< SIMD(VL: 8)
<<< SOFTWARE PIPELINING(IPC: 1.62, ITR: 176,
                           MVE: 6, POL: S)
<<< PREFETCH(HARD) Expected by compiler:
<<<   a, b
<<< Loop-information End >>>
39  1 pp 2v  do i = 1 , n
40  1 p  2v    b(i) = c0 + a(i)*(c1 + a(i)*(c2 + a(i)*(c3 + a(i)*
41  1           & (c4 + a(i)*(c5 + a(i)*(c6 + a(i)*(c7 + a(i)*
42  1           & (c8 + a(i)*c9)))))))
43  1 p  2v    enddo
```

問題点：ループの繰り返し数が少ない

- ・ソフトウェアパイプライニングのループが実行されない
ソフトウェアパイプライニングの条件ITR: 176以下である。
- ・アンローリング展開数が適切でない
アンロール展開数 2×8 (SIMD) = 16
8回転がオリジナルのループで実行されてしまう。
オリジナルのループが実行されることによる性能への影響は大きい。



改善前

繰り返し数に応じたアンロール展開数の指定およびソフトウェアパイプライニングの抑止により、適切な命令スケジューリングが行われました。その結果、浮動小数点演算待ち・整数演算待ちが削減されました。

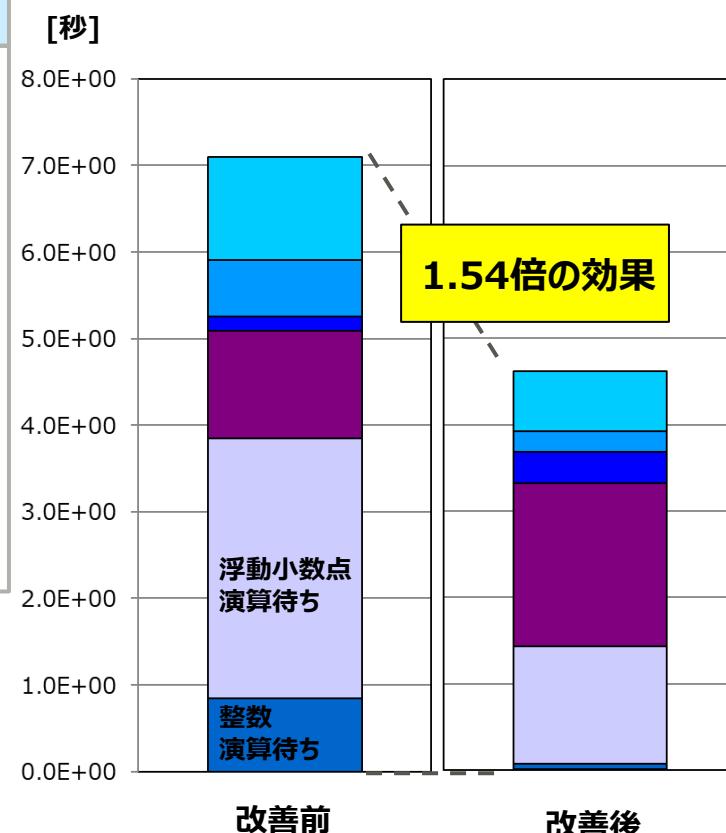
改善後ソース（最適化制御行チューニング）

```

39      !ocl unroll(5)
40      !ocl noswp
41      <<< Loop-information Start >>>
42      <<< [PARALLELIZATION]
43      <<< Standard iteration count: 55
44      <<< [OPTIMIZATION]
45      <<< SIMD(VL: 8)
46      <<< PREFETCH(HARD) Expected by compiler :
47      <<<     a, b
48      <<< Loop-information End >>>
49      1 pp 5v do i = 1 , n
50      1 p 5v   b(i) = c0 + a(i)*(c1 + a(i)*(c2 + a(i)*(c3 + a(i)*
51      &           (c4 + a(i)*(c5 + a(i)*(c6 + a(i)*(c7 + a(i)*
52      &           (c8 + a(i)*c9)))))))
53      1 p 5v enddo

```

アンロール展開数 5指定により
アンロール展開数 5×8 (SIMD) = 40
40回転全てアンローリングされたループが実行される。

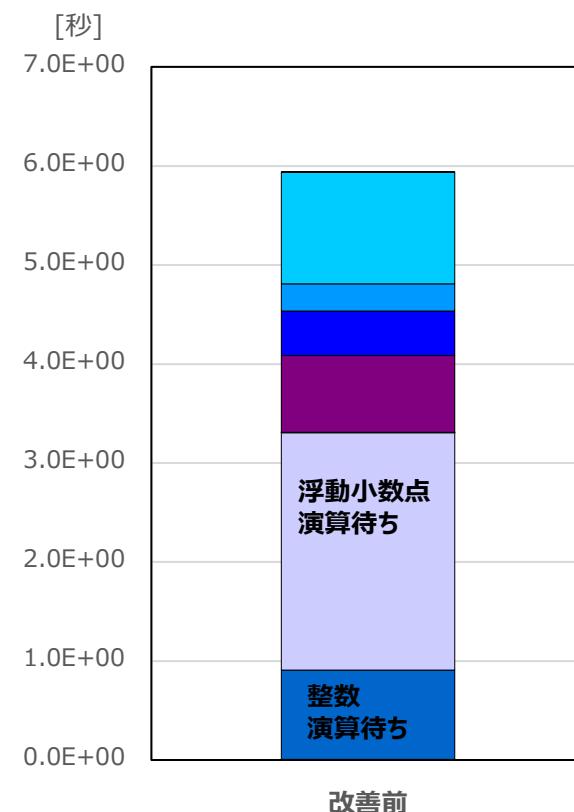


繰り返し数が少なく、アンローリングやソフトウェアパイプライニングが効果的に機能していません。
そのため、浮動小数点演算待ちや整数演算待ちが多くなっています。

改善前ソース

```

<<< Loop-information Start >>>
<<< [OPTIMIZATION]
<<< SIMD(VL: 8)
<<< SOFTWARE PIPELINING(IPC: 1.62, ITR: 176,
                           MVE: 6, POL: S)
<<< PREFETCH(HARD) Expected by compiler :
<<< (unknown)
<<< Loop-information End >>>
31  2v  for (i = 0; i < n; i++) {
32    2v    b[i] = c0 + a[i]*(c1 + a[i]*(c2 + a[i]*(c3 + a[i]*
33      (c4 + a[i]*(c5 + a[i]*(c6 + a[i]*(c7 + a[i]*
34        (c8 + a[i]*c9))))));
35  2v  }
36  }
```



問題点：ループの繰り返し数が少ない

- ・ソフトウェアパイプライニングのループが実行されない
ソフトウェアパイプライニングの条件ITR: 176以下である。
- ・アンローリング展開数が適切でない
アンロール展開数 2×8 (SIMD) = 16
8回転がオリジナルのループで実行されてしまう。
オリジナルのループが実行されることによる性能への影響は大きい。

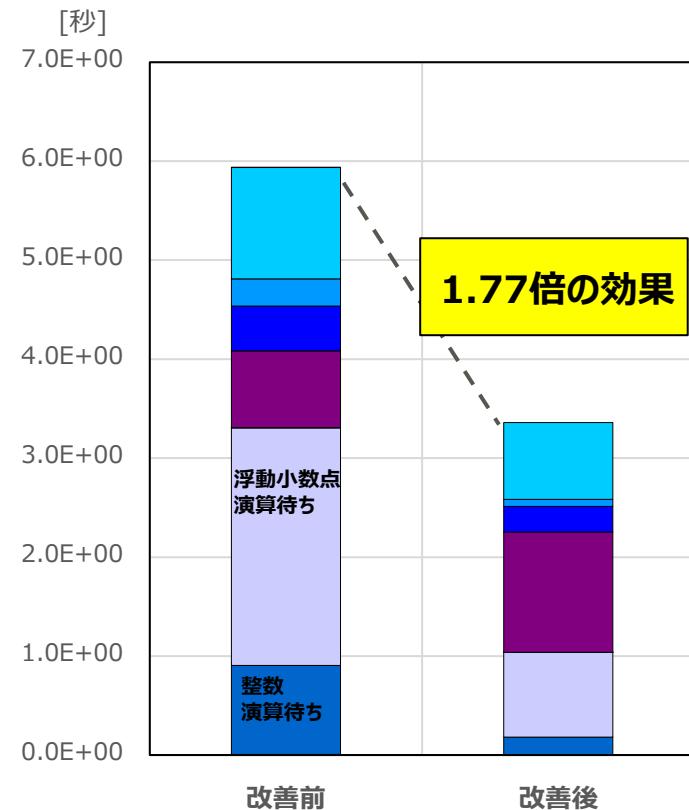
繰り返し数に応じたアンロール展開数の指定およびソフトウェアパイプライニングの抑止により、適切な命令スケジューリングが行われました。その結果、浮動小数点演算待ち・整数演算待ちが削減されました。

改善後ソース（最適化制御行チューニング）

```

31 #pragma loop unroll 5
32 #pragma loop noswp
33 <<< Loop-information Start >>>
34 <<< [OPTIMIZATION]
35 <<< SIMD(VL: 8)
36 <<< PREFETCH(HARD) Expected by compiler:
37 <<< (unknown)
38 <<< Loop-information End >>>
39
40 5v for (i = 0; i < n; i++){
41     b[i] = c0 + a[i]*(c1 + a[i]*(c2 + a[i]*(c3 + a[i]*
42         (c4 + a[i]*(c5 + a[i]*(c6 + a[i]*(c7 + a[i]*
43             (c8 + a[i]*c9)))))));
44 }
45 }
```

アンロール展開数 5指定により
アンロール展開数 5×8 (SIMD) = 40
40回転全てアンローリングされたループが実行される。



適切なアンローリング展開数の指定およびソフトウェアパイプライン ングの抑止（最適化制御行）

FUJITSU

以下の最適化指示子を指定します。または、翻訳オプションで指定することも可能です。

最適化指示子 (Fortran)	意味	指定可能な最適化制御行			
		プログラム単位	DOループ単位	文単位	配列代入文単位
UNROLL(<i>n</i>)	DO ループをアンローリングさせます。 <i>n</i> は展開数（多重度）を表す2～100の10進数です。	×	○	×	×
NOSWP	ソフトウェアパイプラインング機能を無効にします。	○	○	×	○

最適化指示子 (C/C++)	意味	指定可能な最適化制御行			
		global行	procedure行	loop行	statement行
unroll(<i>n</i>)	ループをアンローリングさせます。 <i>n</i> は展開数（多重度）を表す2～100の10進数です。	×	×	○	×
noswp	ソフトウェアパイプラインング機能を無効にします。	○	○	○	×

翻訳オプション	機能説明
-Kunroll[= <i>N</i>] (2≤ <i>N</i> ≤100)	ループアンローリングの最適化を行うことを指示します。 <i>N</i> にはループ展開数の上限を指定します。 <i>N</i> の指定を省略した場合、コンパイラが自動的に最良な値を決定します。 -O0または-O1オプションが有効な場合のデフォルトは-Knunroll、-O2オプション以上が有効な場合のデフォルトは-Kunrollです。
-Knoswp	ソフトウェアパイプラインングの最適化を行わないことを指示します。

■使用例

```
$ frtpx -Kfast,parallel sample.f90 -Kunroll=5,noswp
$ fccpx -Kfast,parallel sample.f90 -Kunroll=5,noswp
```

◆注意事項

clangモードではアンローリングの最適化機能を使用できません

ストライピング（インターリーブ）展開数の指 定および ソフトウェアパライニングの抑止

- ループ展開
- ストライピング（インターリーブ）展開数の指定およびソフトウェアパライニングの抑止（改善前）
- ストライピング（インターリーブ）展開数の指定およびソフトウェアパライニングの抑止（最適化制御行チューニング）
- ストライピング（インターリーブ）展開数の指定およびソフトウェアパライニングの抑止（最適化制御行）

- ループ展開には次の2種類あります。

- アンロール(unroll)
- ストライピング(striping)

アンローリングとストライピングの違い

展開の方法が違う

```
DO I=1,N
  A(I) = B(I) + C(I)
ENDDO
```

赤：1展開目
青：2展開目

■ アンローリング2展開のイメージ

```
DO I=1,N,2
  TMP_B1 = B(I)
  TMP_C1 = C(I)
  TMP_A1 = TMP_B1 + TMP_C1
  A(I) = TMP_A1
  TMP_B2 = B(I+1)
  TMP_C2 = C(I+1)
  TMP_A2 = TMP_B2 + TMP_C2
  A(I+1) = TMP_A2
ENDDO
```

■ ストライピング2展開のイメージ

```
DO I=1,N,2
  TMP_B1 = B(I)
  TMP_B2 = B(I+1)
  TMP_C1 = C(I)
  TMP_C2 = C(I+1)
  TMP_A1 = TMP_B1 + TMP_C1
  TMP_A2 = TMP_B2 + TMP_C2
  A(I) = TMP_A1
  A(I+1) = TMP_A2
ENDDO
```

● ポイント

- SIMD化、ソフトウェアパライニアリングとの調整で効果が見込める可能性があるため、unrollを先に適用して効果がなかった場合にストライピングを適用することをお勧めします。
- unrollよりstripingの方が使用レジスタ数が増加するため、ストライプ長nを多くすると実行性能が低下する場合があります。
- Kstripingオプションと-Kunrollオプションを同時に指定した場合、後に指定した方が有効となります。

繰り返し数が少なく、アンローリングやソフトウェアパイプライニングが効果的に機能していません。
そのため、浮動小数点演算待ちや整数演算待ちが多くなっています。

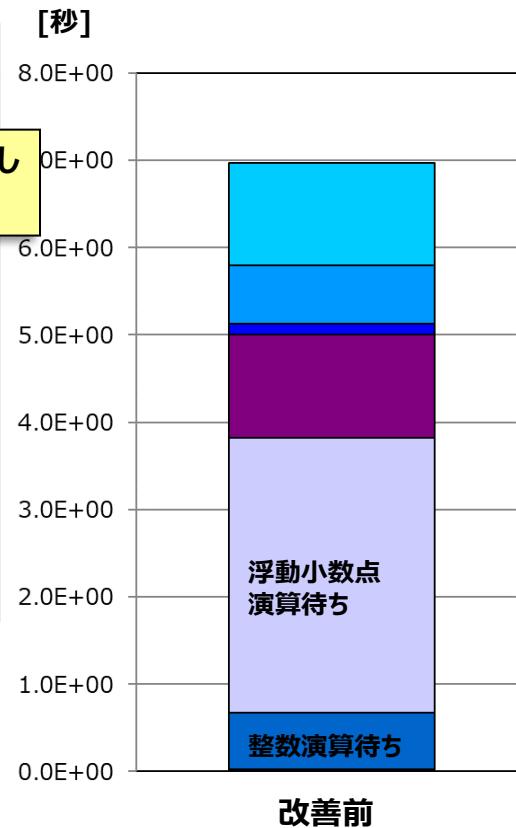
改善前ソース

```

<<< Loop-information Start >>>
<<< [PARALLELIZATION]
<<< Standard iteration count:
<<< [OPTIMIZATION]
<<< SIMD(VL: 8)
<<< SOFTWARE PIPELINING(IPC: 1.62, ITR: 176,
                           MVE: 6, POL: S)
<<< PREFETCH(HARD) Expected by compiler :
<<<     a, b
<<< Loop-information End >>>

39  1 pp 2v  do i = 1 , n
40  1 p   2v    b(i) = c0 + a(i)*(c1 + a(i)*(c2 + a(i)*(c3 + a(i)*
41  1      &      (c4 + a(i)*(c5 + a(i)*(c6 + a(i)*(c7 + a(i)*
42  1      &      (c8 + a(i)*c9)))))))
43  1 p   2v    enddo

```



問題点：ループの繰り返し数が少ない

- ・ソフトウェアパイプライニングのループが実行されない
ソフトウェアパイプライニングの条件ITR: 176以下である。
- ・アンローリング展開数が適切でない
アンロール展開数 2×8 (SIMD) = 16
8回転がオリジナルのループが実行されてしまう。
- オリジナルのループが実行されることによる性能への影響は大きい。

繰り返し数に応じたストライピング展開数の指定およびソフトウェアパイプライニングの抑止により、適切な命令スケジューリングが行われました。
その結果、浮動小数点演算待ち・整数演算待ちが削減されました。

改善後ソース（最適化制御行チューニング）

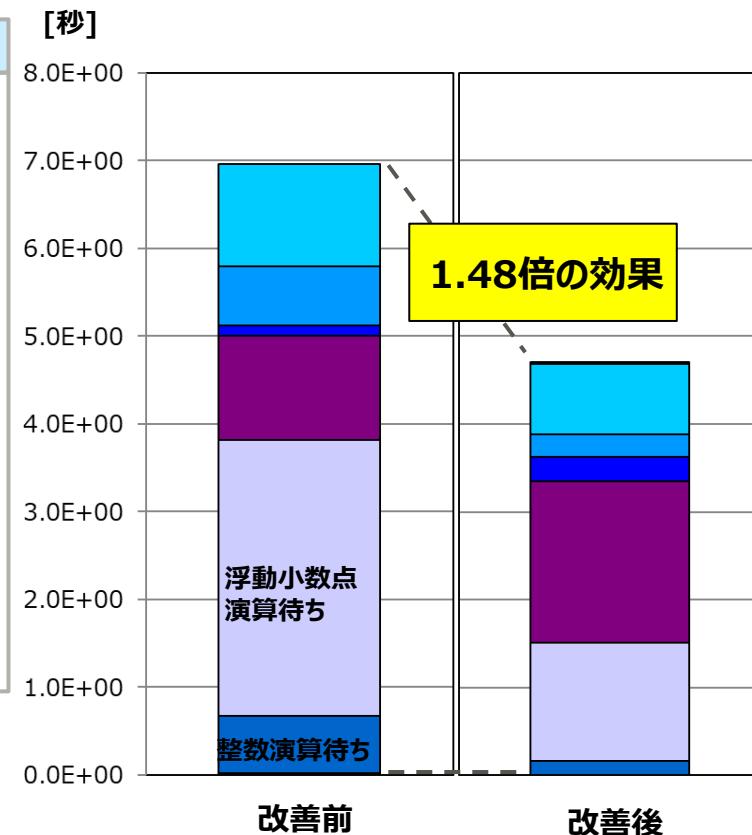
```

39      !ocl striping(5)
40      !ocl nounroll
41      !ocl noswp
42      <<< Loop-information Start >>>
43      <<< [PARALLELIZATION]
44      <<< Standard iteration count: 552
45      <<< [OPTIMIZATION]
46      <<< SIMD(VL: 8)
47      <<< PREFETCH(HARD) Expected by compiler
48      <<< a, b
49      <<< Loop-information End >>>
50      1 pp 5v    do i = 1 , n
51      1 p 5v      b(i) = c0 + a(i)*(c1 + a(i)*(c2 + a(i)*(c3 + a(i)*
52      1           & (c4 + a(i)*(c5 + a(i)*(c6 + a(i)*(c7 + a(i)*
53      1           & (c8 + a(i)*c9)))))))
54      1 p 5v    enddo

```

ストライピング
数5に指定

アンロール、
ソフトウェアパイプラ
イニングの抑止



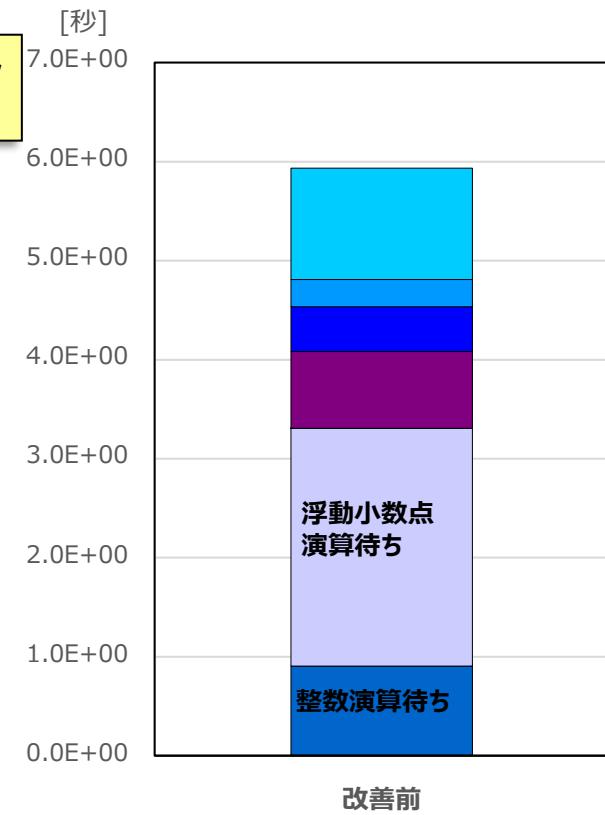
ストライピング展開数 5 指定により
ストライピング展開数 5×8 (SIMD) = 40
40回転全てストライピングされたループが実行される。

繰り返し数が少なく、アンローリングやソフトウェアパイプライニングが効果的に機能していません。
そのため、浮動小数点演算待ちや整数演算待ちが多くなっています。

改善前ソース

```
<<< Loop-information Start >>>
<<< [OPTIMIZATION]
<<< SIMD(VL: 8)
<<< SOFTWARE PIPELINING(IPC: 1.62, ITR: 176,
                           MVE: 6, POL: S)
<<< PREFETCH(HARD) Expected by compiler :
<<< (unknown)
<<< Loop-information End >>>
33   2v  for (i = 0; i < n; i++) {
34     2v    b[i] = c0 + a[i]*(c1 + a[i]*(c2 + a[i]*(c3 + a[i]*
35      (c4 + a[i]*(c5 + a[i]*(c6 + a[i]*(c7 + a[i]*
36        (c8 + a[i]*c9)))))));
37   2v  }
38 }
```

ループ繰り返し数n=40 に対し
アンロール展開数2



問題点：ループの繰り返し数が少ない

- ・ソフトウェアパイプライニングのループが実行されない
ソフトウェアパイプライニングの条件ITR: 176以下である。

・アンローリング展開数が適切でない

アンロール展開数 2×8 (SIMD) = 16

8回転がオリジナルのループが実行されてしまう。

オリジナルのループが実行されることによる性能への影響は大きい。

繰り返し数に応じたストライピング展開数の指定およびソフトウェアパイプライニングの抑止により、適切な命令スケジューリングが行われました。
その結果、浮動小数点演算待ち・整数演算待ちが削減されました。

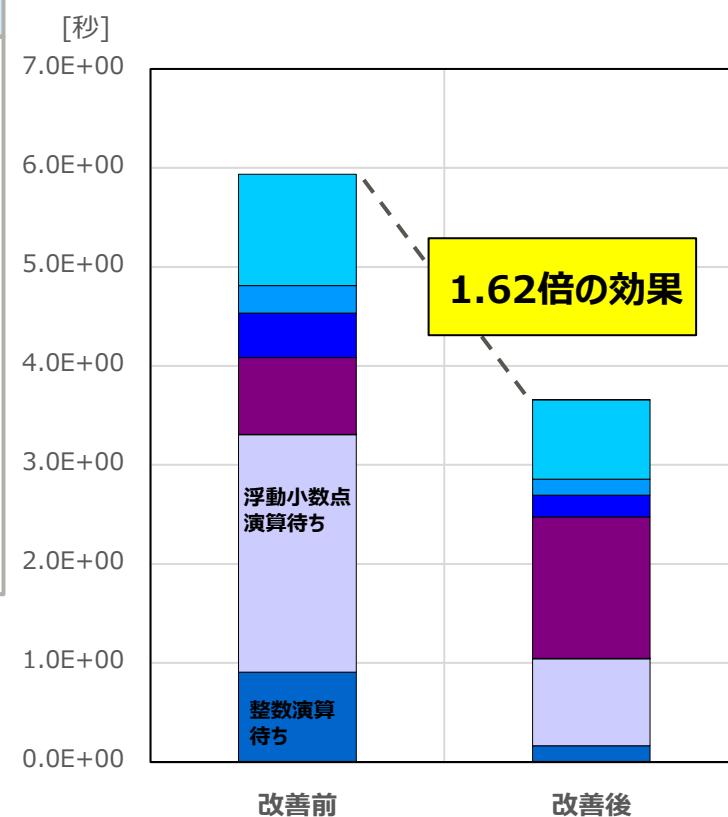
改善後ソース（最適化制御行チューニング）

```

33 #pragma loop striping 5
34 #pragma loop nounroll
35 #pragma loop noswp
36 <<< Loop-information Start >>>
37 <<< [OPTIMIZATION]
38 <<< SIMD(VL: 8)
39 <<< STRIPING
40 <<< PREFETCH(HARD) Expected by compiler :
41 <<< (unknown)
<<< Loop-information End >>>
36 v   for (i = 0; i < n; i++){
37   v     b[i] = c0 + a[i]*(c1 + a[i]*(c2 + a[i]*(c3 + a[i]*
38     (c4 + a[i]*(c5 + a[i]*(c6 + a[i]*(c7 + a[i]*
39     (c8 + a[i]*c9))))));
40   v }
41 }
```

ストライピング
数5に指定

アンロール、
ソフトウェアパイ
プライングの抑止



ストライピング展開数 5 指定により

ストライピング展開数 5×8 (SIMD) = 40

40回転全てストライピングされたループが実行される。

ストライピング（インターリーブ）展開数の指定およびソフトウェア パイプライニングの抑止（最適化制御行）

FUJITSU

以下の最適化指示子を指定します。または、翻訳オプションで指定することも可能です。

最適化指示子 (Fortran)	意味	指定可能な最適化制御行			
		プログラム単位	DOループ単位	文単位	配列代入文単位
STRIPING[(n)]	ループストライピング機能を有効にします。nは展開数(多重度)を表す2～100の10進数です。	○	○	×	○
NOSWP	ソフトウェアパイプライニング機能を無効にします。	○	○	×	○

最適化指示子 (C/C++)	意味	指定可能な最適化制御行			
		global行	procedure行	loop行	statement行
striping[(n)]	ループストライピング機能を有効にします。nは展開数(多重度)を表す2～100の10進数です。	○	○	○	×
noswp	ソフトウェアパイプライニング機能を無効にします。	○	○	○	×

翻訳オプション	機能説明
-Kstriping[=N] (2≤N≤100)	ループストライピングの最適化を行うかどうかを指示します。ストライプ長(展開数)をNで指定することができます。Nは2から100までの値を指定することができます。Nの値が省略された場合、コンパイラが自動的に値を決定します。ソースプログラムでループの繰返し数が既知の場合、Nに繰返し数を超える値を指定しても、コンパイラが自動的に決定した展開数が有効となります。デフォルトは-Knostripingです。
-Knoswp	ソフトウェアパイプライニングの最適化を行わないことを指示します。

■使用例

```
$ frtpx -Kfast,parallel sample.f90 -Kstriping=5,noswp
$ fccpx -Kfast,parallel sample.f90 -Kstriping=5,noswp
```

◆注意事項

clangモードではストライピング機能を使用できません

外側ループでのソフトウェアパイプライン

- ・ 外側ループでのソフトウェアパイプラインとは
- ・ 外側ループでのソフトウェアパイプライン (改善前)
- ・ 外側ループでのソフトウェアパイプライン (ソースチューニング)
- ・ 外側ループでのソフトウェアパイプライン (CLONE利用)
- ・ 外側ループでのソフトウェアパイプライン (CLONE利用) (改善前)
- ・ 外側ループでのソフトウェアパイプライン (CLONE利用) (ソースチューニング)

外側ループでのソフトウェアパイプラインングとは

FUJITSU

- 最内ループの繰り返し数が少ない場合は、ストリップマイニング等を行って最内ループの繰り返し数をSIMD長と等しくすることで、その外側のループでソフトウェアパイプラインングを促進させることができます。

改善前ソース

```
24 3 p      do j=1,M
    <<< Loop-information Start >>>
    <<< [OPTIMIZATION]
    <<< SIMD(VL: 8)
    <<< SOFTWARE PIPELINING(IPC: 3.00, ITR: 192,
                                MVE: 7, POL: S)
        <<< PREFETCH(HARD) Expected by compiler :
        <<< b, a
    <<< Loop-information End >>>
25 4 p 2v    do i=1,N
26 4 p 2v    a(i,j,k)=a(i,j,k)+c*b(i,j,k)
27 4 p 2v    enddo
28 3 p      enddo
```

改善後ソース

```
25 3 p      do ii=1,N,blk
    <<< Loop-information Start >>>
    <<< [OPTIMIZATION]
    <<< SOFTWARE PIPELINING(IPC: 0.31, ITR: 192,
                                MVE: 2, POL: L)
        <<< PREFETCH(HARD) Expected by compiler :
        <<< b, a
    <<< Loop-information End >>>
26 4 p 8     do j=1,M
    <<< Loop-information Start >>>
    <<< [OPTIMIZATION]
    <<< SIMD(VL: 8)
        <<< FULL UNROLLING
    <<< Loop-information End >>>
27 5 p fv    do i=ii,ii+blk-1
28 5 p fv    a(i,j,k)=a(i,j,k)+c*b(i,j,k)
29 5 p fv    enddo
30 4 p 8     enddo
31 3 p      enddo
```

上記の例は固定長SIMDの場合に有効です。
SIMD幅512[bit]の時のSIMD長は以下の通りです。

SIMD幅(Vector Length)=512[bit]の場合のSIMD長

データ型	SIMD長
倍精度型	8
单精度型	16
半精度型	32
1バイト型	64

最内ループの繰り返し数Nがソフトウェアパイプラインングの条件に対して小さいため、ソフトウェアパイプラインングが適用されていません。

改善前ソース

```

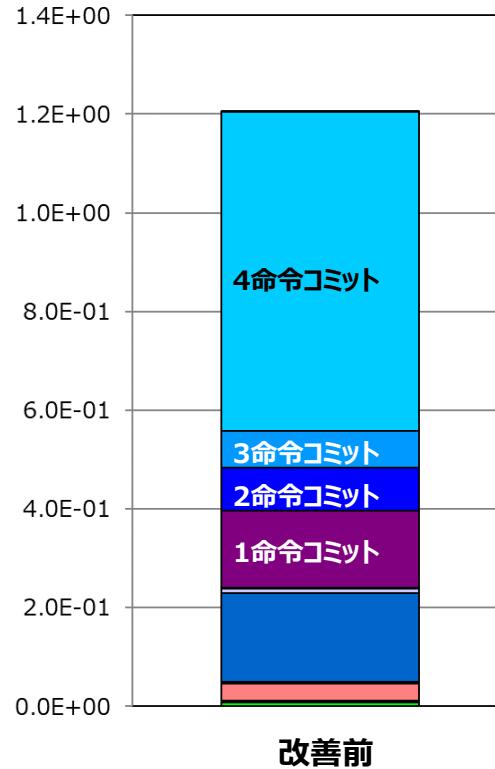
2      integer,parameter::N=16,M=200,L=48
:
18      real(8)::a(N,M,L),b(N,M,L)
19      real(8),parameter::c=0.5
20      !$omp parallel private(iter,i,j,k)
21      1      do iter=1,itmax
22      1      !$omp do
23      2      p      do k=1,L
24      3      p      <<< Loop-information Start >>>
25      4      p      2v      <<< [OPTIMIZATION]
26      4      p      2v      <<< PREFETCH(HARD) :
27      4      p      2v      <<< b, a
28      3      p      <<< Loop-information End >>>
29      2      p      <<< SIMD(VL: 8)
30      1      !$omp enddo nowait
31      1      enddo
32      !$omp end parallel

```

最内ループの繰り返し数Nが
 ITR:192よりも少ないため
 ソフトウェアパイプラインングの
 処理ルートに入らない

N=16

[秒]



改善前

	GFLOPS	Effective instruction
改善前	50.98	7.53E+10

ストリップマイニングを行い最内ループをSIMD長に固定しました。その結果、外側のループでSWPLが促進され、演算効率・有効命令数が改善しました。

C/C++ではループの初期値と終値に変数が含まれている場合に、ループを全展開できずループ中に分岐が残ります。そのため、外側ループにSWPLを適用できません。

改善後ソース

```

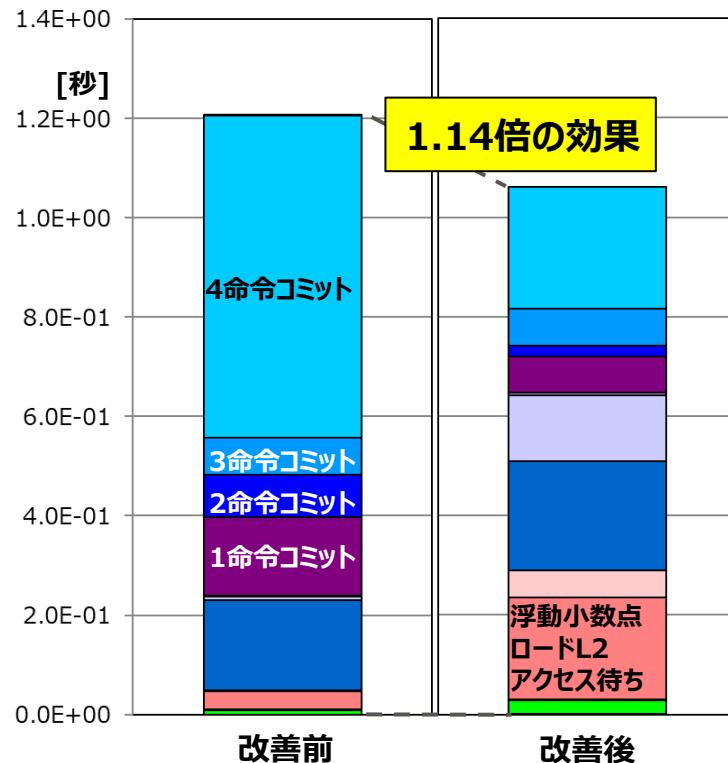
2
:
18      integer,parameter::N=16,M=200,L=48
19      real(8)::a(N,M,L),b(N,M,L)
20      real(8),parameter::c=0.5
21      integer,parameter::blk=8
22      !$omp parallel private(iter,ii,i,j,k)
23      1          do iter=1,itmax
24      1          !$omp do
25      2          do k=1,L
26      <<< Loop-information Start >>>
27      <<< [OPTIMIZATION]
28      <<< PREFETCH(HARD) Expected by compiler :
29      <<< b, a
30      <<< Loop-information End >>>
31      3          do ii=1,N,blk
32      <<< Loop-information Start >>>
33      <<< [OPTIMIZATION]
34      <<< SOFTWARE PIPELINING(IPC: 0.31, ITR: 192,
35      MVE: 2, POL: L)
36      <<< PREFETCH(HARD) Expected by compiler :
37      <<< b, a
38      <<< Loop-information End >>>
39      4          do j=1,M
40      <<< Loop-information Start >>>
41      <<< [OPTIMIZATION]
42      <<< SIMD(VL: 8)
43      <<< FULL UNROLLING
44      <<< Loop-information End >>>
45      5          p fv      do i=ii,ii+blk-1
46      5          p fv      a(i,j,k)=a(i,j,k)+c*b(i,j,k)
47      5          p fv      enddo
48      4          p 8      enddo
49      3          p      enddo
50      2          p      enddo
51      1          !$omp enddo nowait
52      1          enddo
53      !$omp end parallel

```

ソフトウェアパイプラインングが適用される

M=200

最内ループの繰り返し数をSIMD長に固定



	GFLOPS	Effective instruction
改善前	50.98	7.53E+10
改善後	57.95	3.19E+10

外側ループのソフトウェアパイプラインング (CLONE利用)

FUJITSU

最内ループの繰り返し回数が短く固定値の場合には、クローンチューニングを利用することもできます。

● クローンチューニングとは

CLONE指示子を用いて、ループ[°]に対して変数の値による条件分岐を生成し、フルアンローリングなどの最適化を促進するチューニング方法です。

ソース例

```
!ocl clone(N==8)
do i=1,N
    A(i) = i
enddo
```



最適化後 (ソースイメージ)

```
If (N==8) then
    do i=1,8
        A(i)=i
    endo
else
    do i=1,N
        A(i)=i
    enddo
endif
```

N==8の時のループ[°]を複製
ループ長が固定されることで
他の最適化が促進される

最適化指示子 (Fortran)	意味	指定可能な最適化制御行			
		プログラム単位	DOループ単位	文単位	配列代入文単位
CLONE(var= =n1[,n2]…)	ループ内で変数varが不变とみなして、引数で指示した条件分岐を生成し、条件節内にループを複製する最適化を行うことを指示します。条件式は第1引数の変数varと、第2引数以降で指定した値n1[,n2]…の等式とします。	×	○	×	○
最適化指示子 (C/C++)	意味	指定可能な最適化制御行			
		global行	procedure行	loop行	statement行
clone(var== n1[,n2]…)	ループ内で変数varが不变とみなして、引数で指示した条件分岐を生成し、条件節内にループを複製する最適化を行うことを指示します。条件式は第1引数の変数varと、第2引数以降で指定した値n1[,n2]…の等式とします。	×	×	○	×

最内ループの繰り返し数Nが小さいため、ソフトウェアパイプラインングの適用条件を満たしていません。

改善前ソース

```

2      integer,parameter::N=8,M=240,L=48
:
18      real(8)::a(N,M,L),b(N,M,L)
19      real(8),parameter::c=0.5
20      !$omp parallel private(iter,i,j,k)
21      1      do iter=1,itmax
22      1      !$omp do
23      2 p      do k=1,L
<<< Loop-information Start >>>
<<< [OPTIMIZATION]
<<< PREFETCH(HARD) Expected by compiler :
<<< b, a
<<< Loop-information End >>>
24      3 p      do j=1,M
<<< Loop-information Start >>>
<<< [OPTIMIZATION]
<<< SIMD(VL: 8)
<<< SOFTWARE PIPELINING(IPC: 3.00, ITR: 192,
                           MVE: 7, POL: S)
<<< PREFETCH(HARD) Expected by compiler :
<<< b, a
<<< Loop-information End >>> N=8
25      4 p      2v      do i=1,N
26      4 p      2v      a(i,j,k)=a(i,j,k)+c*b(i,j,k)
27      4 p      2v      enddo
28      3 p      enddo
29      2 p      enddo
30      1      !$omp enddo nowait
31      1      enddo
32      !$omp end parallel

```

最内ループの繰り返し数Nが
ITR:192よりも少ないため
ソフトウェアパイプラインングの
処理ルートに入らない

[秒]



	GFLOPS	Effective instruction
改善前	29.51	6.18E+10

CLONEを利用し最内ループの繰り返し数をSIMD長に固定しました。その結果、外側のループでソフトウェアパイプラインングが促進され、演算効率・有効命令数が改善しました。

改善後ソース (ソースチューニング)

```

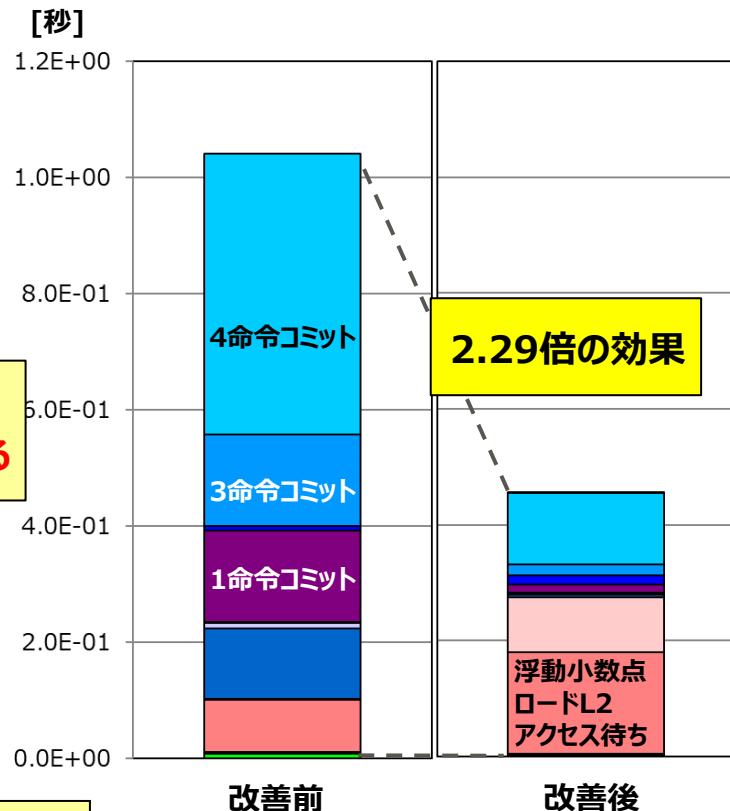
2      integer,parameter::N=8,M=240,L=48
:
18
19      real(8)::a(N,M,L),b(N,M,L)
20      real(8),parameter::c=0.5
21      integer NN
22      NN = N
23      !$omp parallel private(iter,i,j,k)
24      1      do iter=1,itmax
25      1      !$omp do
25      1      !ocl clone(NN==8)
<<< Loop-information Start >>>
<<< [OPTIMIZATION]
<<< CLONE
<<< Loop-information End >>>
26      2 p      do k=1,L
<<< Loop-information Start >>>
<<< [OPTIMIZATION]
<<< SOFTWARE PIPELINING(IPC: 3.75, ITR: 208,
MVE: 6, POL: S)
27      3 p      2      do j=1,M
<<< Loop-information End >>>
27      3 p      2      <<< [OPTIMIZATION]
27      3 p      2      SIMD(VL: 8)
27      3 p      2      <<< SOFTWARE PIPELINING(IPC: 3.00, ITR: 192,
MVE: 7, POL: S)
<<< Loop-information End >>>
28      4 p      2v      do i=1,NN
29      4 p      2v          a(i,j,k)=a(i,j,k)+c*b(i,j,k)
30      4 p      2v          enddo
31      3 p      2      enddo
32      2 p      enddo
33      1      !$omp enddo nowait
34      1      enddo
35      !$omp end parallel

```

ソフトウェアパイフ
ライニングが適用される

M=240

最内ループの繰り返し数を
SIMD長に固定



	GFLOPS	Effective instruction
改善前	29.51	6.18E+10
改善後	67.60	1.44E+10

最内ループの繰り返し数Nが小さいため、ソフトウェアパイプラインングの適用条件を満たしていません。

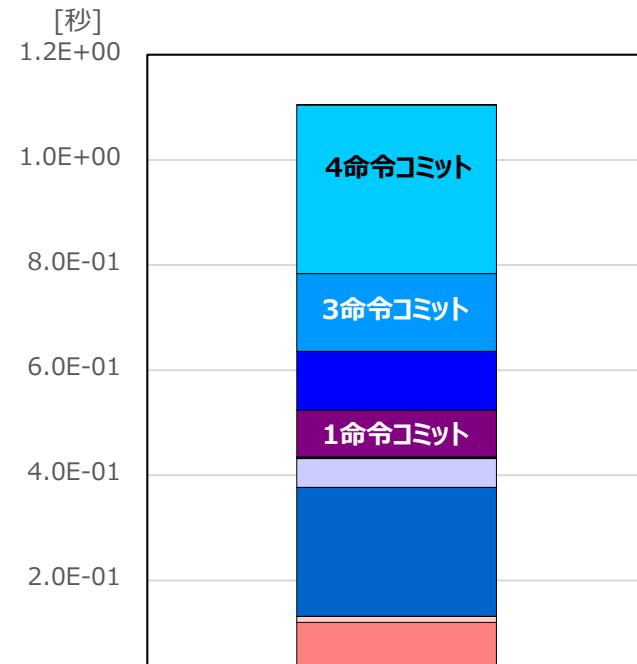
改善前ソース

```

35      #pragma omp parallel private(iter,i,j,k)
36      {
37          for (iter = 0; iter < itmax; iter++){
38              #pragma omp for nowait
39              p           for (k = 0; k < L; k++){
40                  p           <<< Loop-information Sta
41                  p           <<< [OPTIMIZATION]
42                  p           <<< PREFETCH(HARD)
43                  p           <<< (unknown)
44                  p           <<< Loop-information Er
45                  p           for (j = 0; j < M; j++)
46                  p           <<< Loop-information Sta
47                  p           <<< [OPTIMIZATION]
48                  p           <<< SIMD(VL: 8)
49                  p           <<< SOFTWARE PIPELINING(IPC: 3.00, ITR: 192,
50                                MVE: 7, POL: S)
51                  p           <<< PREFETCH(HARD) Expected by compiler :
52                  p           <<< (unknown)
53                  p           <<< Loop-information End>>> N=8
54                  p           2v           for (i = 0; i < N; i++){
55                  p           2v           a[k][j][i] = a[k][j][i] + c * b[k][j][i];
56                  p           2v           }
57                  p           }
58                  p           }
59                  p           }
60                  p           }
61                  p           }
62                  p           }
63                  p           }
64                  p           }
65                  p           }
66                  p           }
67                  p           }
68                  p           }
69                  p           }
70                  p           }
71                  p           }
72                  p           }
73                  p           }
74                  p           }
75                  p           }
76                  p           }
77                  p           }
78                  p           }
79                  p           }
80                  p           }
81                  p           }
82                  p           }
83                  p           }
84                  p           }
85                  p           }
86                  p           }
87                  p           }
88                  p           }
89                  p           }
90                  p           }
91                  p           }
92                  p           }
93                  p           }
94                  p           }
95                  p           }
96                  p           }
97                  p           }
98                  p           }
99                  p           }
100                 p           }
101                 p           }
102                 p           }
103                 p           }
104                 p           }
105                 p           }
106                 p           }
107                 p           }
108                 p           }
109                 p           }
110                 p           }
111                 p           }
112                 p           }
113                 p           }
114                 p           }
115                 p           }
116                 p           }
117                 p           }
118                 p           }
119                 p           }
120                 p           }
121                 p           }
122                 p           }
123                 p           }
124                 p           }
125                 p           }
126                 p           }
127                 p           }
128                 p           }
129                 p           }
130                 p           }
131                 p           }
132                 p           }
133                 p           }
134                 p           }
135                 p           }
136                 p           }
137                 p           }
138                 p           }
139                 p           }
140                 p           }
141                 p           }
142                 p           }
143                 p           }
144                 p           }
145                 p           }
146                 p           }
147                 p           }
148                 p           }
149                 p           }
150                 p           }
151                 p           }
152                 p           }
153                 p           }
154                 p           }
155                 p           }
156                 p           }
157                 p           }
158                 p           }
159                 p           }
160                 p           }
161                 p           }
162                 p           }
163                 p           }
164                 p           }
165                 p           }
166                 p           }
167                 p           }
168                 p           }
169                 p           }
170                 p           }
171                 p           }
172                 p           }
173                 p           }
174                 p           }
175                 p           }
176                 p           }
177                 p           }
178                 p           }
179                 p           }
180                 p           }
181                 p           }
182                 p           }
183                 p           }
184                 p           }
185                 p           }
186                 p           }
187                 p           }
188                 p           }
189                 p           }
190                 p           }
191                 p           }
192                 p           }
193                 p           }
194                 p           }
195                 p           }
196                 p           }
197                 p           }
198                 p           }
199                 p           }
200                 p           }
201                 p           }
202                 p           }
203                 p           }
204                 p           }
205                 p           }
206                 p           }
207                 p           }
208                 p           }
209                 p           }
210                 p           }
211                 p           }
212                 p           }
213                 p           }
214                 p           }
215                 p           }
216                 p           }
217                 p           }
218                 p           }
219                 p           }
220                 p           }
221                 p           }
222                 p           }
223                 p           }
224                 p           }
225                 p           }
226                 p           }
227                 p           }
228                 p           }
229                 p           }
230                 p           }
231                 p           }
232                 p           }
233                 p           }
234                 p           }
235                 p           }
236                 p           }
237                 p           }
238                 p           }
239                 p           }
240                 p           }
241                 p           }
242                 p           }
243                 p           }
244                 p           }
245                 p           }
246                 p           }
247                 p           }
248                 p           }
249                 p           }
250                 p           }
251                 p           }
252                 p           }
253                 p           }
254                 p           }
255                 p           }
256                 p           }
257                 p           }
258                 p           }
259                 p           }
260                 p           }
261                 p           }
262                 p           }
263                 p           }
264                 p           }
265                 p           }
266                 p           }
267                 p           }
268                 p           }
269                 p           }
270                 p           }
271                 p           }
272                 p           }
273                 p           }
274                 p           }
275                 p           }
276                 p           }
277                 p           }
278                 p           }
279                 p           }
280                 p           }
281                 p           }
282                 p           }
283                 p           }
284                 p           }
285                 p           }
286                 p           }
287                 p           }
288                 p           }
289                 p           }
290                 p           }
291                 p           }
292                 p           }
293                 p           }
294                 p           }
295                 p           }
296                 p           }
297                 p           }
298                 p           }
299                 p           }
300                 p           }
301                 p           }
302                 p           }
303                 p           }
304                 p           }
305                 p           }
306                 p           }
307                 p           }
308                 p           }
309                 p           }
310                 p           }
311                 p           }
312                 p           }
313                 p           }
314                 p           }
315                 p           }
316                 p           }
317                 p           }
318                 p           }
319                 p           }
320                 p           }
321                 p           }
322                 p           }
323                 p           }
324                 p           }
325                 p           }
326                 p           }
327                 p           }
328                 p           }
329                 p           }
330                 p           }
331                 p           }
332                 p           }
333                 p           }
334                 p           }
335                 p           }
336                 p           }
337                 p           }
338                 p           }
339                 p           }
340                 p           }
341                 p           }
342                 p           }
343                 p           }
344                 p           }
345                 p           }
346                 p           }
347                 p           }
348                 p           }
349                 p           }
350                 p           }
351                 p           }
352                 p           }
353                 p           }
354                 p           }
355                 p           }
356                 p           }
357                 p           }
358                 p           }
359                 p           }
360                 p           }
361                 p           }
362                 p           }
363                 p           }
364                 p           }
365                 p           }
366                 p           }
367                 p           }
368                 p           }
369                 p           }
370                 p           }
371                 p           }
372                 p           }
373                 p           }
374                 p           }
375                 p           }
376                 p           }
377                 p           }
378                 p           }
379                 p           }
380                 p           }
381                 p           }
382                 p           }
383                 p           }
384                 p           }
385                 p           }
386                 p           }
387                 p           }
388                 p           }
389                 p           }
390                 p           }
391                 p           }
392                 p           }
393                 p           }
394                 p           }
395                 p           }
396                 p           }
397                 p           }
398                 p           }
399                 p           }
400                 p           }
401                 p           }
402                 p           }
403                 p           }
404                 p           }
405                 p           }
406                 p           }
407                 p           }
408                 p           }
409                 p           }
410                 p           }
411                 p           }
412                 p           }
413                 p           }
414                 p           }
415                 p           }
416                 p           }
417                 p           }
418                 p           }
419                 p           }
420                 p           }
421                 p           }
422                 p           }
423                 p           }
424                 p           }
425                 p           }
426                 p           }
427                 p           }
428                 p           }
429                 p           }
430                 p           }
431                 p           }
432                 p           }
433                 p           }
434                 p           }
435                 p           }
436                 p           }
437                 p           }
438                 p           }
439                 p           }
440                 p           }
441                 p           }
442                 p           }
443                 p           }
444                 p           }
445                 p           }
446                 p           }
447                 p           }
448                 p           }
449                 p           }
450                 p           }
451                 p           }
452                 p           }
453                 p           }
454                 p           }
455                 p           }
456                 p           }
457                 p           }
458                 p           }
459                 p           }
460                 p           }
461                 p           }
462                 p           }
463                 p           }
464                 p           }
465                 p           }
466                 p           }
467                 p           }
468                 p           }
469                 p           }
470                 p           }
471                 p           }
472                 p           }
473                 p           }
474                 p           }
475                 p           }
476                 p           }
477                 p           }
478                 p           }
479                 p           }
480                 p           }
481                 p           }
482                 p           }
483                 p           }
484                 p           }
485                 p           }
486                 p           }
487                 p           }
488                 p           }
489                 p           }
490                 p           }
491                 p           }
492                 p           }
493                 p           }
494                 p           }
495                 p           }
496                 p           }
497                 p           }
498                 p           }
499                 p           }
500                 p           }
501                 p           }
502                 p           }
503                 p           }
504                 p           }
505                 p           }
506                 p           }
507                 p           }
508                 p           }
509                 p           }
510                 p           }
511                 p           }
512                 p           }
513                 p           }
514                 p           }
515                 p           }
516                 p           }
517                 p           }
518                 p           }
519                 p           }
520                 p           }
521                 p           }
522                 p           }
523                 p           }
524                 p           }
525                 p           }
526                 p           }
527                 p           }
528                 p           }
529                 p           }
530                 p           }
531                 p           }
532                 p           }
533                 p           }
534                 p           }
535                 p           }
536                 p           }
537                 p           }
538                 p           }
539                 p           }
540                 p           }
541                 p           }
542                 p           }
543                 p           }
544                 p           }
545                 p           }
546                 p           }
547                 p           }
548                 p           }
549                 p           }
550                 p           }
551                 p           }
552                 p           }
553                 p           }
554                 p           }
555                 p           }
556                 p           }
557                 p           }
558                 p           }
559                 p           }
560                 p           }
561                 p           }
562                 p           }
563                 p           }
564                 p           }
565                 p           }
566                 p           }
567                 p           }
568                 p           }
569                 p           }
570                 p           }
571                 p           }
572                 p           }
573                 p           }
574                 p           }
575                 p           }
576                 p           }
577                 p           }
578                 p           }
579                 p           }
580                 p           }
581                 p           }
582                 p           }
583                 p           }
584                 p           }
585                 p           }
586                 p           }
587                 p           }
588                 p           }
589                 p           }
590                 p           }
591                 p           }
592                 p           }
593                 p           }
594                 p           }
595                 p           }
596                 p           }
597                 p           }
598                 p           }
599                 p           }
600                 p           }
601                 p           }
602                 p           }
603                 p           }
604                 p           }
605                 p           }
606                 p           }
607                 p           }
608                 p           }
609                 p           }
610                 p           }
611                 p           }
612                 p           }
613                 p           }
614                 p           }
615                 p           }
616                 p           }
617                 p           }
618                 p           }
619                 p           }
620                 p           }
621                 p           }
622                 p           }
623                 p           }
624                 p           }
625                 p           }
626                 p           }
627                 p           }
628                 p           }
629                 p           }
630                 p           }
631                 p           }
632                 p           }
633                 p           }
634                 p           }
635                 p           }
636                 p           }
637                 p           }
638                 p           }
639                 p           }
640                 p           }
641                 p           }
642                 p           }
643                 p           }
644                 p           }
645                 p           }
646                 p           }
647                 p           }
648                 p           }
649                 p           }
650                 p           }
651                 p           }
652                 p           }
653                 p           }
654                 p           }
655                 p           }
656                 p           }
657                 p           }
658                 p           }
659                 p           }
660                 p           }
661                 p           }
662                 p           }
663                 p           }
664                 p           }
665                 p           }
666                 p           }
667                 p           }
668                 p           }
669                 p           }
670                 p           }
671                 p           }
672                 p           }
673                 p           }
674                 p           }
675                 p           }
676                 p           }
677                 p           }
678                 p           }
679                 p           }
680                 p           }
681                 p           }
682                 p           }
683                 p           }
684                 p           }
685                 p           }
686                 p           }
687                 p           }
688                 p           }
689                 p           }
690                 p           }
691                 p           }
692                 p           }
693                 p           }
694                 p           }
695                 p           }
696                 p           }
697                 p           }
698                 p           }
699                 p           }
700                 p           }
701                 p           }
702                 p           }
703                 p           }
704                 p           }
705                 p           }
706                 p           }
707                 p           }
708                 p           }
709                 p           }
710                 p           }
711                 p           }
712                 p           }
713                 p           }
714                 p           }
715                 p           }
716                 p           }
717                 p           }
718                 p           }
719                 p           }
720                 p           }
721                 p           }
722                 p           }
723                 p           }
724                 p           }
725                 p           }
726                 p           }
727                 p           }
728                 p           }
729                 p           }
730                 p           }
731                 p           }
732                 p           }
733                 p           }
734                 p           }
735                 p           }
736                 p           }
737                 p           }
738                 p           }
739                 p           }
740                 p           }
741                 p           }
742                 p           }
743                 p           }
744                 p           }
745                 p           }
746                 p           }
747                 p           }
748                 p           }
749                 p           }
750                 p           }
751                 p           }
752                 p           }
753                 p           }
754                 p           }
755                 p           }
756                 p           }
757                 p           }
758                 p           }
759                 p           }
760                 p           }
761                 p           }
762                 p           }
763                 p           }
764                 p           }
765                 p           }
766                 p           }
767                 p           }
768                 p           }
769                 p           }
770                 p           }
771                 p           }
772                 p           }
773                 p           }
774                 p           }
775                 p           }
776                 p           }
777                 p           }
778                 p           }
779                 p           }
780                 p           }
781                 p           }
782                 p           }
783                 p           }
784                 p           }
785                 p           }
786                 p           }
787                 p           }
788                 p           }
789                 p           }
790                 p           }
791                 p           }
792                 p           }
793                 p           }
794                 p           }
795                 p           }
796                 p           }
797                 p           }
798                 p           }
799                 p           }
800                 p           }
801                 p           }
802                 p           }
803                 p           }
804                 p           }
805                 p           }
806                 p           }
807                 p           }
808                 p           }
809                 p           }
810                 p           }
811                 p           }
812                 p           }
813                 p           }
814                 p           }
815                 p           }
816                 p           }
817                 p           }
818                 p           }
819                 p           }
820                 p           }
821                 p           }
822                 p           }
823                 p           }
824                 p           }
825                 p           }
826                 p           }
827                 p           }
828                 p           }
829                 p           }
830                 p           }
831                 p           }
832                 p           }
833                 p           }
834                 p           }
835                 p           }
836                 p           }
837                 p           }
838                 p           }
839                 p           }
840                 p           }
841                 p           }
842                 p           }
843                 p           }
844                 p           }
845                 p           }
846                 p           }
847                 p           }
848                 p           }
849                 p           }
850                 p           }
851                 p           }
852                 p           }
853                 p           }
854                 p           }
855                 p           }
856                 p           }
857                 p           }
858                 p           }
859                 p           }
860                 p           }
861                 p           }
862                 p           }
863                 p           }
864                 p           }
865                 p           }
866                 p           }
867                 p           }
868                 p           }
869                 p           }
870                 p           }
871                 p           }
872                 p           }
873                 p           }
874                 p           }
875                 p           }
876                 p           }
877                 p           }
878                 p           }
879                 p           }
880                 p           }
881                 p           }
882                 p           }
883                 p           }
884                 p           }
885                 p           }
886                 p           }
887                 p           }
888                 p           }
889                 p           }
890                 p           }
891                 p           }
892                 p           }
893                 p           }
894                 p           }
895                 p           }
896                 p           }
897                 p           }
898                 p           }
899                 p           }
900                 p           }
901                 p           }
902                 p           }
903                 p           }
904                 p           }
905                 p           }
906                 p           }
907                 p           }
908                 p           }
909                 p           }
910                 p           }
911                 p           }
912                 p           }
913                 p           }
914                 p           }
915                 p           }
916                 p           }
917                 p           }
918                 p           }
919                 p           }
920                 p           }
921                 p           }
922                 p           }
923                 p           }
924                 p           }
925                 p           }
926                 p           }
927                 p           }
928                 p           }
929                 p           }
930                 p           }
931                 p           }
932                 p           }
933                 p           }
934                 p           }
935                 p           }
936                 p           }
937                 p           }
938                 p           }
939                 p           }
940                 p           }
941                 p           }
942                 p           }
943                 p           }
944                 p           }
945                 p           }
946                 p           }
947                 p           }
948                 p           }
949                 p           }
950                 p           }
951                 p           }
952                 p           }
953                 p           }
954                 p           }
955                 p           }
956                 p           }
957                 p           }
958                 p           }
959                 p           }
960                 p           }
961                 p           }
962                 p           }
963                 p           }
964                 p           }
965                 p           }
966                 p           }
967                 p           }
968                 p           }
969                 p           }
970                 p           }
971                 p           }
972                 p           }
973                 p           }
974                 p           }
975                 p           }
976                 p           }
977                 p           }
978                 p           }
979                 p           }
980                 p           }
981                 p           }
982                 p           }
983                 p           }
984                 p           }
985                 p           }
986                 p           }
987                 p           }
988                 p           }
989                 p           }
990                 p           }
991                 p           }
992                 p           }
993                 p           }
994                 p           }
995                 p           }
996                 p           }
997                 p           }
998                 p           }
999                 p           }
1000                p           }
1001               p           }
1002              p           }
1003             p           }
1004            p           }
1005           p           }
1006          p           }
1007         p           }
1008        p           }
1009       p           }
1010      p           }
1011     p           }
1012    p           }
1013   p           }
1014  p           }
1015 
```

最内ループの繰り返し数Nが
ITR:192よりも少ないため
ソフトウェアパイプラインングの
処理ルートに入らない

N=8



Statistics	GFLOPS	Effective instruction
改善前	33.39	4.87E+10

CLONEを利用し最内ループの繰り返し数をSIMD長に固定しました。その結果、外側のループでソフトウェアパイプラインングが促進され、演算効率・有効命令数が改善しました。

改善後ソース (ソースチューニング)

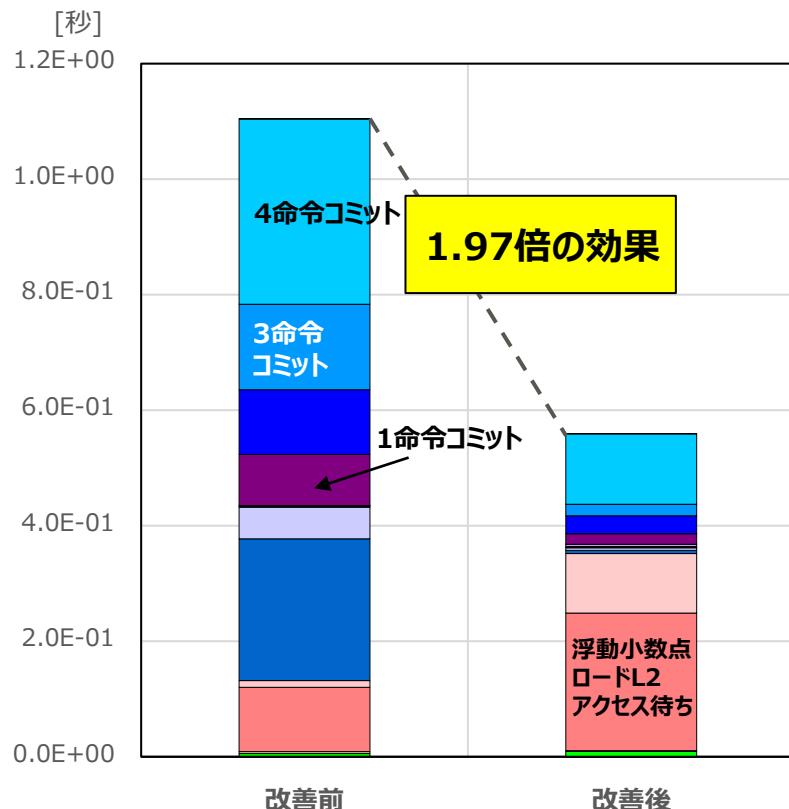
```

37 #pragma omp parallel private(iter,i,j,k)
38 {
39     for (iter = 0; iter < itmax; iter++){
40         #pragma omp for nowait
41         #pragma loop clone NN==8
42         <<< Loop-information Start >>>
43         <<< [OPTIMIZATION]
44         <<< CLONE
45         <<< PREFETCH(HARD) E
46         <<< (unknown)
47         <<< Loop-information End
48         for (k = 0; k < L; k++){
49             <<< Loop-information Start >>>
50             <<< [OPTIMIZATION]
51             <<< SOFTWARE PIPELINING(IPC: 3.25, ITR: 176,
52                 MVE: 6, POL: S)
53             <<< PREFETCH(HARD) Expected by compiler :
54             <<< (unknown)
55             <<< Loop-information End >>>
56             for (j = 0; j < M; j++){
57                 <<< Loop-information Start >>>
58                 <<< [OPTIMIZATION]
59                 <<< SIMD(VL: 8)
60                 <<< SOFTWARE PIPELINING(IPC: 3.00, ITR: 192,
61                     MVE: 7, POL: S)
62                 <<< PREFETCH(HARD) Expected by compiler :
63                 <<< (unknown)
64                 <<< Loop-information End >>>
65                 for (i = 0; i < NN; i++){
66                     a[k][j][i] = a[k][j][i] + c * b[k][j][i];
67                 }
68             }
69         }
70     }
71 }
```

ソフトウェアパイプラインングが適用される

M=240

最内ループの繰り返し数をSIMD長に固定



Statistics	GFLOPS	Effective instruction
改善前	33.39	4.87E+10
改善後	65.89	1.50E+11

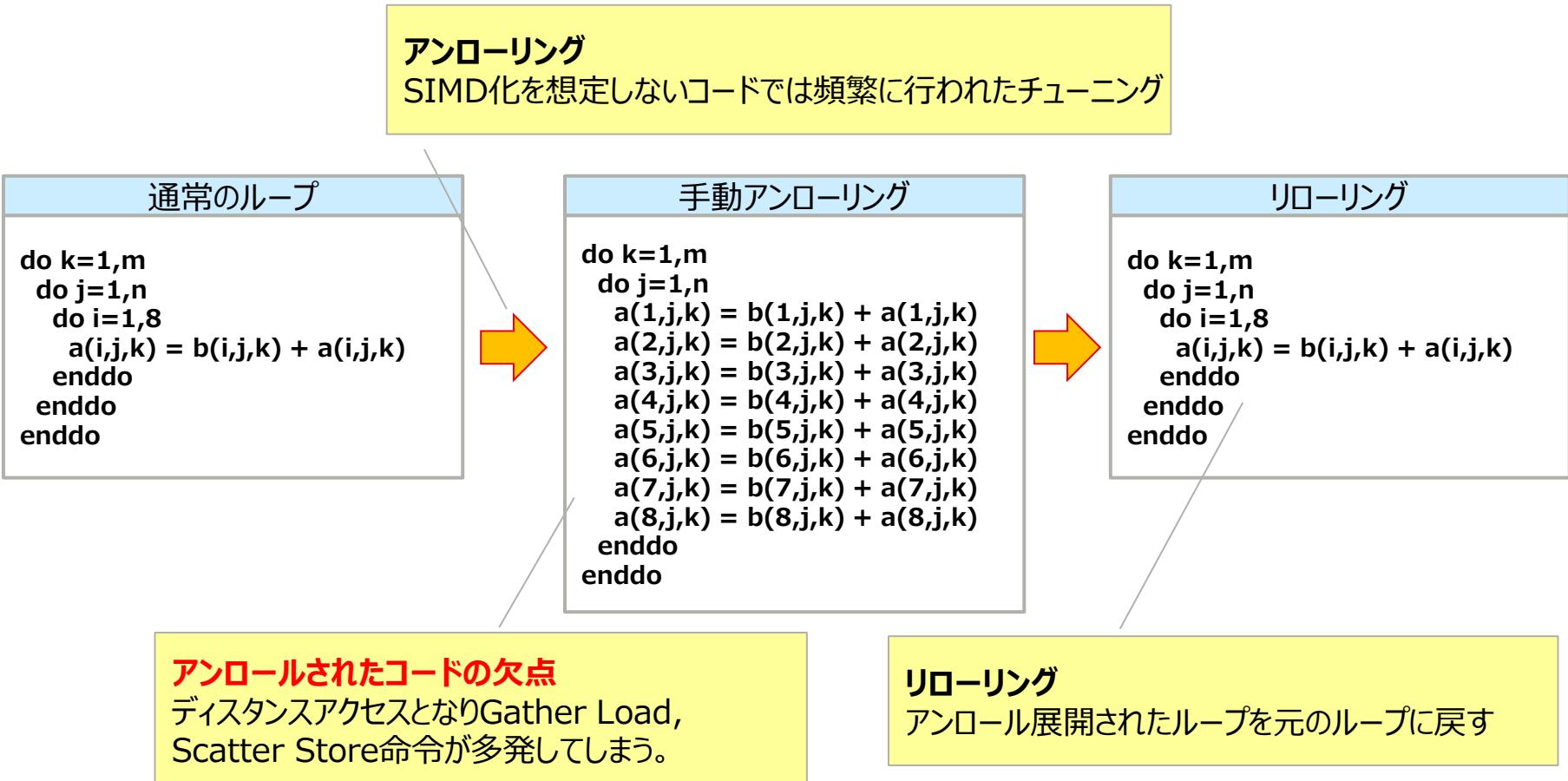
リローリング

- ・ リローリングとは
- ・ リローリング（改善前）
- ・ リローリング（ソースチューニング）

リローリングとは

FUJITSU

リローリングとはアンローリング展開された文をループ文に戻すことで、ループに対する最適化を促進するチューニングです。



手動アンロール展開されています。Gather Load, Scatter Storeが多発し、浮動小数点ロードL1Dアクセス待ちが多くなっています。

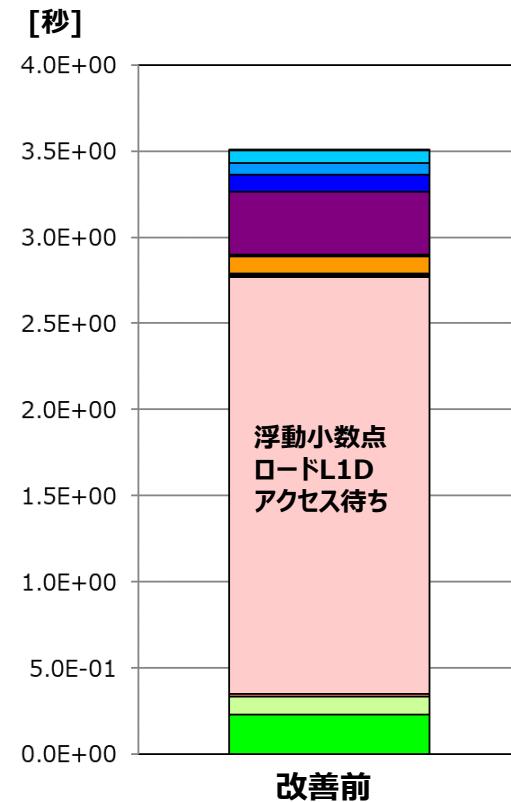
改善前ソース

```

40      real(8) :: a(8,N,M),b(8,N,M)
41
42      <<< Loop-information Start >>>
43      <<< [PARALLELIZATION]
44      <<< Standard iteration count: 2
45      <<< [OPTIMIZATION]
46      <<< PREFETCH(HARD) Expected by compiler :
47      <<< b, a
48      <<< Loop-information End >>>
49      1 pp      DO K=1,M
50          <<< Loop-information Start >>>
51          <<< [OPTIMIZATION]
52          <<< SIMD(VL: 8)
53          <<< SOFTWARE PIPELINING(IPC: 2.04, ITR: 40,
54                           MVE: 3, POL: S)
55          <<< PREFETCH(HARD) Expected by compiler :
56          <<< b, a
57          <<< Loop-information End >>>
58
59      2 p v      DO J=1,N
60      2 p v          a(1,J,K) = b(1,J,K) + a(1,J,K)
61      2 p v          a(2,J,K) = b(2,J,K) + a(2,J,K)
62      2 p v          a(3,J,K) = b(3,J,K) + a(3,J,K)
63      2 p v          a(4,J,K) = b(4,J,K) + a(4,J,K)
64      2 p v          a(5,J,K) = b(5,J,K) + a(5,J,K)
65      2 p v          a(6,J,K) = b(6,J,K) + a(6,J,K)
66      2 p v          a(7,J,K) = b(7,J,K) + a(7,J,K)
67      2 p v          a(8,J,K) = b(8,J,K) + a(8,J,K)
68
69      2 p v          ENDDO
70
71      1 p          ENDDO

```

N=128
M=250



リローリングを行う（ループに戻す）ことでデータが連続アクセスとなり、効果的なソフトウェアパイプラインング、SIMD化、ループアンロールなどの最適化が促進されました。

改善後ソース（ソースチューニング）

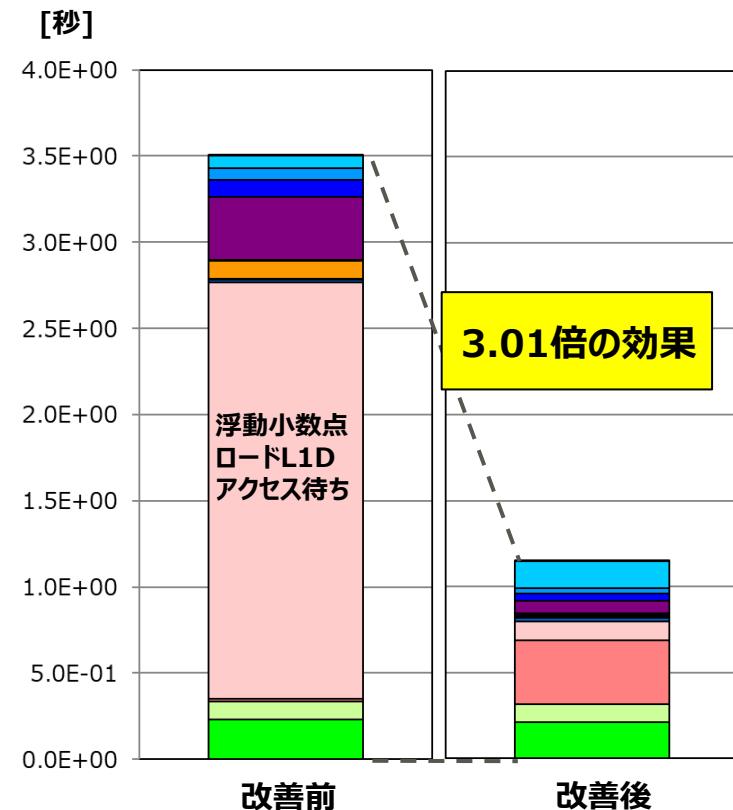
```

40      real(8) :: a(8,N,M),b(8,N,M)
41
42      <<< Loop-information Start >>>
43      <<< [PARALLELIZATION]
44      <<< Standard iteration c
45      <<< [OPTIMIZATION]
46      <<< COLLAPSED
47      <<< SIMD(VL: 8)
48      <<< SOFTWARE PIPELINING(IPC: 3.00, ITR: 192,
49          MVE: 7, POL: S)
50      <<< PREFETCH(HARD) Expected by compiler :
51      <<< a, b
52      <<< Loop-information End >>>
53      1 pp 2v      DO K=1,M
54      <<< Loop-information Start >>>
55      <<< [OPTIMIZATION]
56      <<< COLLAPSED
57      <<< Loop-information End >>>
58      2 p  2       DO J=1,N
59      <<< Loop-information Start >>>
60      <<< [OPTIMIZATION]
61      <<< COLLAPSED
62      <<< Loop-information End >>>
63      3 p  2       DO I=1,8
64      3 p  2v      a(I,J,K) = b(I,J,K) + a(I,J,K)
65      3 p  2v      ENDDO
66      2 p          ENDDO
67      1 p          ENDDO

```

ループの1重化

SIMD化・
ループアンロールが
促進されている



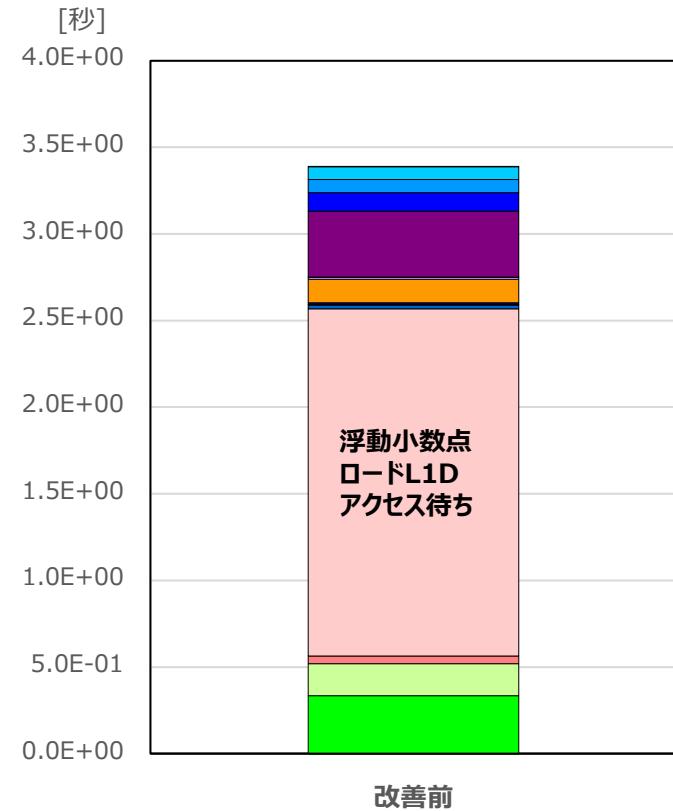
手動アンロール展開されています。Gather Load, Scatter Storeが多発し、浮動小数点ロードL1Dアクセス待ちが多くなっています。

改善前ソース

```

50    #pragma omp parallel for private(j)
      <<< Loop-information Start >>>
      <<< [OPTIMIZATION]
      <<< PREFETCH(HARD) Expected by compiler :
      <<< (unknown)
      <<< Loop-information End >>>
51  p    for( k=0;k<M; k++)
52  p    {
      <<< Loop-information Start >>>
      <<< [OPTIMIZATION]
      <<< SIMD(VL: 8)
      <<< SOFTWARE PIPELINING(IPC: 2.04, ITR: 40,
                                MVE: 3, POL: S)
      <<< PREFETCH(HARD) Expected by compiler :
      <<< (unknown)
      <<< Loop-information End >>>
53  p    v    for(j=0; j<N; j++)
54  p    v    {
55  p    v      a[k][j][0] = b[k][j][0] + a[k][j][0];
56  p    v      a[k][j][1] = b[k][j][1] + a[k][j][1];
57  p    v      a[k][j][2] = b[k][j][2] + a[k][j][2];
58  p    v      a[k][j][3] = b[k][j][3] + a[k][j][3];
59  p    v      a[k][j][4] = b[k][j][4] + a[k][j][4];
60  p    v      a[k][j][5] = b[k][j][5] + a[k][j][5];
61  p    v      a[k][j][6] = b[k][j][6] + a[k][j][6];
62  p    v      a[k][j][7] = b[k][j][7] + a[k][j][7];
63  p    v    }
64  p    }
65
66    return;
67  }
```

N=128
M=250



リローリングを行う（ループに戻す）ことでデータが連続アクセスとなり、効果的なソフトウェアパイプラインング、SIMD化、ループアンロールなどの最適化が促進されました。

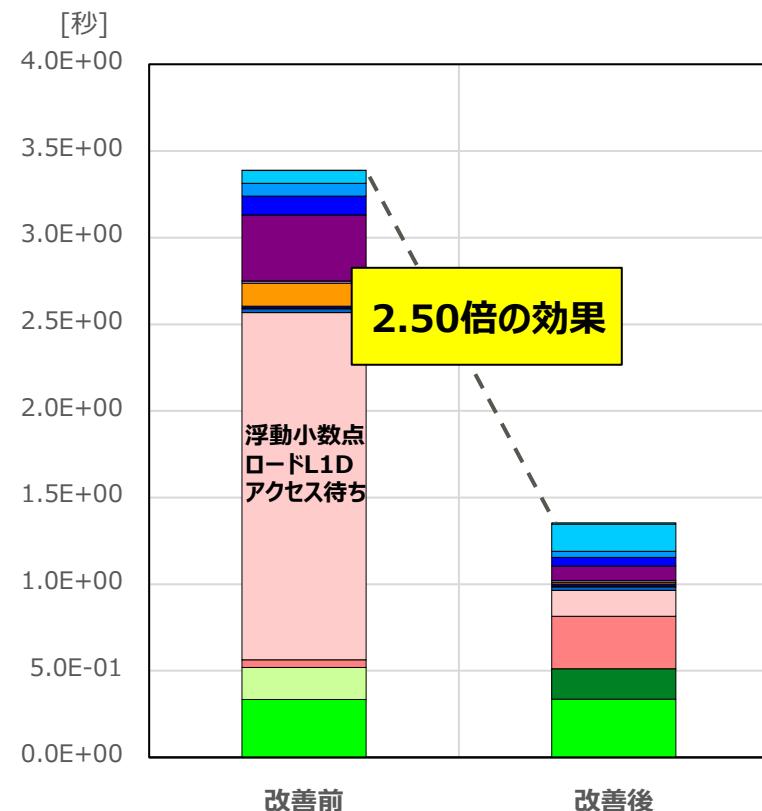
改善後ソース（ソースチューニング）

```

50    #pragma omp parallel for private(i,j) collapse(3)
      <<< Loop-information Start >>>
      <<< [OPTIMIZATION]
      <<< SIMD(VL: 8)
      <<< SOFTWARE PIPELINING(IPC: 3.00, ITR: 192,
                                MVE: 7, POL: S)
      <<< PREFETCH(HARD) Expected by compiler :
      <<< (unknown)
      <<< Loop-information End >>>
51 p  2v  for (k=0; k< M; k++)
52   2  {
53   p  2  for(j=0; j<N; j++)
54   2  {
55   p  2  for(i=0; i<8; i++)
56   p  2v  {
57   p  2v  a[k][j][i] = b[k][j][i] + a[k][j][i];
58   p  2v  }
59   p  2v  }
60   p  2v  }
61
62     return;
63 }
```

SIMD化・
ループアンロールが
促進されている

ループの1重化



ループアンスイッチング

- ループアンスイッチングとは
- ループアンスイッチングの効果（改善前）
- ループアンスイッチングの効果（改善後）

ループアンスイッチングとは

FUJITSU

ループ内に不变な状態で分岐するIF構文が存在する場合に、そのIF構文をループの外側に出し、IF構文の条件が成立した場合のループと、成立しない場合のループを作成する最適化です。

ソースコード

```
!$omp do
  do i=1,n1
!ocl unswitching
  if (n1 >= q) then
    处理①
  endif
!ocl unswitching
  if(n1 > r) then
    处理②
  endif
!ocl unswitching
  if(n1 < s) then
    处理③
  endif
enddo
 !$omp enddo
```

最適化イメージ	
!パターン(1) if((条件①真).and.(条件②真).and.(条件③真))then do i=1,n1 处理① 处理② 处理③ enddo endif	!パターン(5) if((条件①偽).and.(条件②真).and.(条件③真))then do i=1,n1 处理② 处理③ enddo endif
!パターン(2) if((条件①真).and.(条件②真).and.(条件③偽))then do i=1,n1 处理① 处理② enddo endif	!パターン(6) if((条件①偽).and.(条件②真).and.(条件③偽))then do i=1,n1 处理② enddo endif
!パターン(3) if((条件①真).and.(条件②偽).and.(条件③真))then do i=1,n1 处理① 处理③ enddo endif	!パターン(7) if((条件①偽).and.(条件②偽).and.(条件③真))then do i=1,n1 处理③ enddo endif
!パターン(4) if((条件①真).and.(条件②偽).and.(条件③偽))then do i=1,n1 处理① enddo endif	!パターン(8) if((条件①偽).and.(条件②偽).and.(条件③偽))then do i=1,n1 enddo endif

8つのif文(do文)に展開される

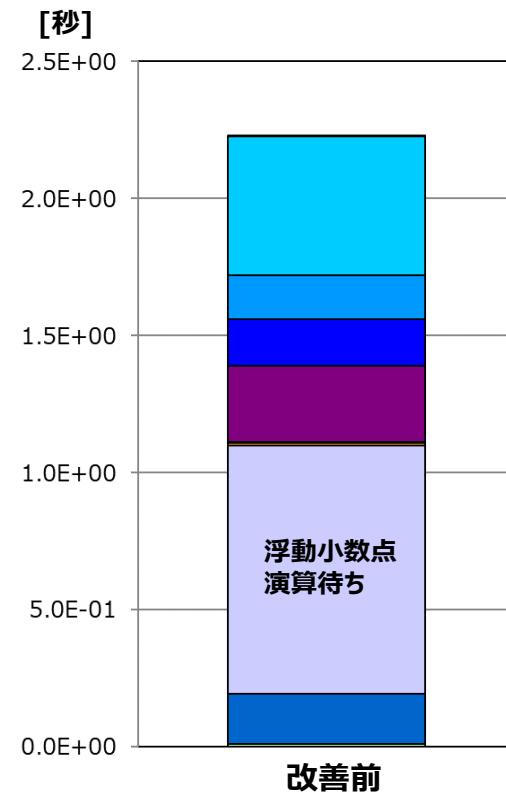
最内ループにIF構文があり、SIMD化やソフトウェアパイプラインが効果的に行われず、浮動小数点演算待ちが多くなっています。

改善前ソース

```

97   1      !$omp do
98     <<< Loop-information Start >>>
99       <<< [OPTIMIZATION]
100      <<< SIMD(VL: 8)
101      <<< SOFTWARE PIPELINING(IPC: 1.31, ITR: 96,
102          MVE: 2, POL: L)
103      <<< UNSWITCHING
104      <<< PREFETCH(HARD) Expected by compiler :
105          a, b
106      <<< Loop-information End >>>
107
108     2 p 2v    do i=1,n1
109
110     3 p 2v    if (n1 >= q) then
111       3 p 2v        a(i) = c0+b(i)*(c1+b(i)*(c2+b(i)*(c3+b(i)*c4)))
112       3 p 2v        endif
113
114     3 p 2v    if(n1 > r) then
115       3 p 2v        a(i) = c0*b(i)/(c1*b(i)/(c2*b(i)/(c3*b(i)/c4)))
116       3 p 2v        endif
117
118     3 p 2v    if(n1 < s) then
119       3 p 2v        a(i) = c0+b(i)/(c1+b(i)/(c2+b(i)/(c3+b(i)/c4)))
120       3 p 2v        endif
121     2 p 2v    enddo
122   1      !$omp enddo nowait

```



SIMD

	Effective instruction	SIMD instruction rate (%) (/Effective instruction)
改善前	7.47E+10	94.18%

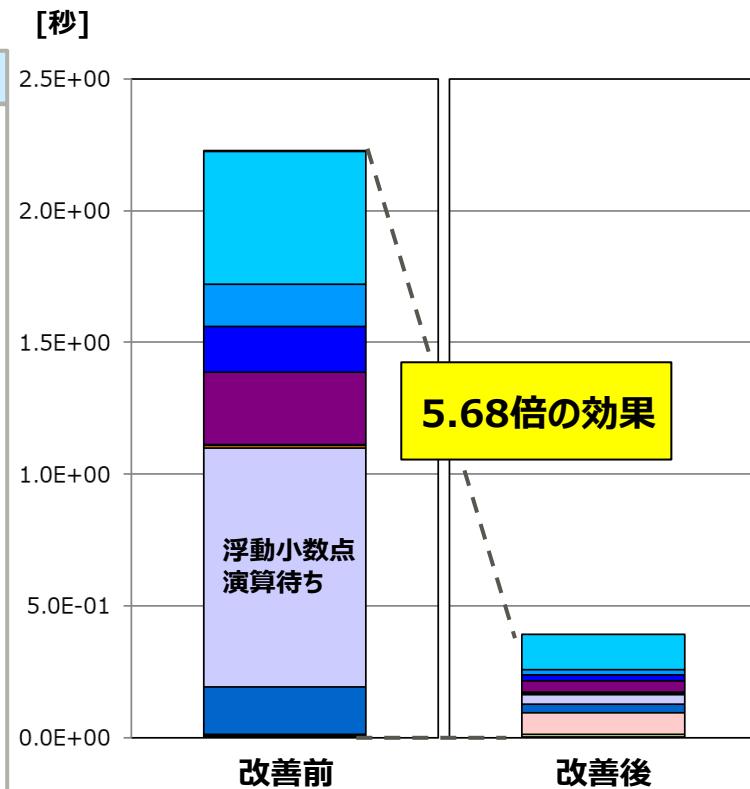
IF構文に対してループアンスイッチングを指示することで分岐がなくなり、SIMD化およびソフトウェアパイプラインングが促進されました。その結果、浮動小数点演算待ちが大幅に改善されました。

改善後ソース（最適化制御行チューニング）

```

97   1      !$omp do
<<< Loop-information Start >>>
<<< [OPTIMIZATION]
<<< SIMD(VL: 8)
<<< SOFTWARE PIPELINING(IPC: 1.04, ITR: 80,
                           MVE: 3, POL: S)
<<< UNSWITCHING
<<< PREFETCH(HARD) Expected by compiler :
<<< b, a
<<< Loop-information End >>>
98   2 p 2v  do i=1,n1
99   2      !ocl unswitching
100  3 p 2v  if (n1 >= q) then
101  3 p 2v  a(i) = c0+b(i)*(c1+b(i)*(c2+b(i)*(c3+b(i)*c4)))
102  3 p 2v  endif
103  2      !ocl unswitching
104  3 p 2v  if(n1 > r) then
105  3 p 2v  a(i) = c0*b(i)/(c1*b(i)/(c2*b(i)/(c3*b(i)/c4)))
106  3 p 2v  endif
107  2      !ocl unswitching
108  3 p 2v  if(n1 < s) then
109  3 p 2v  a(i) = c0+b(i)/(c1+b(i)/(c2+b(i)/(c3+b(i)/c4)))
110  3 p 2v  endif
111  2 p 2v  enddo
112  1      !$omp enddo nowait

```



SIMD

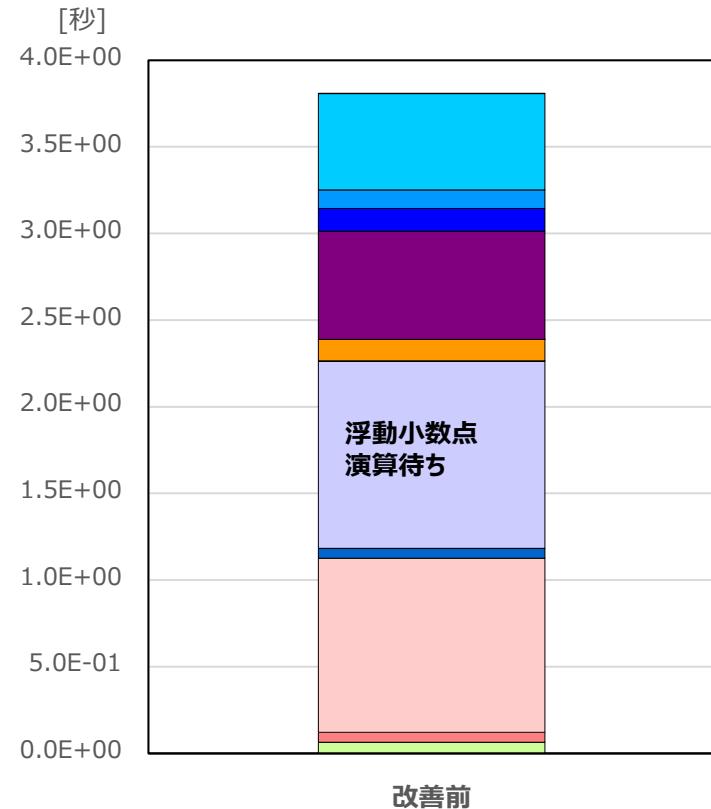
	Effective instruction	SIMD instruction rate (%) (/Effective instruction)
改善前	7.47E+10	94.18%
改善後	1.60E+10	75.83%

最内ループにif文があり、SIMD化やソフトウェアパイプラインングが効果的に行われず、浮動小数点演算待ちが多くなっています。

改善前ソース

```

91      #pragma omp for nowait
<<< Loop-information Start >>>
<<< [OPTIMIZATION]
<<< SIMD(VL: 8)
<<< UNSWITCHING
<<< PREFETCH(HARD) Expected by compiler :
<<< (unknown)
<<< Loop-information End >>>
92 p 2v   for(i=0; i<n1; i++)
93 p 2v   {
94 p 2v     if (n1 >= q)
95 p 2v     {
96 p 2v       a[i] = c0+b[i]*(c0+b[i]*(c0+b[i]*(c0+b[i]*c0)));
97 p 2v     }
98
99 p 2v     if(n1 > r)
100 p 2v     {
101 p 2v       a[i] = c0*b[i]/(c0*b[i]/(c0*b[i]/(c0*b[i]/c0)));
102 p 2v     }
103
104 p 2v     if(n1 < s)
105 p 2v     {
106 p 2v       a[i] = c0+b[i]/(c0+b[i]/(c0+b[i]/(c0+b[i]/c0)));
107 p 2v     }
108 p 2v   }
109 }
```



Statistics	Effective instruction	SIMD instruction rate (%) (/Effective instruction)
改善前	8.54E+10	89.49%

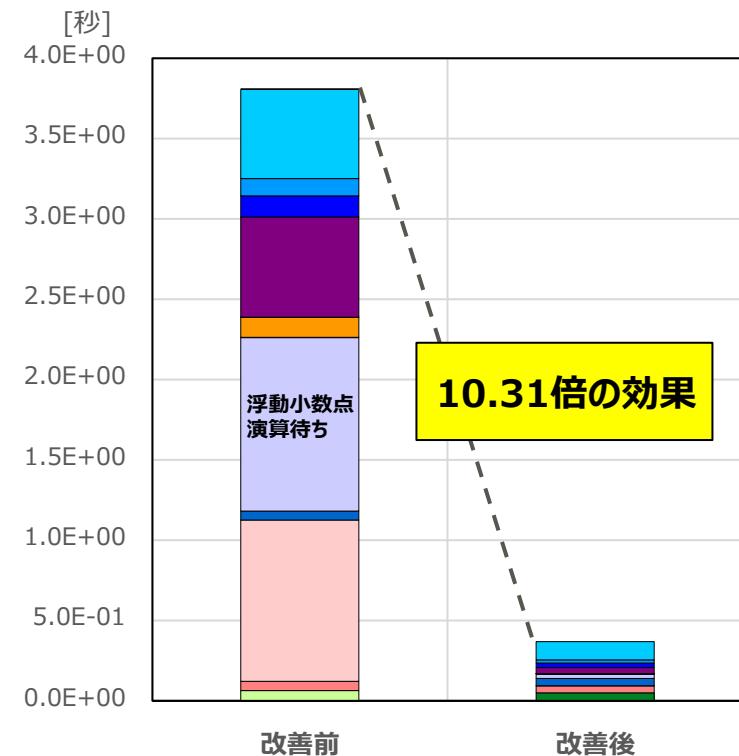
if文に対してループアンスイッチングを指示することで分岐がなくなり、SIMD化およびソフトウェアパイプラインングが促進されました。その結果、浮動小数点演算待ちが大幅に改善されました。

改善後ソース（最適化制御行チューニング）

```

89     for( j=0; j<iter; j++ )
90     {
91         #pragma omp for nowait
92         <<< Loop-information Start >>>
93         <<< [OPTIMIZATION]
94         <<< SIMD(VL: 8)
95         <<< SOFTWARE PIPELINING(IPC: 1.12, ITR: 96, MVE: 3, POL: S)
96         <<< UNSWITCHING
97         <<< PREFETCH(HARD) Expected by compiler :
98         <<< (unknown)
99         <<< Loop-information End >>>
100        p 2v    for(i=0; i<n1; i++)
101        p 2v    {
102        p 2v        #pragma statement unswitching
103        p 2v        if (n1 >= q)
104        p 2v        {
105        p 2v            a[i] = c0+b[i]*(c0+b[i]*(c0+b[i]*(c0+b[i]*c0)));
106        p 2v        }
107        p 2v        #pragma statement unswitching
108        p 2v        if(n1 > r)
109        p 2v        {
110        p 2v            a[i] = c0+b[i]/(c0*b[i]/(c0*b[i]/(c0*b[i]/c0)));
111        p 2v        }
112    }

```



Statistics	Effective instruction	SIMD instruction rate (%) (/Effective instruction)
改善前	8.54E+10	89.49%
改善後	1.68E+10	71.58%

マイクロアーキ依存の改善

- Scatter Store命令の回避
- Gather Load命令の纏め上げ機能
- 過剰SFIの改善
- Multiple Structures命令の活用
- ハードウェアプリフェッチの距離調整
- SVEのベクトルレジスタサイズ（SIMD幅）
- 半精度実数型の活用

Scatter Store命令の回避

- Scatter Store命令の回避（改善前）
- Scatter Store命令の回避（ソースチューニング）

Scatter Store（連続でないストア） 命令数が多い場合、性能が低下します。浮動小数点ロードL1Dアクセス待ちが多くなっています。

改善前ソース

```

42      !$omp parallel
43 1      DO K=1, ITER
44 1      !$omp do
45 <<< Loop-information Start >>>
46 <<< [OPTIMIZATION]
47 <<< PREFETCH(HARD) Expected by compiler :
48 <<< b
49 <<< Loop-information End >>>
50 2 p      DO J=1,M
51 <<< Loop-information Start >>>
52 <<< [OPTIMIZATION]
53 <<< SIMD(VL: 8)
54 <<< SOFTWARE PIPELINING(IPC: 2.60, ITR: 80,
55                               MVE: 3, POL: S)
56 <<< PREFETCH(HARD) Expected by compiler :
57 <<< b
58 <<< Loop-information End >>>
59 3 p 2v      DO I=1,M
60 <<< a(J,I) = b(I,J)
61 3 p 2v      ENDDO
62 2 p      ENDDO
63 1      !$omp end do nowait
64 1      ENDDO
65 1      !$omp end parallel

```

L1Dミスが高い



Store instruction									
SIMD									
	L1I miss rate (/Effective instruction)	Load-store instruction	L1D miss	L1D miss rate (/Load-store instruction)	L1D miss demand rate (%) (/L1D miss)	hardware prefetch rate (%) (/L1D miss)	software prefetch rate (%) (/L1D miss)	contiguous store instruction	non-contiguous structure store instruction
改善前	0.00	1.62E+09	1.50E+09	0.92	99.98%	0.02%	0.00%	1.60E+01	0.00E+00
									1.30E+08
									1.07E+04
									5.49E+03

Cache

	L1I miss rate (/Effective instruction)	Load-store instruction	L1D miss	L1D miss rate (/Load-store instruction)	L1D miss demand rate (%) (/L1D miss)	hardware prefetch rate (%) (/L1D miss)	software prefetch rate (%) (/L1D miss)	contiguous store instruction	non-contiguous structure store instruction	Floating-point register spill instruction	Predicate register spill instruction
改善前	0.00	1.62E+09	1.50E+09	0.92	99.98%	0.02%	0.00%	1.60E+01	0.00E+00	1.30E+08	1.07E+04

ループインデックスを変更することでScatter Store命令の発生を回避しました。
その結果、L1Dミスが大幅に改善されました。

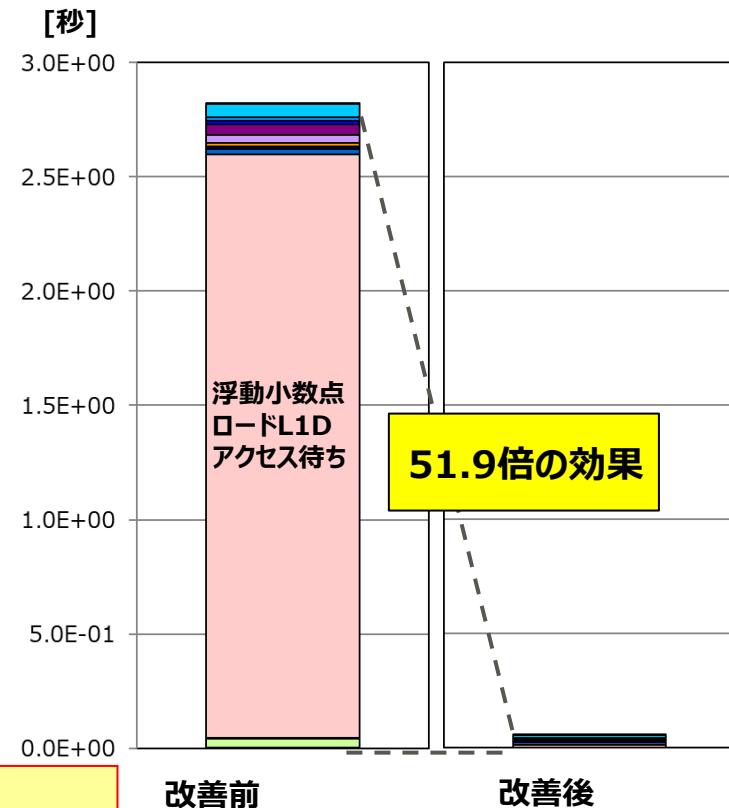
改善後ソース

```

42      !$omp parallel
43      1      DO K=1, ITER
44      1      !$omp do
45      2      p      <<< Loop-information Start >>>
46      3      p      4v      <<< [OPTIMIZATION]
47      3      p      4v      <<< PREFETCH(HARD) Expected by compiler :
48      3      p      4v      <<< a
49      2      p      <<< Loop-information End >>>
50      1      !$omp end do nowait
51      1      ENDDO
52      !$omp end parallel

```

L1Dミスが大幅に改善された

scatter store命令数を
低減できている

	L1I miss rate (/Effective instruction)	Load-store instruction	L1D miss	L1D miss rate (/Load-store instruction)	demand rate (%) (/L1D miss)	prefetch rate (%) (/L1D miss)	prefetch rate (%) (/L1D miss)
改善前	0.00	1.62E+09	1.50E+09	0.92	99.98%	0.02%	0.00%
改善後	0.00	3.64E+08	3.03E+06	0.01	99.97%	0.02%	0.01%

	Store instruction				
	SIMD				
	Single vector contiguous store instruction	Multiple vector contiguous store instruction	Non-contiguous scatter store instruction	Floating-point register spill instruction	Predicate register spill instruction
改善前	1.60E+01	0.00E+00	1.30E+08	1.07E+04	5.49E+03
改善後	1.30E+08	0.00E+00	0.00E+00	6.40E+01	2.38E+02

Scatter Store（連続でないストア）命令数が多い場合、性能が低下します。浮動小数点ロードL1Dアクセス待ちが多くなっています。

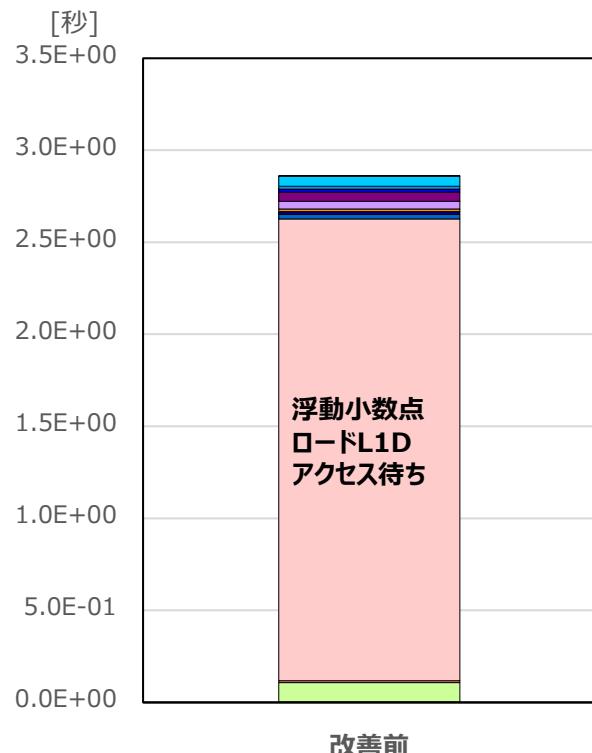
改善前ソース

```

44     #pragma omp parallel
45     {
46         for(k=0; k<iter; k++)
47         {
48             #pragma omp for nowait
49             p     <<< Loop-information Start >>>
50             p     <<< [OPTIMIZATION]
51             p     <<< PREFETCH(HARD) Expected by compiler :
52             p     <<< b
53             p     <<< Loop-information End >>>
54             p     for(j=0; j<m; j++)
55             p     {
56                 p     <<< Loop-information Start >>>
57                 p     <<< [OPTIMIZATION]
58                 p     <<< SIMD(VL: 8)
59                 p     <<< SOFTWARE PIPELINING(IPC: 2.40, ITR: 80,
60                                     MVE: 3, POL: S)
61                 p     <<< PREFETCH(HARD) Expected by compiler :
62                 p     <<< b
63                 p     <<< Loop-information End >>>
64                 p     2v     for(i=0;i<m;i++)
65                 p     2v     {
66                 p     2v     a[i][j] = b[j][i];
67                 p     2v     }
68                 p     }
69                 p     }
70                 p     }

```

L1Dミスが高い



scatter store命令数が他の命令に比べて多い

Cache	L1I miss rate (/Effective instruction)	Load-store instruction	L1D miss	L1D miss rate (/Load-store instruction)	L1D miss demand rate (%) (/L1D miss)	L1D miss hardware prefetch rate (%) (/L1D miss)	L1D miss software prefetch rate (%) (/L1D miss)	Single vector contiguous store instruction	Multiple vector contiguous structure store instruction	Non-contiguous scatter store instruction	Floating-point register spill instruction	Predicate register spill instruction
改善前	0.00	1.58E+09	1.50E+09	0.95	100.00%	0.00%	0.00%	0.00E+00	0.00E+00	1.30E+08	1.17E+04	5.88E+03

ループインデックスを変更することでScatter Store命令の発生を回避しました。
その結果、L1Dミスが大幅に改善されました。

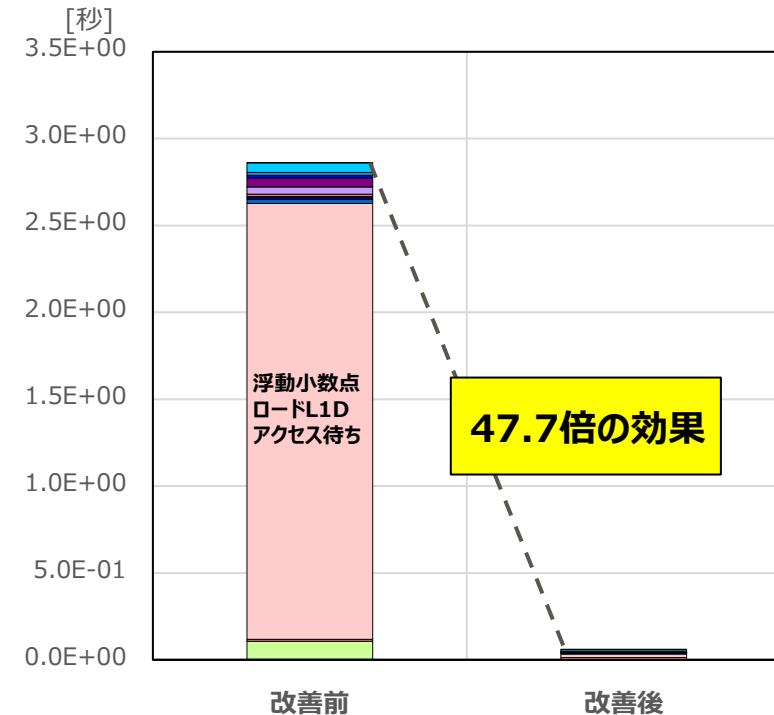
改善後ソース

```

44  #pragma omp parallel
45  {
46      for(k=0; k<iter; k++)
47  {
48      #pragma omp for nowait
49      p
50      p
51      p
52      p
53      p
54      p
55      p
56      }
57      }

    <<< Loop-information Start >>>
    <<< [OPTIMIZATION]
    <<< PREFETCH(HARD) Expected by compiler :
    <<< a
    <<< Loop-information End >>>
        for(j=0; j<m; j++)
    {
        <<< Loop-information Start >>>
        <<< [OPTIMIZATION]
        <<< SIMD(VL: 8)
        <<< SOFTWARE PIPELINING(IPC: 1.28, ITR: 96,
                                MVE: 2, POL: S)
        <<< PREFETCH(HARD) Expected by compiler :
        <<< a
        <<< Loop-information End >>>
4v            for(i=0;i<m;i++)
4v            {
4v                a[j][i] = b[i][j];
4v            }
    }
```

L1Dミスが大幅に改善された

scatter store命令数を
低減できている

Cache	L1I miss rate (/Effective instruction)	Load-store instruction	L1D miss	L1D miss rate (/Load-store instruction)	L1D miss demand rate (%) (/L1D miss)	L1D miss hardware prefetch rate (%) (/L1D miss)	L1D miss software prefetch rate (%) (/L1D miss)	Single vector contiguous store instruction	Multiple vector contiguous structure store instruction	Non-contiguous scatter store instruction	Floating-point register spill instruction	Predicate register spill instruction
改善前	0.00	1.58E+09	1.50E+09	0.95	100.00%	0.00%	0.00%	0.00E+00	0.00E+00	1.30E+08	1.17E+04	5.88E+03
改善後	0.00	4.13E+08	3.67E+06	0.01	99.35%	0.67%	-0.01%	1.30E+08	0.00E+00	0.00E+00	3.20E+01	1.70E+01

Gather Load命令の纏め上げの促進

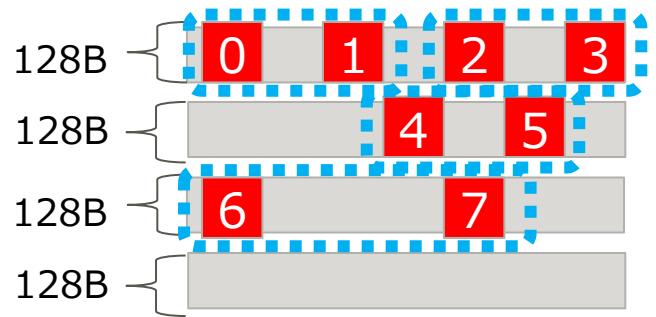
- Gather Load命令の纏め上げ機能について
- Gather Load命令の纏め上げの促進（改善前）
- Gather Load命令の纏め上げの促進（ソースチューニング）

■ Gather命令の纏め上げ機能とは

Gather命令は、1つのFPから同時発行される2要素が隣接する場合に高速処理ができます。隣接する2要素のアドレスが128バイト内で一致する場合には、まとめて1フローで処理することで高速化できます。

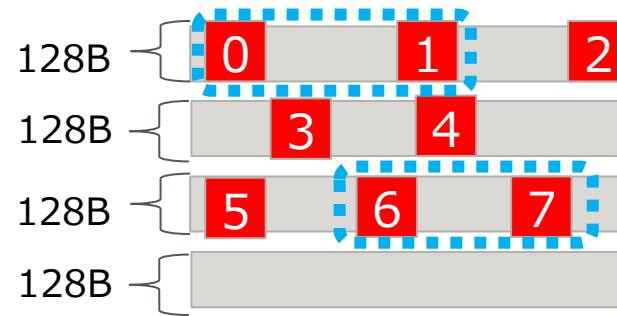
■ 隣接2要素が、同一の128バイトブロックに属する場合、それらを纏めて1フローで処理します。以下のアドレスパターン例で図の部分が纏め上げされて、2要素をL1D\$パイプライン1フローで処理します。

◆ アドレスパターン例1



纏め上げで1フロー処理

◆ アドレスパターン例2



Gather命令の纏め上げ機能を活用するために配列の先頭アドレスには注意して実装して下さい。

ストライドアクセスによりGather Load, Scatter Store命令が発生し、浮動小数点ロードL1Dアクセス待ちが多くなっています。

改善前ソース

```

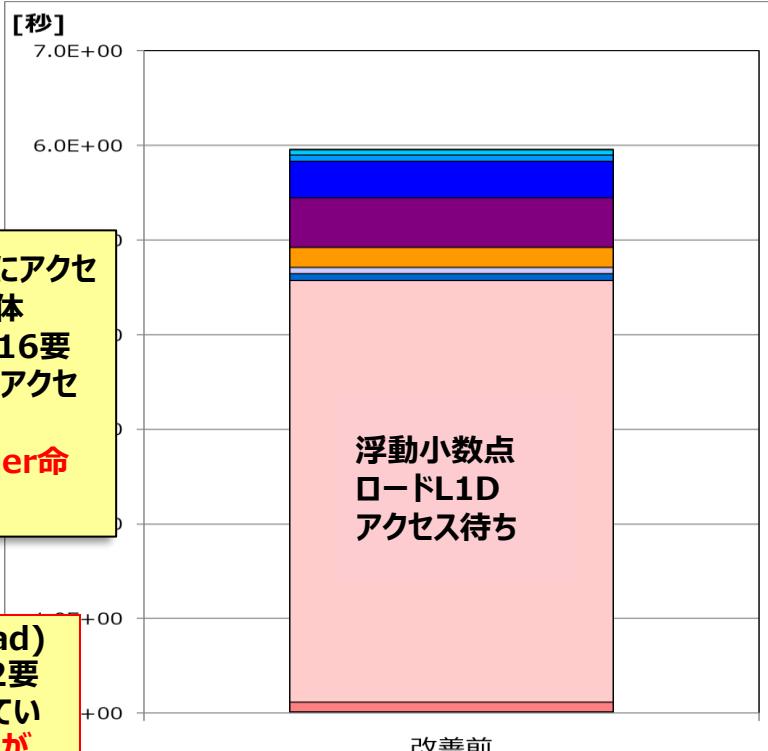
30      real(kind=8),dimension(n,m) :: a
31      real(kind=8),dimension(n,m) :: b,c
32
<<< Loop-information Start >>>
<<< [OPTIMIZATION]
<<< SIMD(VL: 8)
<<< SOFTWARE PIPELINING(IPC)
<<< Loop-information End >>>
33 1    v    do i = 1, m
34 1    v      a( 1,i) = b( 1,i) + c( 1,i)
35 1    v      a( 2,i) = b( 2,i) + c( 2,i)
36 1    v      a( 3,i) = b( 3,i) + c( 3,i)
37 1    v      a( 4,i) = b( 4,i) + c( 4,i)
38 1    v      a( 5,i) = b( 5,i) + c( 5,i)
39 1    v      a( 6,i) = b( 6,i) + c( 6,i)
40 1    v      a( 7,i) = b( 7,i) + c( 7,i)
41 1    v      a( 8,i) = b( 8,i) + c( 8,i)
42 1    v      a( 9,i) = b( 9,i) + c( 9,i)
43 1    v      a(10,i) = b(10,i) + c(10,i)
44 1    v      a(11,i) = b(11,i) + c(11,i)
45 1    v      a(12,i) = b(12,i) + c(12,i)
46 1    v      a(13,i) = b(13,i) + c(13,i)
47 1    v      a(14,i) = b(14,i) + c(14,i)
48 1    v      a(15,i) = b(15,i) + c(15,i)
49 1    v      a(16,i) = b(16,i) + c(16,i)
50 1    v    end do

```

n = 16

配列a, 配列b, 配列cは連続にアクセスされるが、それぞれの配列自体(a(1,i)など)は、1回転あたり16要素(128バイト)飛ぶストライドアクセスである
⇒離散アクセスのため、Gather命令が使用される

非連続ギャザーロード(Gather Load)命令が使用されているが、隣接する2要素のアドレスが128バイト以上離れているためGather命令の纏め上げ機能が動作せずL1ビジー率が高い



Instruction

	Load-store instruction	
	Load instruction	Store instruction
	Non-contiguous gather load instruction	Non-contiguous scatter store instruction
改善前	9.60E+08	4.80E+08

Busy	L1 busy rate (%)	L2 busy rate (%)	Memory busy rate (%)
改善前	76.29%	2.15%	0.00%

Extra	Gather instruction rate (%)		
	0 flow rate (%)	1 flow rate (%)	2 flow rate (%)
改善前	0.00%	0.00%	100.00%

配列を分割して隣接する2要素のアドレスを128バイト内にすることでGather命令の纏め上げ機能が動作し、浮動小数点ロードL1Dアクセス待ちが改善されました。

改善後ソース (ソースチューニング)

```

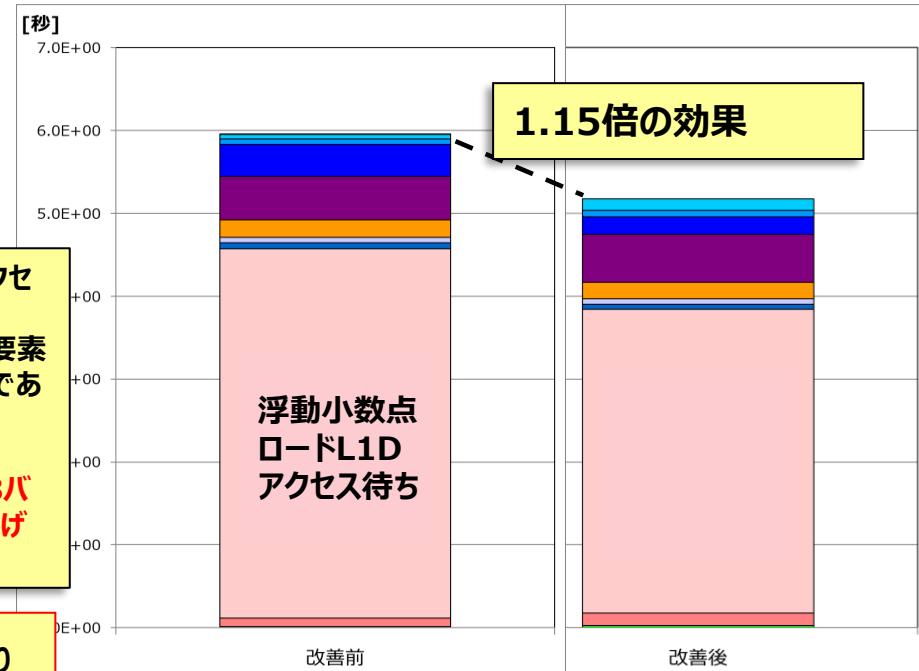
30      real(kind=8),dimension(n,m) :: a1,a2,a3,a4
31      real(kind=8),dimension(n,m) :: b1,b2,b3,b4,c1,c2,c3,c4
32
<<< Loop-information
<<< [OPTIMIZATION]
<<< SIMD(VL: 8)
<<< SOFTWARE PIPELINING
<<< Loop-information End
33 1   v     do i = 1, m
34 1   v       a1( 1,i) = b1( 1,i) +
35 1   v       a1( 2,i) = b1( 2,i) +
36 1   v       a1( 3,i) = b1( 3,i) +
37 1   v       a1( 4,i) = b1( 4,i) +
38 1   v       a2( 1,i) = b2( 1,i) +
39 1   v       a2( 2,i) = b2( 2,i) +
40 1   v       a2( 3,i) = b2( 3,i) +
41 1   v       a2( 4,i) = b2( 4,i) +
42 1   v       a3( 1,i) = b3( 1,i) +
43 1   v       a3( 2,i) = b3( 2,i) +
44 1   v       a3( 3,i) = b3( 3,i) +
45 1   v       a3( 4,i) = b3( 4,i) +
46 1   v       a4( 1,i) = b4( 1,i) +
47 1   v       a4( 2,i) = b4( 2,i) +
48 1   v       a4( 3,i) = b4( 3,i) +
49 1   v       a4( 4,i) = b4( 4,i) + c4( 4,i)
50 1   v   end do

```

n = 4

配列a, 配列b, 配列cは連続にアクセスされるが、それぞれの配列自体 (a1(1,i)など)は、1回転あたり4要素 (32バイト)飛ぶストライドアクセスである
⇒ 隣接する2要素のアドレスが128バイト内となりGather命令の纏め上げ機能が動作する

Gather命令の纏め上げ機能により隣接2要素がL1D\$パイプライン1ポートで処理されるようになった。



Busy

	L1 busy rate (%)	L2 busy rate (%)	Memory busy rate (%)
改善前	76.29%	2.15%	0.00%
改善後	66.21%	2.84%	0.00%

Instruction

	Load-store instruction	
	Load instruction	Store instruction
	Non-contiguous gather load instruction	Non-contiguous scatter store instruction
改善前	9.60E+08	4.80E+08
改善後	9.60E+08	4.80E+08

Extra

	Gather instruction rate (%)		
	0 flow rate (%)	1 flow rate (%)	2 flow rate (%)
改善前	0.00%	0.00%	100.00%
改善後	0.00%	75.00%	25.00%

ストライドアクセスによりGather Load, Scatter Store命令が発生し、浮動小数点ロードL1Dアクセス待ちが多くなっています。

改善前ソース

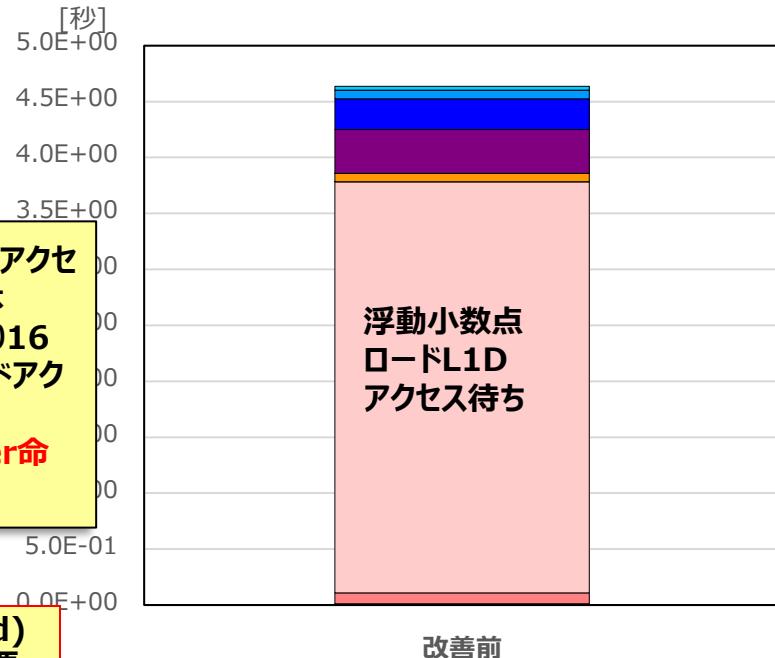
```

50 void sub(int n, int m, double (* restrict a)[n],
51     double (* restrict b)[n], double (* restrict c)[n])
52 {
53     int i;
54
55     <<< Loop-information Start >>>
56     <<< [OPTIMIZATION]
57     <<< SIMD(VL: 8)
58     <<< SOFTWARE PIPELINING(IPC: 2)
59     <<< Loop-information End >>>
60
61     for(i=0; i<m; i++)
62     {
63         a[i][ 0] = b[i][ 0] + c[i][ 0];
64         a[i][ 1] = b[i][ 1] + c[i][ 1];
65         a[i][ 2] = b[i][ 2] + c[i][ 2];
66         a[i][ 3] = b[i][ 3] + c[i][ 3];
67         a[i][ 4] = b[i][ 4] + c[i][ 4];
68         a[i][ 5] = b[i][ 5] + c[i][ 5];
69         a[i][ 6] = b[i][ 6] + c[i][ 6];
70         a[i][ 7] = b[i][ 7] + c[i][ 7];
71         a[i][ 8] = b[i][ 8] + c[i][ 8];
72         a[i][ 9] = b[i][ 9] + c[i][ 9];
73         a[i][10] = b[i][10] + c[i][10];
74         a[i][11] = b[i][11] + c[i][11];
75     }
76     return;
77 }
```

n = 16

配列a, 配列b, 配列cは連続にアクセスされるが、それぞれの配列自体(a[i][1]など)は、1回転あたり16要素(128バイト)飛ぶストライドアクセスである
⇒離散アクセスのため、Gather命令が使用される

非連続ギャザーロード(Gather Load)命令が使用されているが、隣接する2要素のアドレスが128バイト以上離れているためGather命令の纏め上げ機能が動作せずL1ビジー率が高い



Instruction	Load-store instruction	
	Load instruction	Store instruction
Non-contiguous gather load instruction		Non-contiguous scatter store instruction
改善前	7.20E+08	3.60E+08

Busy	L1 busy rate (%)	L2 busy rate (%)	Memory busy rate (%)
改善前	77.75%	2.77%	0.00%
Extra	Gather instruction rate (%)		
改善前	0.00%	0.00%	8.33%

配列を分割して隣接する2要素のアドレスを128バイト内にすることでGather命令の纏め上げ機能が動作し、浮動小数点ロードL1Dアクセス待ちが改善されました。

改善後ソース (ソースチューニング)

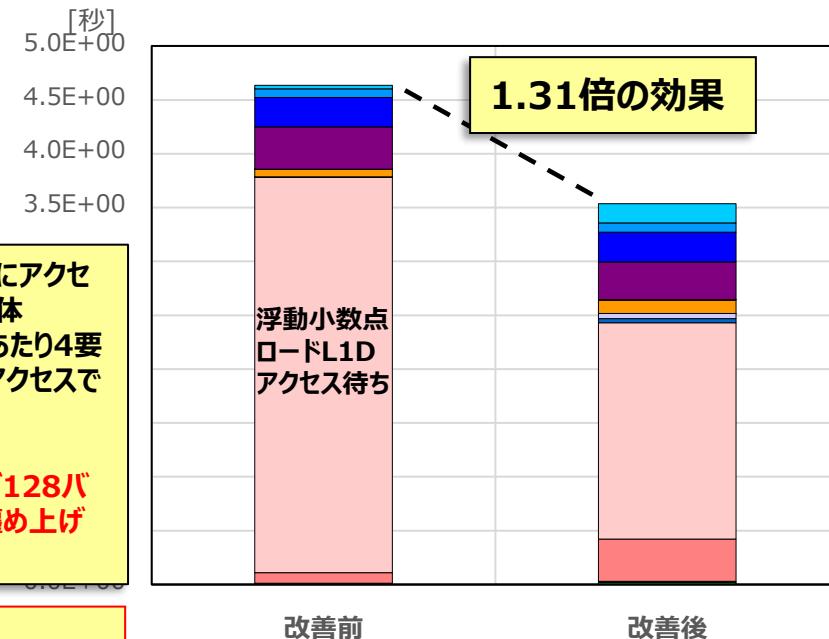
```

73 void sub(int n, int m, double (* restrict a1)[n], double (* restrict a2)[n],
74         double (* restrict a3)[n],
75         double (* restrict b1)[n], double (* restrict b2)[n],
76         double (* restrict b3)[n],
77         double (* restrict c1)[n], double (* restrict c2)[n],
78         double (* restrict c3)[n])
79     {
80         int i;
81
82         <<< Loop-information Start >>>
83         <<< [OPTIMIZATION]
84         <<< SIMD(VL: 8)
85         <<< SOFTWARE PIPELINING(IPC: 2.61,
86         <<< Loop-information End >>>
87         v for(i=0; i<m; i++)
88         v {
89             v a1[i][ 0] = b1[i][ 0] + c1[i][ 0];
90             v a1[i][ 1] = b1[i][ 1] + c1[i][ 1];
91             v a1[i][ 2] = b1[i][ 2] + c1[i][ 2];
92             v a1[i][ 3] = b1[i][ 3] + c1[i][ 3];
93             v a2[i][ 0] = b2[i][ 0] + c2[i][ 0];
94             v a2[i][ 1] = b2[i][ 1] + c2[i][ 1];
95             v a2[i][ 2] = b2[i][ 2] + c2[i][ 2];
96             v a2[i][ 3] = b2[i][ 3] + c2[i][ 3];
97             v a3[i][ 0] = b3[i][ 0] + c3[i][ 0];
98             v a3[i][ 1] = b3[i][ 1] + c3[i][ 1];
99             v a3[i][ 2] = b3[i][ 2] + c3[i][ 2];
100            v a3[i][ 3] = b3[i][ 3] + c3[i][ 3];
101        }
102    return;
103 }
```

n = 4

配列a, 配列b, 配列cは連続にアクセスされるが、それぞれの配列自体 (a1[i][1]など)は、1回転あたり4要素(32バイト)飛びストライドアクセスである
⇒ 隣接する2要素のアドレスが128バイト内となりGather命令の纏め上げ機能が動作する

Gather命令の纏め上げ機能により隣接2要素がL1D\$パイプライン1フローで処理されるようになった。



	Busy	L1 busy rate (%)	L2 busy rate (%)	Memory busy rate (%)
		改善前	改善後	改善前
		77.75%	2.77%	0.00%
		59.49%	1.07%	0.00%

Instruction	Gather instruction rate (%)		
	Extra	0 flow rate (%)	1 flow rate (%)
改善前		0.00%	0.00%
改善後		0.00%	100.00%

Instruction	Load-store instruction	
	Load instruction	Store instruction
Non-contiguous gather load instruction		Non-contiguous scatter store instruction
改善前	7.20E+08	3.60E+08
改善後	7.20E+08	3.60E+08

過剰SFIの回避

- ・ 過剰SFIとは
- ・ 過剰SFIが発生するケース1
- ・ 過剰SFIが発生するケース2
- ・ 過剰SFIの回避（改善前）
- ・ 過剰SFIの回避（ソースチューニング）

- SFI(Store Fetch Interlock)とは

- 先行するstore命令と後続のload命令のアドレスが異なっている場合、loadがstoreを追い越さないようにインターロックをかける制御機構です。
- 基本的にはstoreが行われるアドレスにのみロックがかかります。

- 過剰SFIについて

過剰SFIは上述の制御により過剰にロックかかる現象です。以下のケースで発生します。

- マスク付きSIMDのケースで、マスク判定が0(storeしない)のアドレスがロックされる。
- GatherLoadの纏め上げ機能が動作した場合、GatherLoadでのSFI確認対象がキャッシュラインに含まれる全要素となり、実際にloadを行わないアドレスのSFIを検出しロックされていると判断する場合がある。

過剰SFIが発生するケース1 (1/2)

右の例でX=8のときにはSFIは発生しませんが、
X=6のときには、マスク値が0のアドレスがロックされ
過剰SFIが発生します。

ソースコード

```
real*8 y(X,n), x1(X,n)    <- X=8または6
Do k = 1, iter
  Do j = 1, n
    Do i = 1, X      <- X=8または6
      y(i,j) = y(i,j) + x1(i,j)
    End Do
  End Do
End Do
```

■ X=8 (SFIは発生しない)

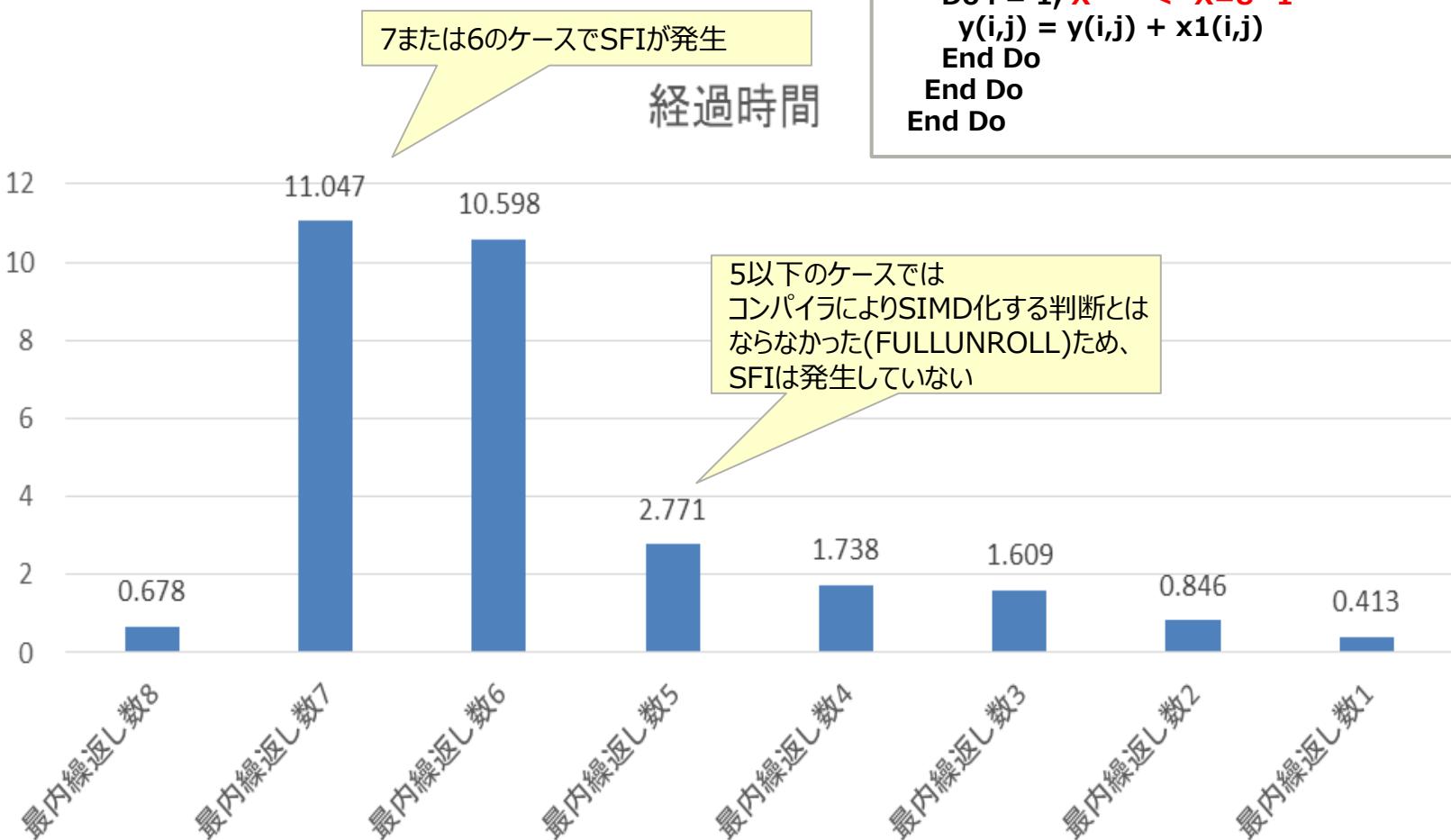
配列y	1, 2, 3, 4, 5, 6, 7, 8, 1, 2, 3, 4, 5, 6, 7, 8, 3, ...
j=1 loadマスク	1 1 1 1 1 1 1 1
j=1 storeマスク	1 1 1 1 1 1 1 1
重なりがないため SFIは発生しない	
j=2 loadマスク	1 1 1 1 1 1 1 1
j=2 storeマスク	1 1 1 1 1 1 1 1

■ X=6 (SFIが発生)

配列y	1, 2, 3, 4, 5, 6, 1, 2, 3, 4, 5, 6, 3, 3, 3, 3, 3, 3, ...
j=1 loadマスク	1 1 1 1 1 1 0 0
j=1 storeマスク	1 1 1 1 1 1 0 0
マスク値が0でもロックされ、 j=2の回転でSFIが発生する	
j=2 loadマスク	1 1 1 1 1 1 0 0
j=2 storeマスク	1 1 1 1 1 1 0 0

過剰SFIが発生するケース1 (2/2)

右の例で、X=7または6のケースではSFIが発生しています。X=5以下のケースではSIMD化されなかつたため、SFIは発生していません。



ソースコード

```
real*8 y(X,n), x1(X,n)  <- X=8~1
Do k = 1, iter
  Do j = 1, n
    Do i = 1, X      <- X=8~1
      y(i,j) = y(i,j) + x1(i,j)
    End Do
  End Do
End Do
```

過剰SFIが発生するケース2

右の例で、X=1のときにはSFIは発生しませんが、X=16のときには、SFI対象のstoreが存在するため過剰SFIが発生します。

X=1 (SFIは発生しない)

配列y

1,1	2,1	3,1	4,1	…	…	15,1	16,1	1,2	2,2	…
i=1~16 load	1,1	1,2	1,3	1,4	…	…	1,15	1,16		
i=1~16 store	1,1	1,2	1,3	1,4	…	…	1,15	1,16		
Gather縫め上げ	x		x	…	…		x		x	…

128バイト飛びのアクセスで、縫め上げなし

i=17~ load

1,17	1,18	…
1,17	1,18	…

i=17~ store

SFI対象のstoreは存在せず、過剰SFIは発生しない

GatherLD縫め上げ機能は動作しない。

X=16 (SFIが発生)

配列y

1,1	1,2	1,3	1,4	…	…	1,15	1,16	1,17	1,18	…
i=1~16 load	1,1	1,2	1,3	1,4	…	…	1,15	1,16		
i=1~16 store	1,1	1,2	1,3	1,4	…	…	1,15	1,16		
Gather縫め上げ	○		○	…	…		○		○	…

8バイト飛びのアクセスで、縫め上げあり

i=17~ load

1,17	1,18	…
1,17	1,18	…

i=17~ store

SFI対象のstoreが存在、過剰SFIが発生する

GatherLD縫め上げ機能の動作により、SFIを確認する対象がキャッシュライン(に含まれる全要素)になる

最内ループがケース1の状態になっており、マスク値が0のアドレスがロックされ過剰SFIが発生しています。ポイント：ループ繰り返し回数が少なく、SWPLされていない場合

改善前ソース	
39	real(4) :: a(20 ,M), b(20 ,M)
40	integer(4) :: M, ITER
41	real(4),parameter :: c=0.5
:	
	<<< Loop-information Start >>>
	<<< [OPTIMIZATION]
	<<< SOFTWARE PIPELINING(IPC: 3.25, ITR: 304, MVE: 7, POL: S)
	<<< PREFETCH(HARD) Expected by compiler :
	<<< a, b
	<<< Loop-information End >>>
46 2 p	DO J=1,M
	<<< Loop-information Start >>>
	<<< [OPTIMIZATION]
	<<< SIMD(VL: 16)
	<<< Loop-information End >>>
47 3 p v	DO I=1,20
48 3 p v	a(I,J) = a(I,J) + c * b(I,J)
49 3 p v	ENDDO
50 2 p	ENDDO



Busy	SFI(Store Fetch Interlock) rate
改善前	0.44

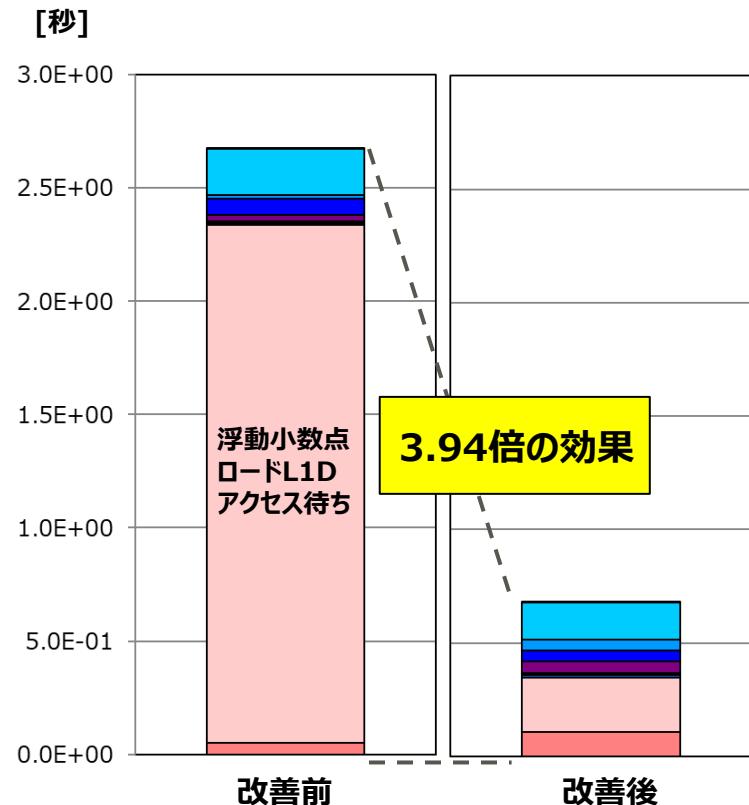
パディングにより配列要素数を変更し、SIMD長の倍数とすることで、過剰SFIを回避することができました。

改善後ソース

```

39      real(4) :: a(32,M), b(32,M)
40      integer(4) :: M, ITER
41      real(4),parameter :: c=0.5
:
<<< Loop-information Start >>>
<<< [OPTIMIZATION]
<<< SOFTWARE PIPELINING(IPC: 3.25, ITR: 304,
                           MVE: 7, POL: S)
<<< PREFETCH(HARD) Expected by compiler :
<<<     a, b
<<< Loop-information End >>>
46   2 p      DO J=1,M
<<< Loop-information Start >>>
<<< [OPTIMIZATION]
<<< SIMD(VL: 16)
<<< Loop-information End >>>
47   3 p      v          DO I=1,20
48   3 p      v          a(I,J) = a(I,J) + c * b(I,J)
49   3 p      v          ENDDO
50   2 p      ENDDO

```



Busy

	SFI(Store Fetch Interlock) rate
改善前	0.43
改善後	0.01

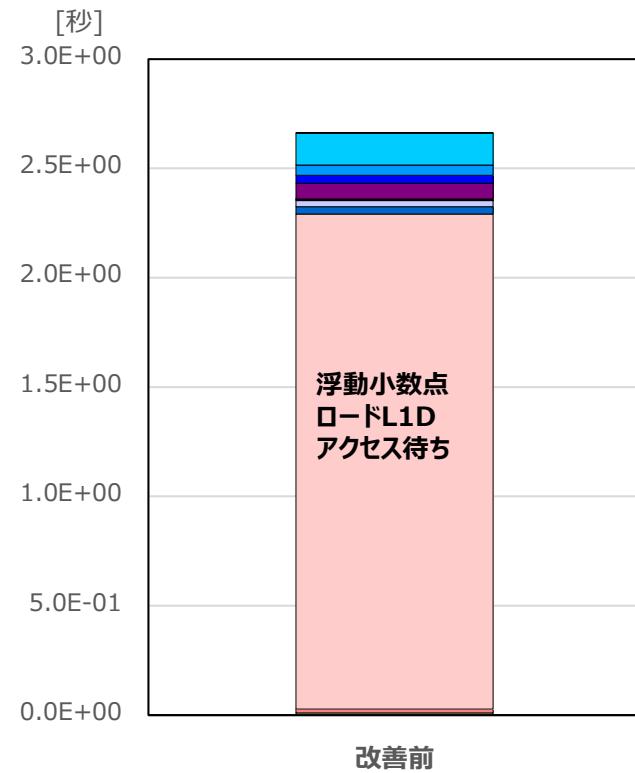
最内ループがケース1の状態になっており、マスク値が0のアドレスがロックされ過剰SFIが発生しています。ポイント：ループ繰り返し回数が少なく、SWPLされていない場合

改善前ソース

```

42 void sfil1(float (* restrict a)[20], float (* restrict b)[20],
        int m, int iter)
43     {
44         float c=0.5;
45         int i,j,k;
46
47         #pragma omp parallel
48         {
49             <<< Loop-information Start >>>
50             :
51             <<< Loop-information End >>>
52             for(k=0; k<iter; k++)
53             {
54                 #pragma omp for nowait
55                 <<< Loop-information Start >>>
56                 :
57                 <<< Loop-information End >>>
58                 p           for(j=0; j<m; j++)
59                 p           {
60                     v           for(i=0; i< 20; i++)
61                     v           {
62                         v               a[j][i] = a[j][i] + c * b[j][i];
63                     }
64                 }
65             }
66             return;
67         }

```



Busy	SFI(Store Fetch Interlock) rate
改善前	0.44%

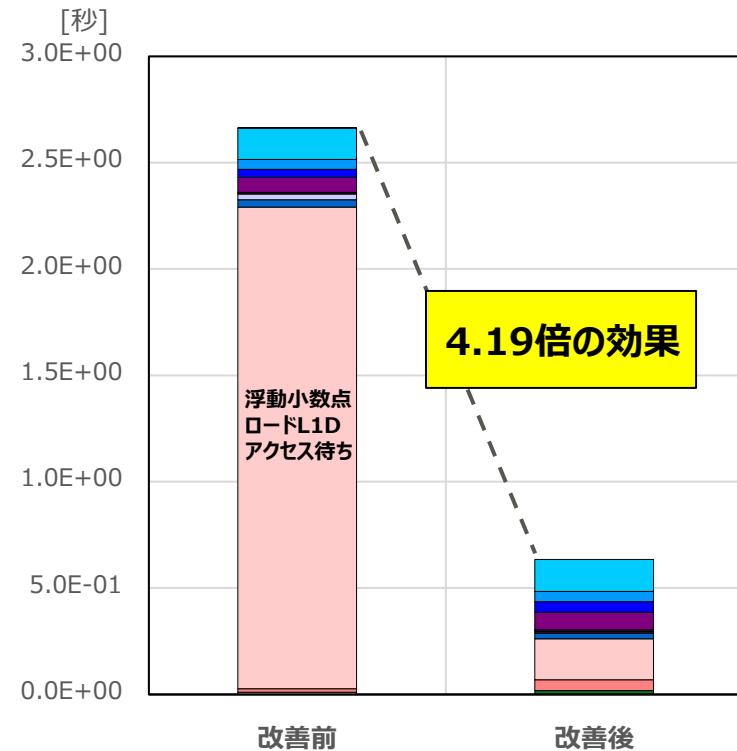
パディングにより配列要素数を変更し、SIMD長の倍数とすることで、過剰SFIを回避することができました。

改善後ソース

```

42 void sfil1(float (* restrict a)[32], float (* restrict b)[32],
             int m, int iter)
43     {
44         float c=0.5;
45         int i,j,k;
46
47         #pragma omp parallel
48         {
49             <<< Loop-information Start >>>
50             :
51             <<< Loop-information End >>>
52             for(k=0; k<iter; k++)
53             {
54                 #pragma omp for nowait
55                 <<< Loop-information Start >>>
56                 :
57                 <<< Loop-information End >>>
58                 p           for(j=0; j<m; j++)
59                 p           {
60                     <<< Loop-information Start >>>
61                     :
62                     <<< Loop-information End >>>
63                     v           for(i=0; i< 20; i++)
64                     v           {
65                         v           a[j][i] = a[j][i] + c * b[j][i];
66                     v           }
67                     }
68                     return;
69                 }
70             }
71         }
72     }

```



Busy	SFI(Store Fetch Interlock) rate
改善前	0.44%
改善後	0.01%

Multiple Structures命令の活用

- Multiple Structures命令の適用条件
- Multiple Structures命令の活用（改善前）
- Multiple Structures命令の活用（ソースチューニング）

Multiple Structures命令の適用条件

- Fortran、C/C++(Tradモード/Clangモード)でサポート
- Fotran、C/C++(Tradモード)はオプションによりON/OFFが可能です。

```
-Ksimd_use_multiple_structures  
-Ksimd_nouse_multiple_structures
```

デフォルトは-Ksimd_use_multiple_structures、
-Ksimd_nouse_multiple_structuresで抑止が可能

- 構造体配列(AoS)の場合、最内次元が2,3,4個の要素で形成されている、かつ、その要素すべてにアクセスしている場合、適用されます。最内次元が5個以上の要素の場合、適用されません。

- 適用可能例(Fortran) :

```
real*8 y(n),x(4,n)  
  
Do j = 1, iter  
  Do i = 1, n  
    y(i) = x(1,i) + x(2,i) + x(3,i) + x(4,i)  
  End Do  
End Do
```

- 適用可能例(C/C++) :

```
double y[n], x[N][4];  
  
for (j = 0; j < iter; j++) {  
  for (i = 0; i < N; i++) {  
    y(i) = x[i][0] + x[i][1] + x[i][2] + x[i][3];  
  }  
}
```

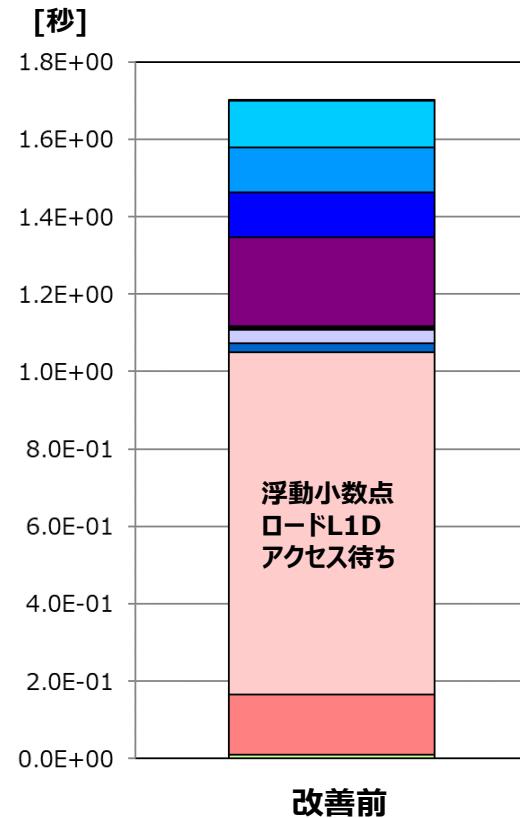
- なお、構造体配列(AoS)を配列構造体(SoA)に書き換えることで、連続ロードとなり性能向上が見込める場合は、書き換えをお勧めします。

構造体配列の配列が6個の要素でGatherロード命令が使用（Multiple Structures命令が適用されない）されているため、浮動小数点ロードL1Dアクセス待ちが多くなっています。

改善前ソース

```
<<< Loop-information Start >>>
<<< [OPTIMIZATION]
<<< SIMD(VL: 8)
<<< SOFTWARE PIPELINING(IPC: 2.35, ITR: 96,
                           MVE: 3, POL: S)
<<< PREFETCH(HARD) Expected by compiler :
<<< a, b
<<< Loop-information End >>>
22  2 p 2v      do i=1,N
23  2 p 2v      b(i)=a(1,i) + a(2,i) + a(3,i) + &
                  a(4,i) + a(5,i) + a(6,i)
24  2 p 2v      enddo
```

配列aが6個の要素であるため、
Multiple Structures命令が適用され
ない。

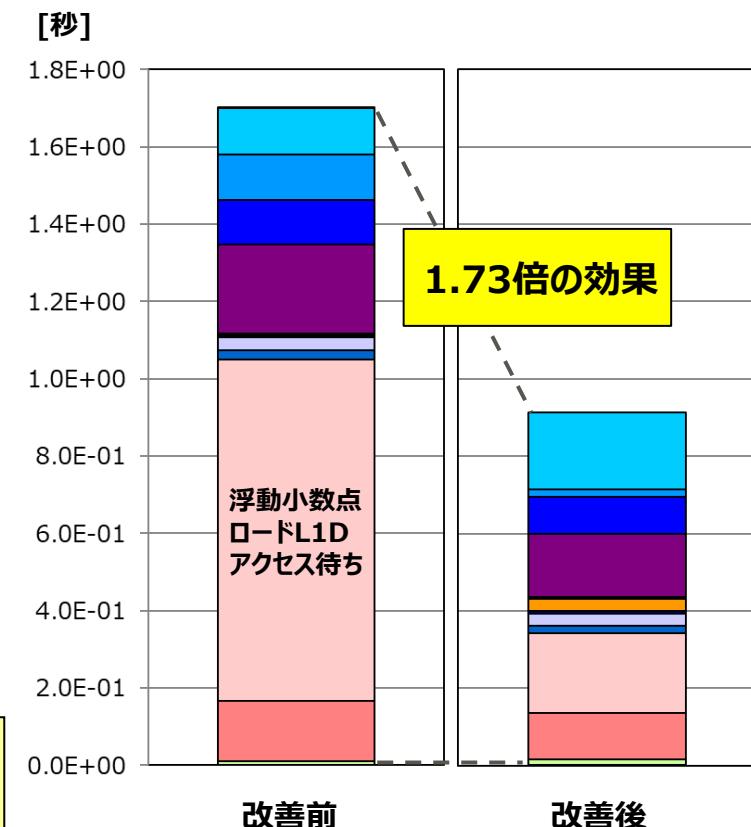


配列を3つずつの要素に分けることでMultiple Structures命令が適用され、浮動小数点ロードL1Dアクセス待ちが削減されました。

改善後ソース（ソースチューニング）

```
<<< Loop-information Start >>>
<<< [OPTIMIZATION]
<<< SIMD(VL: 8)
<<< SOFTWARE PIPELINING(IPC: 1.33, ITR: 56,
                           MVE: 2, POL: S)
<<< PREFETCH(HARD) Expected by compiler :
<<<   c, a, b
<<< Loop-information End >>>
22  2 p v      do i=1,N
23  2 p v          b(i) = a(1,i) + a(2,i) + a(3,i) + &
                           c(1,i) + c(2,i) + c(3,i)
24  2 p v      enddo
```

Multiple Structures命令適用
配列aが6個の要素であるため、
配列cを追加して3個ずつの要素にした。

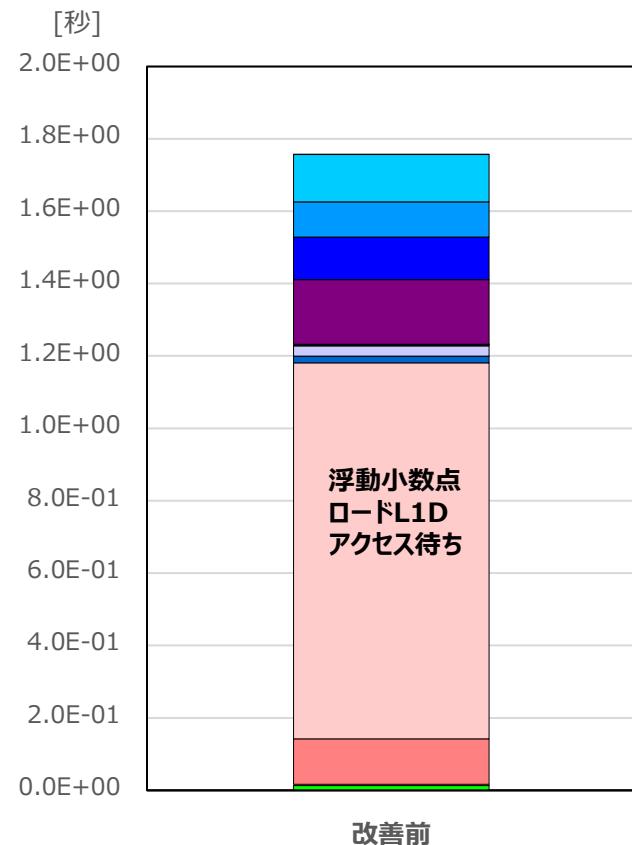


構造体配列の配列が6個の要素でGatherロード命令が使用（Multiple Structures命令が適用されない）されているため、浮動小数点ロードL1Dアクセス待ちが多くなっています。

改善前ソース

```
33     void MultiStructure(int n, int m, int iter)
34     {
35         int i,k;
36
37         #pragma omp parallel
38         {
39             <<< Loop-information Start >>>
39             :
40             <<< Loop-information End >>>
41             for(k=0; k< iter; k++)
42             {
43                 #pragma omp for nowait
44                 <<< Loop-information Start >>>
45                 :
46                 <<< Loop-information End >>>
47                 p 2v      for(i=0; i< n; i++)
48                 p 2v      {
49                 p 2v          b[i]=a[i][0] + a[i][1] + a[i][2]
50                         + a[i][3] + a[i][4] + a[i][5];
51                 p 2v      }
52                 }
53             }
54         return;
55     }
```

配列aが6個の要素であるため、
Multiple Structures命令が適用され
ない。



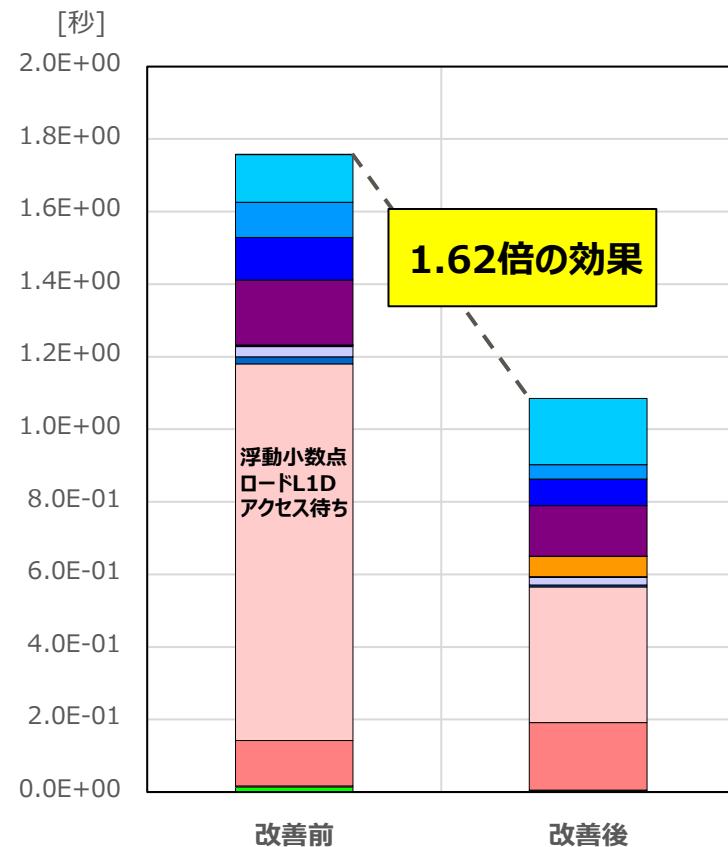
配列を3つずつの要素に分けることでMultiple Structures命令が適用され、浮動小数点ロードL1Dアクセス待ちが削減されました。

改善後ソース（ソースチューニング）

```

32     void MultiStructure(int n, int m, int iter)
33     {
34         int i,k;
35
36         #pragma omp parallel
37         {
38             <<< Loop-information Start >>>
39             :
40             <<< Loop-information End >>>
41             for(k=0; k<iter; k++)
42             {
43                 #pragma omp for nowait
44                 <<< Loop-information Start >>>
45                 :
46                 <<< Loop-information End >>>
47                 p   v   for(i=0; i<N; i++)
48                 p   v   {
49                 p   v       b[i] = a[i][0] + a[i][1] + a[i][2]
49                                         + c[i][0] + c[i][1] + c[i][2];
50                 p   v   }
51             }
52         }
53     return;
54 }
```

Multiple Structures命令適用
配列aが6個の要素であるため、
配列cを追加して3個ずつの要素にした。



ハードウェアプリフェッчの距離調整

- ・ プリフェッчの距離について
- ・ ハードウェアプリフェッчの距離調整機能
- ・ ハードウェアプリフェッчの距離調整機能(onL2)
- ・ ハードウェアプリフェッчの距離調整機能(onメモリ)

■ ハードウェアプリフェッヂ・ソフトウェアプリフェッヂの距離について

- ハードウェアプリフェッヂ・ソフトウェアプリフェッヂでは、以下のライン先のデータに対しプリフェッヂを行います。

ハードウェアプリフェッヂ		ソフトウェアプリフェッヂ	
L1プリフェッヂ	L2プリフェッヂ	L1プリフェッヂ	L2プリフェッヂ
最大6ライン	最大40ライン	自動	自動

ハードウェアプリフェッヂの
距離はユーザの設定も
可能になる

ソフトウェアプリフェッヂは
距離の自動調整がある

■ プリフェッヂの距離はアプリケーションのアクセス依存です。

プリフェッヂするキャッシュラインを遠くにするとスラッシングする場合があるため、
近くのキャッシュラインをプリフェッヂすることをお勧めします。

■ ハードウェアプリフェッч距離設定用コマンド(hwpfctl)

項目	内容
書式	<pre>hwpfctl [--disableL1] [--disableL2] [--distL1 lines_l1] [--distL2 lines_l2] [--weakL1] [--weakL2] [--verbose] command {arguments ...} hwpfctl --default [--verbose] command {arguments ...} hwpfctl --reset [--verbose] hwpfctl --help</pre>
説明	<p>hwpfctlコマンドは、A64FXに搭載されているハードウェアプリフェッч(stream detect mode)の振る舞いを変更する。変更対象となるCPUコアは、プロセスアフィニティに準ずる。</p>
オプション	<ul style="list-style-type: none"> --disableL1 --disableL2 <p>L1/L2キャッシュに対するハードウェアプリフェッчを無効とする。省略時はハードウェアプリフェッчが有効である。</p> --distL1=lines_l1 --distL2=lines_l2 <p>L1/L2キャッシュに対してプリフェッчするキャッシュラインを、キャッシュミスが発生したキャッシュラインを起点とするキャッシュラインの数で指定する。lines_l1にはL1キャッシュに対するプリフェッч対象ラインを1から15までの値を指定できる。また、lines_l2にはL2キャッシュに対するプリフェッч対象ラインを1から60までの値を指定できる。ただし、lines_l2に指定した値は4の倍数に切り上げられてシステムレジスタに書き込まれる。0を指定した場合はCPUのデフォルト値で動作する。省略時ならびに無効な値が指定された場合は0が指定されたものとみなす。</p> --weakL1 --weakL2 <p>L1/L2キャッシュに対するプリフェッч要求の優先度をweakとすることを指示する。省略時はstrongである。</p> --default <p>デフォルト設定でcommandを起動する。--verbose以外のオプションは無視する。</p> --reset <p>システムレジスタの値を初期化する。--verbose以外のオプションは無視する。</p> --verbose <p>システムレジスタの変更前後の値を出力する。</p> --help <p>使用方法を表示する。</p>

■ ハードウェアプリフェッч距離設定用コマンド(hwpfctl)の使用例

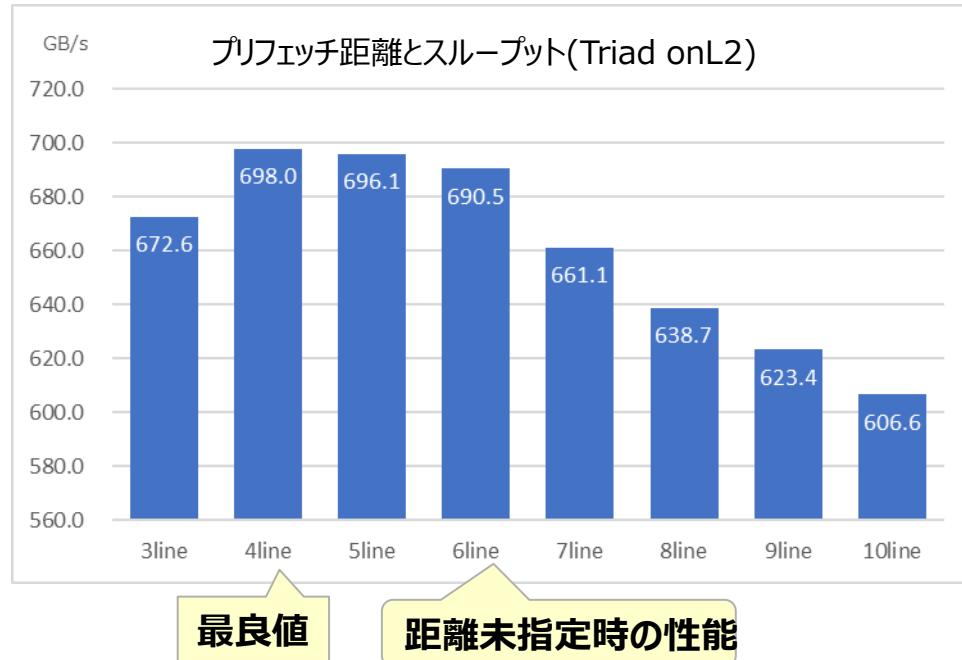
```
hwpfctl --distL1=6 --distL2=40 a.out
```

ハードウェアプリフェッчの距離調整の実行結果(onL2)

FUJITSU

以下にハードウェアプリフェッчの距離調整の実行結果を示します。

■ Triad(L2キャッシュアクセスケース)によるL1プリフェッч距離評価



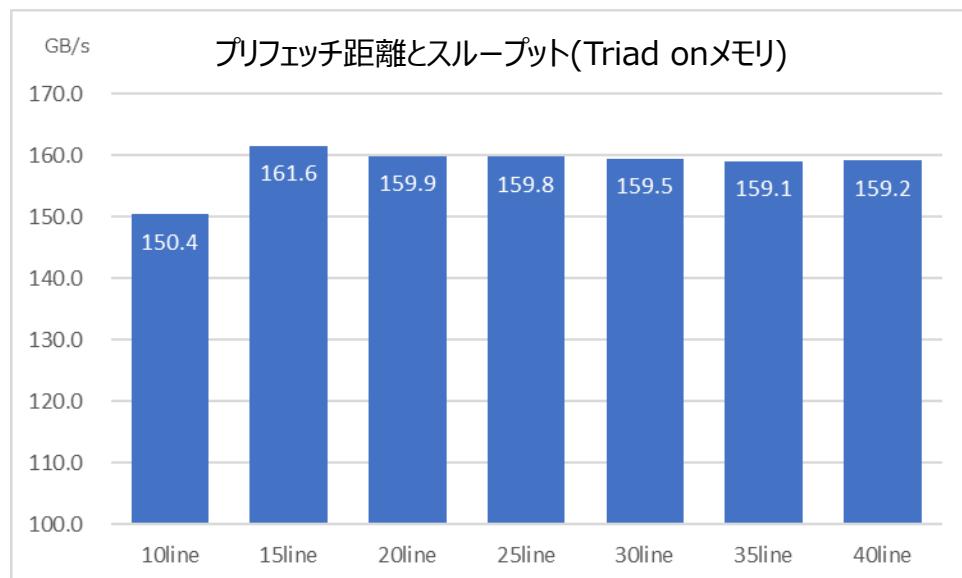
Triad

```
!$omp parallel
Do j = 1, iter
!$omp do
Do i = 1, n
y(i)=x1(i) + c0 * x2(i)
End Do
!$omp end do nowait
End Do
!$omp end parallel
```

■ 設定コマンド

```
hwpfctl -distL1=3~10 a.out
```

■ Triad(メモリアクセスケース)によるL2プリフェッч距離評価



最良値

距離未指定時の性能

Triad

```
!$omp parallel
Do j = 1, iter
 !$omp do
 Do i = 1, n
 y(i)=x1(i) + c0 * x2(i)
 End Do
 !$omp end do nowait
 End Do
 !$omp end parallel
```

※スループットは、データ書き込み先キャッシュラインの読み込み分を含めない値になります

■ 設定コマンド

```
hwpfctl -distL2=10~40 a.out
```

SVEのベクトルレジスタサイズ(SIMD幅)

- SVEのベクトルレジスタサイズ(SIMD幅)について
- -Ksimd_reg_size=agnostic指定時の最適化影響(注意点)

- プロセッサでサポートしているSIMD幅

ARMのSVE拡張はVector Length(SIMD幅)が128bitから2048bitの範囲の128bit単位で実装者が自由に決めることが可能である。

A64FXではVector Lengthを512bitで実装を行っている。

- A64FXでサポートするVector Lengthは512bit, 256bit, 128bit

- コンパイラオプション

- -Ksimd_reg_size={ 128 | 256 | 512 | agnostic }

デフォルトは、-Ksimd_reg_size=512です。

- simd_reg_size={ 128 | 256 | 512 }

SVEのベクトルレジスタサイズを指定します。単位はビットです。翻訳時に、本オプションで指定した値がSVEのベクトルレジスタサイズであるとみなして、最適化を実施します。ただし、生成した実行可能プログラムは、本オプションで指定したサイズのSVEのベクトルレジスタを実装しているCPUアーキテクチャでのみ正常に動作します。

指定したベクトルレジスタのサイズが、実行するCPUアーキテクチャのベクトルレジスタのサイズより小さいことが明らかな場合、prctl(2)システムコールなどをを利用して有効なベクトルレジスタのサイズを設定する必要があります。

- simd_reg_size=agnostic

SVEのベクトルレジスタを特定のサイズとみなさず翻訳を行い、実行時にSVEのベクトルレジスタサイズを決定する実行可能プログラムを作成します。この実行可能プログラムは、CPUアーキテクチャに実装されたSVEのベクトルレジスタサイズによらず実行可能です。

ただし、-Ksimd_reg_size={128|256|512}オプションを指定した場合と比べて、実行性能が低下する場合があります。

-Ksimd_reg_size=agnostic指定時に-Ksimd_reg_size=512(デフォルト)とは同様の最適化が行われないことがあります。その場合、実行性能低下する可能性があります。

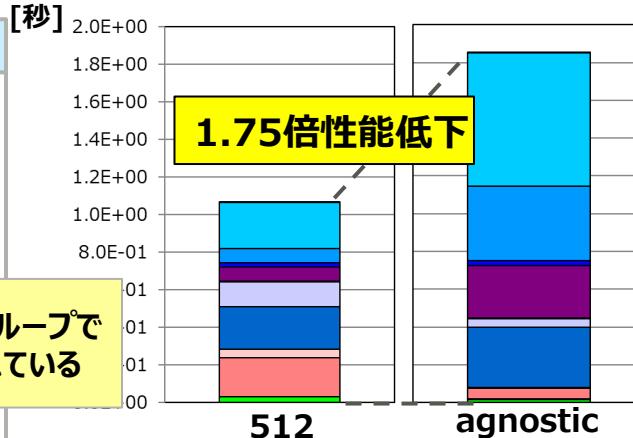
-Ksimd_reg_size=512(デフォルト)ソース

```

24 2 p      do k=1,L
25 3 p      do ii=1,N,blk
26 4 p 8    blk=8
             <<< Loop-information Start >>>
             <<< [OPTIMIZATION]
             <<< SOFTWARE PIPELINING(IPC: 0.31, ITR: 192,
                           MVE: 2, POL: L)
             <<< PREFETCH(HARD)
             <<< b, a
             <<< Loop-information End >>>
             do j=1,M
             <<< Loop-information Start >>>
             <<< [OPTIMIZATION]
             <<< SIMD(VL: 8)
             <<< FULL UNROLLING
             <<< Loop-information End >>>
27 5 p fv   do i=ii,ii+blk-1
28 5 p fv   a(i,j,k)=a(i,j,k)+c*b(i,j,k)
29 5 p fv   enddo
30 4 p 8    enddo
31 3 p      enddo
32 2 p      enddo

```

最内ループがSIMD長であり、外側のループで
ソフトウェアパイプラインングが適用されている



-Ksimd_reg_size=agnostic指定時ソース

```

24 2 p      do k=1,L
25 3 p      do ii=1,N,blk
26 4 p      do j=1,M
             <<< Loop-information Start >>>
             <<< [OPTIMIZATION]
             <<< SIMD(VL: AGNOSTIC; VL: 2 in 128-bit)
             <<< PREFETCH(HARD) Expected by compiler :
             <<< b, a
             <<< Loop-information End >>>
27 5 p 2v    do i=ii,ii+blk-1
28 5 p 2v    a(i,j,k)=a(i,j,k)+c*b(i,j,k)
29 5 p 2v    enddo
30 4 p      enddo
31 3 p      enddo
32 2 p      enddo

```

ソフトウェアパイプラインングが適用されない

■ 注意事項

- Ksimd_reg_size=agnostic指定で翻訳した場合にはSIMD幅に依存した最適化
(上記の例ではソフトウェアパイプラインング) が行われません。

-Ksimd_reg_size=agnostic指定時に-Ksimd_reg_size=512(デフォルト)とは同様の最適化が行われないことがあります。その場合、実行性能低下する可能性があります。

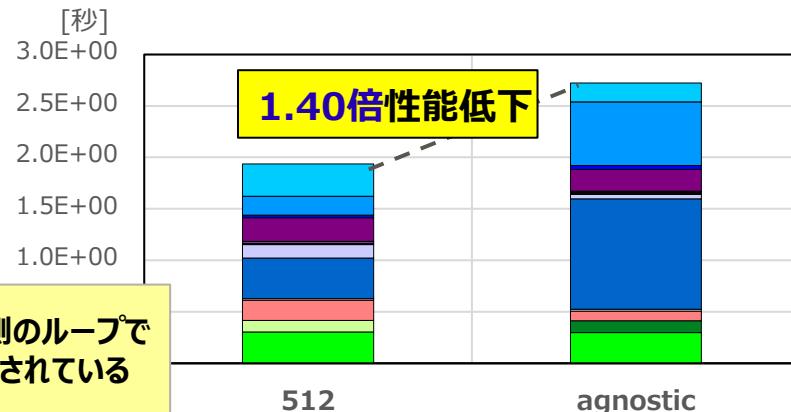
-Ksimd_reg_size=512(デフォルト)ソース

```

46     for(iter=1; iter<=itmax; iter++) {
48         #pragma omp for
49         p         for(k=0;k<L; k++) {
50             for(ii=0;ii<N;ii+=blk) {
51                 <<< Loop-information Start >>> blk=8
52                 <<< [OPTIMIZATION]
53                 : 
54                 <<< Loop-information End
55                 for(j=0;j<M;j++)
56                 <<< Loop-information Start
57                 for(k=0;k<L; k++) {
58                     p         2v         for(i=ii; i<ii+blk-1; i++) {
59                     p         2v         a[k][j][i]=a[k][j][i]+c*b[k][j][i];
60                     p         2v         }
61                     p         }
62                     }

```

最内ループがSIMD長であり、外側のループでソフトウェアパイプラインングが適用されている



-Ksimd_reg_size=agnostic指定時ソース

```

46     for(iter=1; iter<=itmax; iter++) {
48         #pragma omp for
49         p         for(k=0;k<L; k++) {
50             for(ii=0;ii<N;ii+=blk) {
51                 for(j=0;j<M;j++) {
52                     <<< Loop-information Start >>>
53                     <<< [OPTIMIZATION]
54                     <<< SIMD(VL: AGNOSTIC; VL: 2 in 128-bit)
55                     <<< PREFETCH(HARD) Expected by compiler :
56                     (unknown)
57                     <<< Loop-information End >>>
58                     p         2v         for(i=ii; i<ii+blk-1; i++) {
59                     p         2v         a[k][j][i]=a[k][j][i]+c*b[k][j][i];
60                     p         2v         }
61                     p         }
62                     }

```

ソフトウェアパイプラインングが適用されない

■ 注意事項

- Ksimd_reg_size=agnostic指定で翻訳した場合にはSIMD幅に依存した最適化（上記の例ではソフトウェアパイプラインング）が行われません。

半精度実数型の活用

- ・ 半精度実数型の活用（改善前）
- ・ 半精度実数型の活用（ソースチューニング）

倍精度実数型のデータの場合、SIMD長は8ですが、データの精度を落とすことで SIMD長を長くしバンド幅・演算器を有効活用できます。

改善前ソース

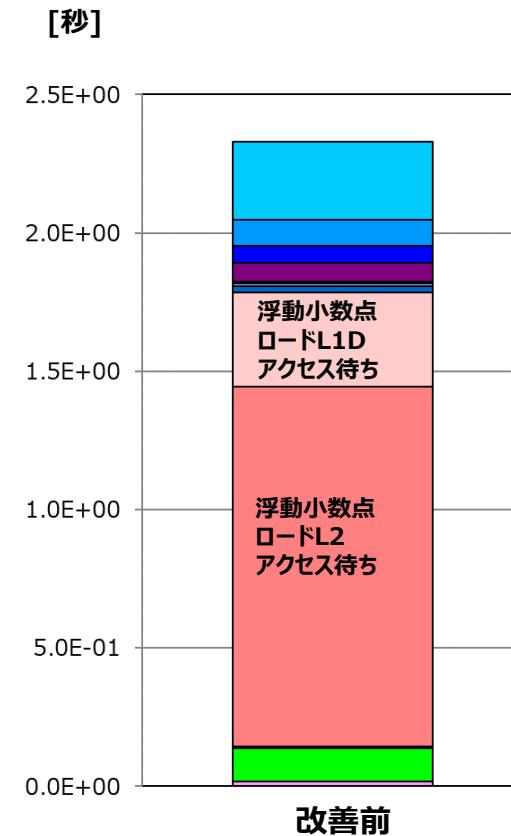
```

2      integer,parameter::N=60000
:
19
20      real(8)::x1(N),x2(N),y(N)
      real(8),parameter::c=0.5
:
      <<< Loop-information Start >>>
      <<< [OPTIMIZATION]
      <<< SIMD(VL: 8)
      <<< SOFTWARE PIPELINING(IPC: 3.25, ITR: 144,
                                MVE: 4, POL: S)
      <<< PREFETCH(HARD) Expected by compiler :
      <<<     x2, x1, y
      <<< Loop-information End >>>
24  2 p 2v      do i=1,N
25  2 p 2v      y(i) = x1(i) + c * x2(i)
26  2 p 2v      enddo

```

SIMD幅(Vector Length)=512[bit]の場合のSIMD長

データ型	SIMD長
倍精度型	8
単精度型	16
半精度型	32
1バイト型	64



	GFLOPS	Floating-point operation peak ratio (%)
改善前	51.52	6.71%

桁数を落とした半精度実数型を使うことでSIMD長を32に拡張することができます。

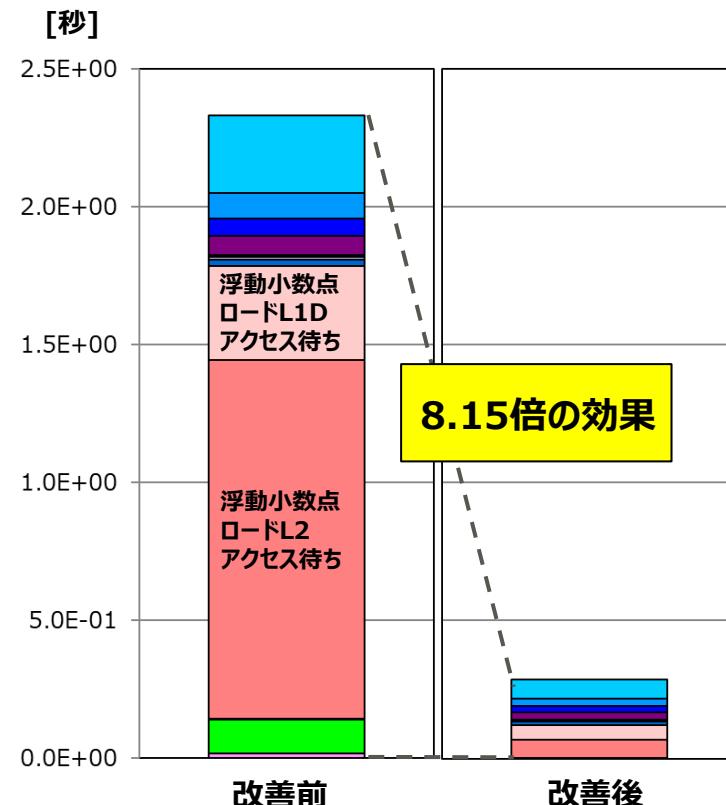
データ容量が小さくなることで浮動小数点ロードアクセス待ちが改善されました。

改善後ソース

```

2      integer,parameter::N=60000
:
19      real(2)::x1(N),x2(N),y(N)
20      real(2),parameter::c=0.5
:
<<< Loop-information Start >>>
<<< [OPTIMIZATION]
<<< SIMD(VL: 32)
<<< SOFTWARE PIPELINING(IPC: 3.25, ITR: 576,
                           MVE: 4, POL: S)
<<< PREFETCH(HARD) Expected by compiler :
<<<   x2, x1, y
<<< Loop-information End >>>
24  2 p 2v      do i=1,N
25  2 p 2v      y(i) = x1(i) + c * x2(i)
26  2 p 2v      enddo

```



	GFLOPS	Floating-point operation peak ratio (%)
改善前	51.52	6.71%
改善後	421.57	54.89%

倍精度実数型のデータの場合、SIMD長は8ですが、データの精度を落とすことでSIMD長を長くしバンド幅・演算器を有効活用できます。

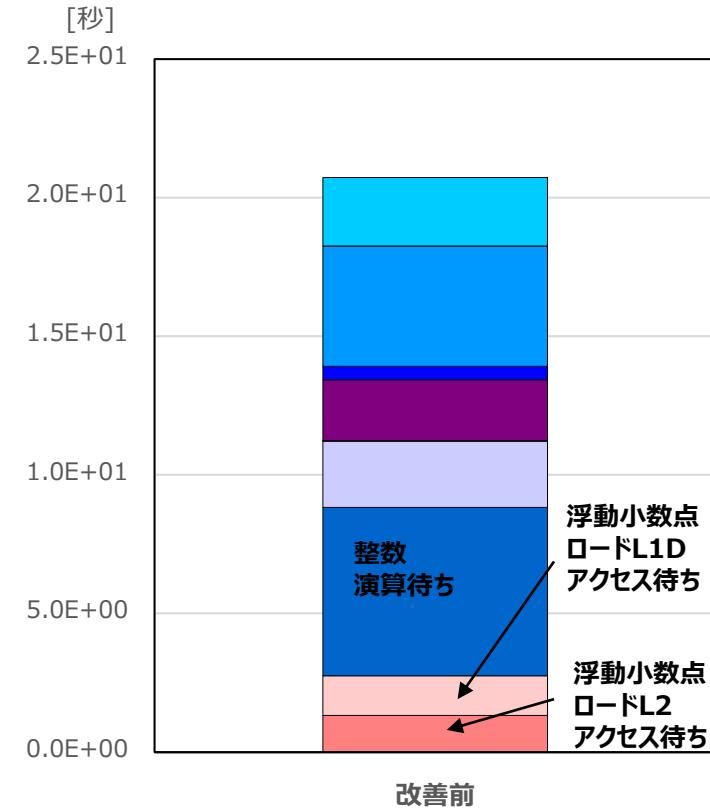
改善前ソース

```

34     double const c=0.5;
35     int i,k;
36
37     for(k=0;k<itmax; k++)
38     {
39         <<< Loop-information Start >>>
40         <<< [OPTIMIZATION]
41         <<< SIMD(VL: AGNOSTIC;
42             VL: 2 in 128-bit Interleave: 1)
43         <<< Loop-information End >>>
44         for(i=0;i<N; i++)
45         {
46             y[i] = x1[i] + c * x2[i];
47         }
48     }

```

N = 60000
itmax = 1000000
配列宣言
double x1[N], x2[N], y[N];



SIMD幅(Vector Length)=512[bit]の場合のSIMD長

データ型	SIMD長
倍精度型	8
単精度型	16
半精度型	32
1バイト型	64

Statistics	GFLOPS	Floating-point operation
改善前	5.79	1.20E+11

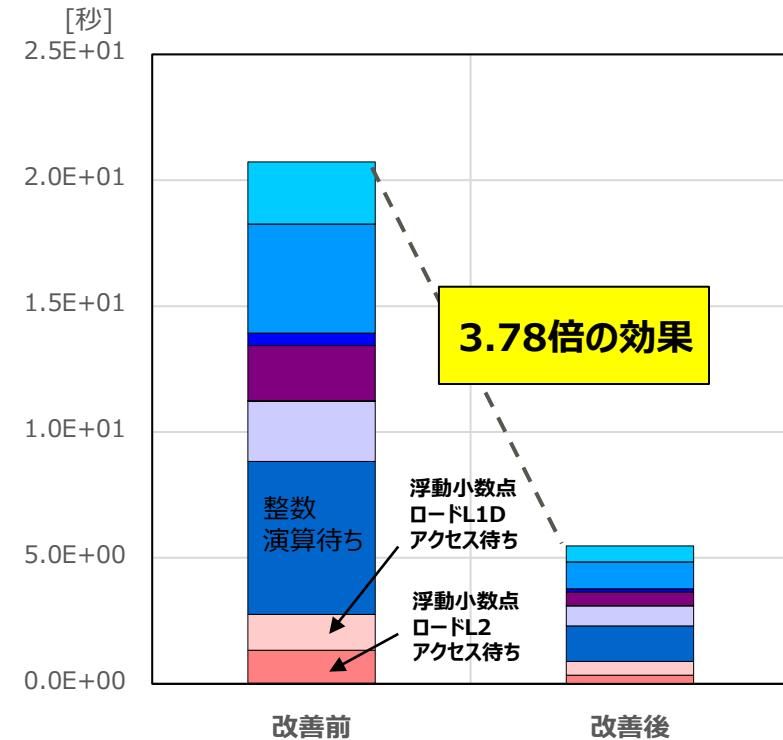
桁数を落とした半精度実数型を使うことでSIMD長を32に拡張することができます。
 データ容量が小さくなることで浮動小数点ロードアクセス待ちが改善されました。

改善後ソース

```

34      _Float16 c = 0.5f16;
35      int i,k;
36
37  4    for(k=0;k<itmax; k++)
38  {
39      <<< Loop-information Start >>>
40      <<< [OPTIMIZATION]
41      <<< SIMD(VL: AGNOSTIC;
42          VL: 8 in 128-bit Interleave: 1)
43      <<< Loop-information End >>>
44      v    for(i=0;i<N; i++)
45      {
46          y[i] = x1[i] + c * x2[i];
47      }
48  }
```

N = 60000
itmax = 1000000
配列宣言
double x1[N], x2[N], y[N];



■ 注意事項

- C/C++のtradモードでは半精度実数型機能がないため、本手法によるチューニングはできません。

Statistics	GFLOPS	Floating-point operation
改善前	5.79	1.20E+11
改善後	21.91	1.20E+11

スレッド並列チューニング

- ・ スレッド並列化率の改善
- ・ スレッド並列化率の向上
- ・ スレッド並列実行効率の改善
- ・ ラージページの設定による実行効率の改善
- ・ メモリ使用量の削減

スレッド並列化率の改善

- ・ スレッド並列化率とは
- ・ スレッド並列化率の向上

スレッド並列化率とは

並列化率とは 1 並列実行時の並列実行可能部分の占める割合を意味します。

1 並列実行時

逐次実行
部分

並列実行可能部分

2 並列実行時

逐次実行
部分

並列実行可能部分

2 並列実行時は、並列実行可能部分が 1 並列実行の 1/2 となる。

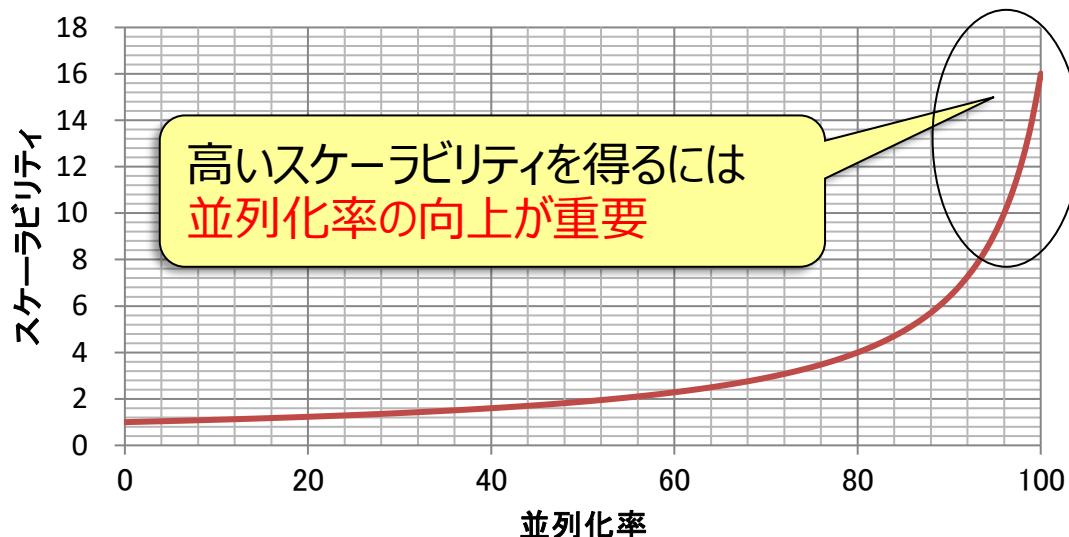
n 並列実行時のスレッド並列化率とスケーラビリティの関係を表す法則としてアムダールの法則があります。

■ アムダールの法則

$$\text{スケーラビリティ} = \frac{1}{(1-p) + \frac{p}{n}}$$

- p : 並列化率
- n : 並列数

$n = 16$ (理想的な 1 6 並列処理) の場合



スレッド並列化率の向上

- 定義引用関係が明らかでないループ
- ポインタ変数が含まれるループ
- データ依存関係があるループ

● !OCL NORECURRENCE

以下のプログラムで、配列aの添字式が別の配列要素y(j) のため、本処理系は配列aをループスライスしても問題ないか判断できません。配列a をループスライスしても問題ないとプログラマが分かっている場合に**NORECURRENCE指示子**を指定することで並列化することができます。

改善前ソース	改善後ソース
<pre> <<< Loop-information Start >>> <<< [OPTIMIZATION] <<< SOFTWARE PIPELINING(IPC: 0.32, ITR: 6, MVE: 2, POL: S) <<< PREFETCH(HARD) Expected by compiler : <<< b, y <<< Loop-information End >>> 6 1 s 2s do i=1,160000 7 1 m 2m a(y(i))=a(y(i))+b(i) 8 1 p 2v end do : jwd5228p-i "a.f90", line 7: データの定義引用の順序が 逐次実行と異なるため、このDOループは並列化できません。 jwd6228s-i "a.f90", line 7: データの定義引用の順序が 逐次実行と異なる可能性があるため、このDOループはSIMD化でき ません。 </pre>	<pre> 5 !ocl norecurrence(a) <<< Loop-information Start >>> <<< [PARALLELIZATION] <<< Standard iteration count: 762 <<< [OPTIMIZATION] <<< SIMD(VL: 16) <<< SOFTWARE PIPELINING(IPC: 2.33, ITR: 416, MVE: 7, POL: S) <<< PREFETCH(HARD) Expected by compiler : <<< y, b <<< Loop-information End >>> 6 1 pp 2v do i=1,160000 7 1 p 2v a(y(i))=a(y(i))+b(i) 8 1 p 2v end do </pre>

！注意！

- NORECURRENCE指示子をループスライス不可能な配列に指定した場合、本処理系は誤ったループスライスを行うことがあります。
- 配列名を省略した場合、対象範囲内のすべての配列に有効となります。

● !OCL NOALIAS

ポインタ変数が記憶領域のどこを占めるかは実行時に決るために、データの依存関係が不明となり並列化されません。ポインタ変数が同一の記憶領域を指すことはないとプログラマが分かっている場合に**NOALIAS**指示子を指定することで並列化することができます。

改善前ソース	改善後ソース
<pre> 1 real,dimension(100000),target::x 2 real,dimension(:),pointer::a,b 3 a=>x(1:10000) 4 b=>x(10001:20000) 5 6 7 <<< Loop-information Start >>> 8 <<< [OPTIMIZATION] 9 <<< PREFETCH(HARD) Expected by compiler : 10 <<< x 11 <<< Loop-information End >>> 12 1 s 4s do i=1,100000 13 1 s 4s b(i) = a(i)+1.0 14 1 s 4s end do jwd5228p-i "a.f90", line 8: データの定義引用の順序が 逐次実行と異なるため、このDOループは並列化できません。 jwd6228s-i "a.f90", line 8: データの定義引用の順序が 逐次実行と異なる可能性があるため、このDOループはSIMD化でき ません。 </pre>	<pre> 1 real,dimension(100000),target::x 2 real,dimension(:),pointer::a,b 3 a=>x(1:10000) 4 b=>x(10001:20000) 5 6 !ocl noalias 7 <<< Loop-information Start >>> 8 <<< [PARALLELIZATION] 9 <<< Standard iteration count: 1143 10 <<< [OPTIMIZATION] 11 <<< SIMD(VL: 16) 12 <<< SOFTWARE PIPELINING(IPC: 2.75, 13 ITR: 288, MVE: 4, POL: S) 14 <<< PREFETCH(HARD) Expected by compiler : 15 <<< (unknown) 16 <<< Loop-information End >>> 17 1 pp 2v do i=1,100000 18 1 p 2v b(i) = a(i)+1.0 19 1 p 2v end do </pre>

● ピーリングを行い並列化

以下のループ[°]は配列aについて、 $i=1$ の時と $i=n$ の時に依存性を持っているので並列化できません。ループの最初や最後をループ[°]外に出して依存をなくすことで並列化が促進されます。

改善前ソース	改善後ソース
<pre> <<< Loop-information Start >>> <<< [OPTIMIZATION] <<< PREFETCH(HARD) Expected by compiler : <<< b, a <<< Loop-information End >>> 4 1 s 2s do i=1,n 5 1 s 2m a(i)=a(1)+b(i)+a(n) 6 1 s 2v end do jwd5202p-i "a.f90", line 5: データの定義引用の順序が 逐次実行と異なるため、このDOループは並列化できません。(名前:a) jwd5208p-i "a.f90", line 5: 定義引用の順序が分からな いため、定義引用順序が逐次実行と異なる可能性があり、 このDOループは並列化できません。(名前:a) </pre>	<pre> 4 a(1)=a(1)+b(1)+a(n) <<< Loop-information Start >>> <<< [PARALLELIZATION] <<< Standard iteration count: 843 <<< [OPTIMIZATION] <<< SIMD(VL: 16) <<< SOFTWARE PIPELINING(IPC: 3.00, ITR: 384, MVE: 4, POL: S) <<< PREFETCH(HARD) Expected by compiler : <<< a, b <<< Loop-information End >>> 5 1 pp 2v do i=2,n-1 6 1 p 2v a(i)=a(1)+b(i)+a(n) 7 1 p 2v enddo 8 a(n)=a(1)+b(n)+a(n) </pre>

スレッド並列実行効率の改善

- False Sharingの改善
- 処理量が不規則なループ
- 適切な並列化次元での並列化

False Sharingの改善

- False Sharingとは
- False Sharing（改善前）
- False Sharing（ソースチューニング）

False Sharingとは

FUJITSU

False Sharingとは、スレッド間でキャッシュラインのInvalidateおよびCopy Backが頻発する現象のことです。

4スレッド並列と仮定した場合の例

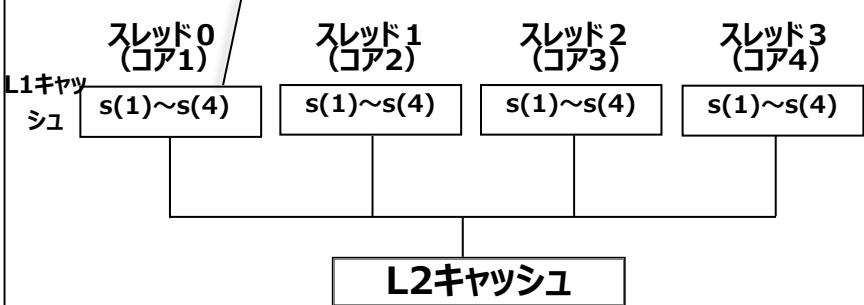
改善前ソース

```
1 subroutine sub(s,a,b,ni,nj)
2   real*8 a(ni,nj),b(ni,nj)
3   real*8 s(nj)          nj=4
4
5   1 pp
6   1 p
7   2 p 8v
8   2 p 8v           ni=2000
9   2 p 8v
10  1 p
11
12 end
```

初期状態

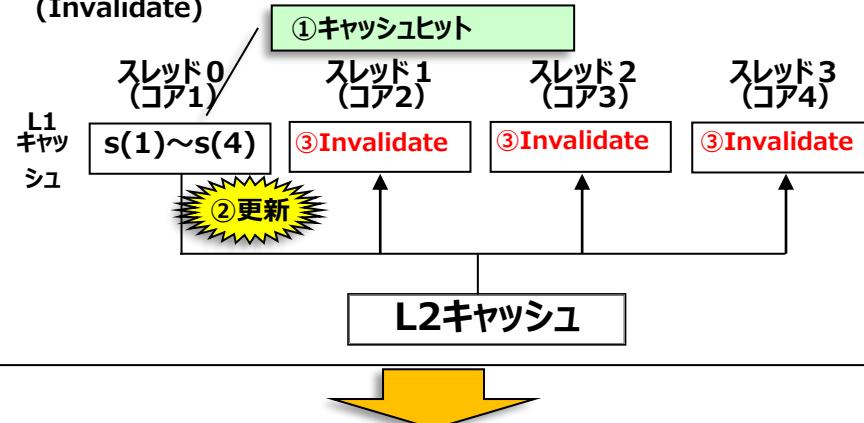
キャッシュにはキャッシュライン単位でデータが載る

各スレッドは、s(1)～s(4)を含む同一キャッシュラインを読み込む



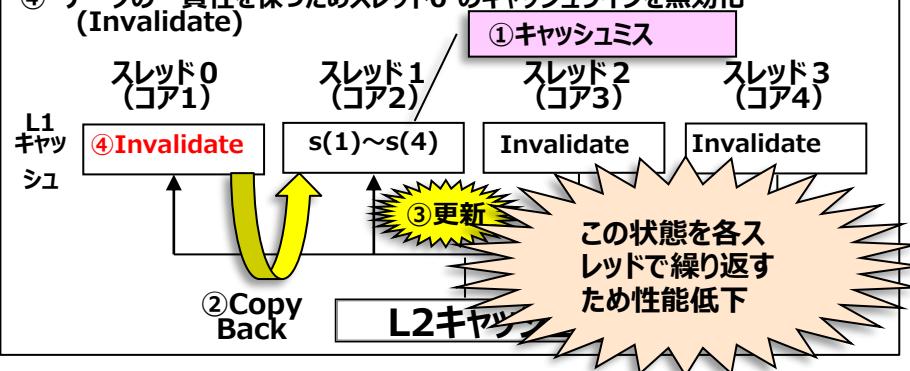
スレッド0が s(1) の更新を指示

- ① キャッシュヒット
- ② スレッド0 が s(1) の更新を完了
- ③ データの一貫性を保つためスレッド 1～3 のキャッシュラインを無効化 (Invalidate)



スレッド1が s(2) の更新を指示

- ① キャッシュミス
- ② スレッド0 からスレッド1へキャッシュラインを Copy Back
- ③ スレッド1 が s(2) の更新を完了
- ④ データの一貫性を保つためスレッド0 のキャッシュラインを無効化 (Invalidate)



Fortran False Sharing (改善前)

FUJITSU

並列化次元のjの繰返し数が16回転と少なく、配列aのデータがスレッド間でキャッシュラインを共有してしまうため、FalseSharingが発生します。
そのため、データアクセス待ちが多くなっています。

改善前ソース

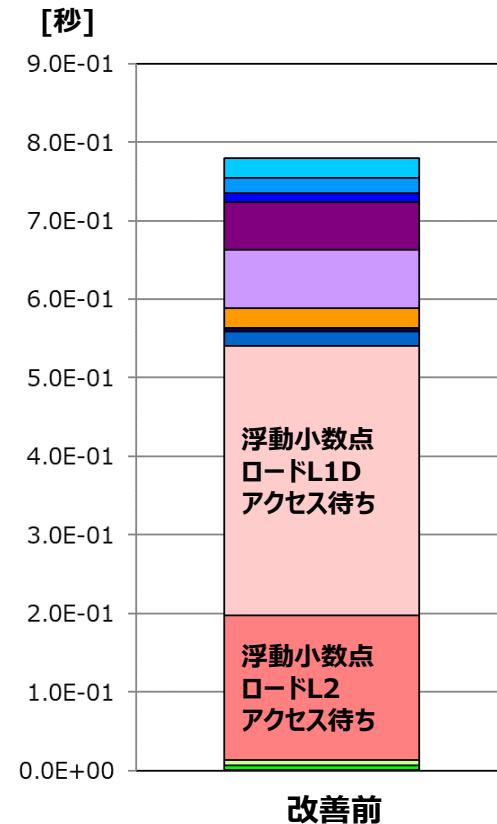
```

22 subroutine sub(flag)
23 integer*8 i,j,n
24 parameter(n=60000)
25 parameter(m=16)
26 real*8 a(m,n),b(m,n)
27 integer flag(m,n)
28 common /com/a,b
29 :
32 1 p      do i=1,m
33 2 p 2v    <<< Loop-information Start >>>
34 3 p 2v    <<< [OPTIMIZATION]
35 3 p 2v    <<< SIMD(VL: 8)
36 3 p 2v    <<< SOFTWARE PIPELINING(IPC: 2.28, ITR: 128,
37 2 p 2v    MVE: 3, POL: S)
38 1 p      <<< Loop-information End >>>
33 2 p 2v    do j=1,n
34 3 p 2v    if(flag(i,j).eq.1)then
35 3 p 2v      a(i,j)=b(j,i)
36 3 p 2v    endif
37 2 p 2v    enddo
38 1 p      enddo

```

**並列化次元の繰
返し数が16回転**

False Sharing発生



Cache

	L1I miss rate (/Effective instruction)	Load-store instruction	L1D miss	L1D miss rate (/Load-store instruction)	L1D miss demand rate (%) (/L1D miss)	L1D miss hardware prefetch rate (%) (/L1D miss)	L1D miss software prefetch rate (%) (/L1D miss)	L2 miss	L2 miss rate (/Load-store instruction)	L2 miss demand rate (%) (/L2 miss)	L2 miss hardware prefetch rate (%) (/L2 miss)	L2 miss software prefetch rate (%) (/L2 miss)
改善前	0.00	3.57E+09	7.27E+08	0.20	9.39%	90.85%	0.00%	1.75E+04	0.00	21.82%	100.00%	0.00%

Fortran False Sharing (ソースチューニング)

FUJITSU

ループ交換を行い、外側で並列化することでFalse Sharingを回避できます。その結果、L1キャッシュミス数が削減され、データアクセス待ちが改善されました。

改善後ソース

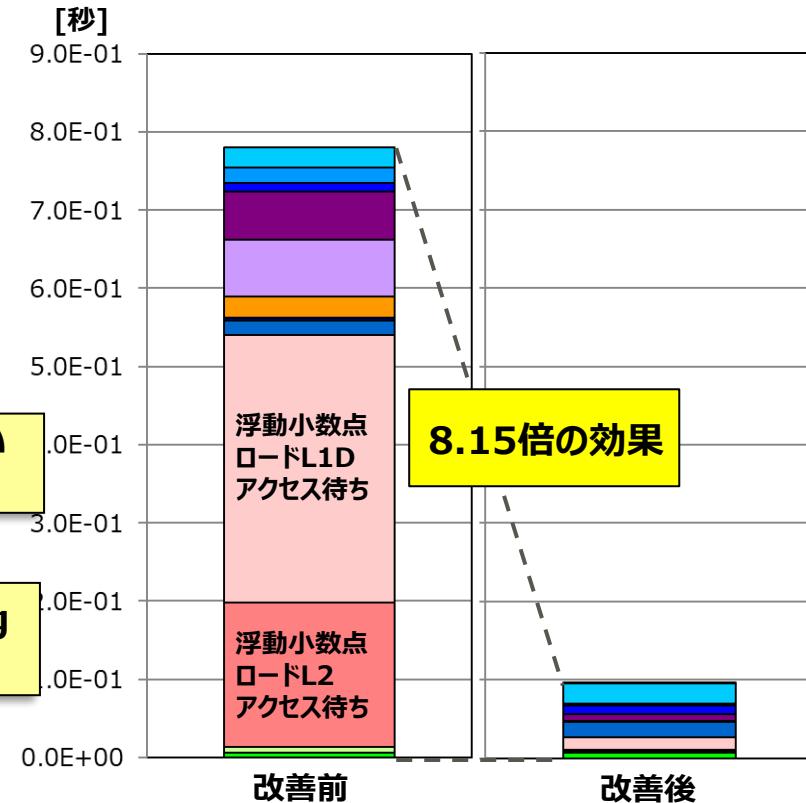
```

23      integer*8 i,j,n
24      parameter(n=60000)
25      parameter(m=16)
26      real*8 a(m,n),b(m,n)
27      integer flag(m,n)
:
<<< Loop-information Start >>>
<<< [OPTIMIZATION]
<<< SOFTWARE PIPELINING(IPC: 1.32, ITR: 48,
                           MVE: 2, POL: S)
<<< PREFETCH(SOFT) : 4
<<< SEQUENTIAL : 4
<<< b: 4
<<< Loop-information End >>>
32   1 p 2      do j=1,n
<<< Loop-information Start >>>
<<< [OPTIMIZATION]
<<< SIMD(VL: 8)
<<< FULL UNROLLING
<<< Loop-information End >>>
33   2 p fv      do i=1,m
34   3 p fv          if(flag(i,j),eq.1)then
35   3 p fv              a(i,j)=b(j,i)
36   3 p fv          endif
37   2 p fv      enddo
38   1 p 2      enddo

```

ループ交換を行い
外側で並列化

False Sharing
回避



False Sharing 回避によりL1Dミス数が
削減され性能向上した

	L1I miss rate (/Effective instruction)	Load-store instruction	L1D miss	L1D miss rate (/Load-store instruction)	L1D miss demand rate (%) (/L1D miss)	L1D miss hardware prefetch rate (%) (/L1D miss)	L1D miss software prefetch rate (%) (/L1D miss)	L2 miss	L2 miss rate (/Load-store instruction)	L2 miss demand rate (%) (/L2 miss)	hardware prefetch rate (%) (/L2 miss)	software prefetch rate (%) (/L2 miss)
改善前	0.00	3.57E+09	7.27E+08	0.20	9.39%	90.85%	0.00%	1.75E+04	0.00	21.82%	100.00%	0.00%
改善後	0.00	1.19E+09	5.29E+07	0.04	4.40%	84.91%	10.70%	1.61E+04	0.00	7.89%	85.83%	6.28%

並列化次元のjの繰返し数が16回転と少なく、配列aのデータがスレッド間でキャッシュラインを共有してしまうため、FalseSharingが発生します。
そのため、データアクセス待ちが多くなっています。

```

    改善前ソース

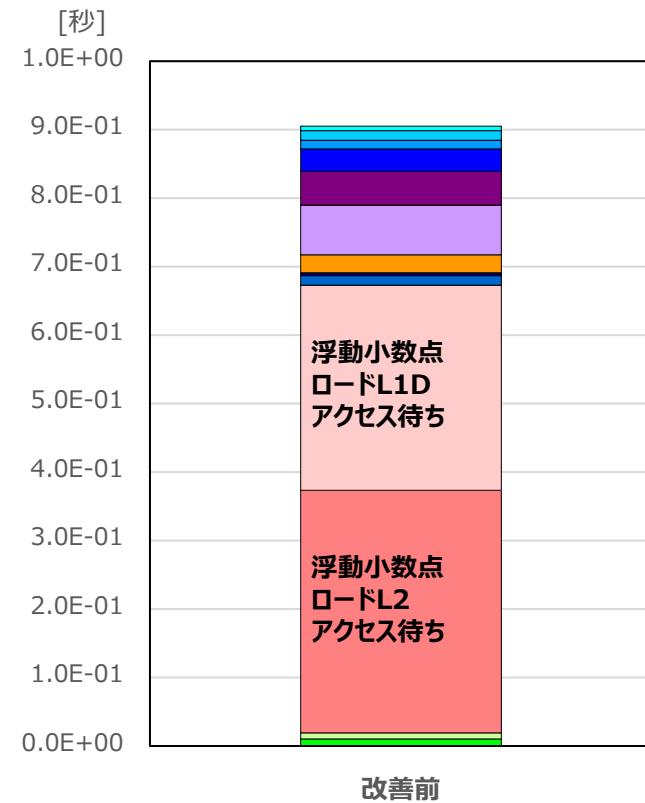
42      #pragma omp for
        <<< Loop-information Start >>>
        <<< [OPTIMIZATION]
        <<< PREFETCH(HARD) Expect :
        <<< b, a, (unknown)
        <<< Loop-information End >>>
        for (i=0;i<M;i++)
            {
43    p        <<< Loop-information Start >>>
44    p        <<< [OPTIMIZATION]
45    p        <<< SIMD(VL: 8)
46    p        <<< SOFTWARE PIPELINING(IPC: 2.00, ITR: 128,
47    p        MVE: 3, POL: S)
48    p        <<< PREFETCH(HARD) Expected by compiler :
49    p        <<< b, a, (unknown)
50    p        <<< Loop-information End >>>
51    p    2v    for(j=0;j<N;j++)
52    p    2v    {
53    p    2v    if(flag[j][i]==1)
54    p    2v    {
55    p    2v    a[j][i]=b[i][j];
56    p    2v    }
57    p    2v    }

```

M=16
N=60000
配列宣言
double a[N][M],
b[N][M];

並列化次元の繰返
し数が16回転

False Sharing発生



Cache	L1I miss rate (/Effective instruction)	Load-store instruction	L1D miss	L1D miss rate (/Load-store instruction)	L1D miss demand rate (%) (/L1D miss)	L1D miss hardware prefetch rate (%) (/L1D miss)	L1D miss software prefetch rate (%) (/L1D miss)	L2 miss	L2 miss rate (/Load-store instruction)	L2 miss demand rate (%) (/L2 miss)	L2 miss hardware prefetch rate (%) (/L2 miss)	L2 miss software prefetch rate (%) (/L2 miss)
改善前	0.00	4.04E+09	7.28E+08	0.18	30.45%	69.61%	0.00%	3.94E+04	0.00	48.46%	58.66%	0.00%

ループ交換を行い、外側で並列化することで False Sharing を回避できます。その結果、L1 キャッシュミス数が削減され、データアクセス待ちが改善されました。

改善後ソース

```

42      #pragma omp for
<<< Loop-information Start >>
<<< [OPTIMIZATION]
<<< SOFTWARE PIPELINING(
    MVE: 3, POL: S)
<<< PREFETCH(HARD) Expect
<<< a, (unknown)
<<< Loop-information End >>>
43 p   for(j=0;j<N;j++)
44 p   {
<<< Loop-information Start >>
<<< [OPTIMIZATION]
<<< SIMD(VL: 8)
<<< FULL UNROLLING
<<< Loop-information End >>>
45 p   fv   for (i=0;i<M;i++)
46 p   fv   {
47 p   fv   if(flag[j][i]==1)
48 p   fv   {
49 p   fv   a[j][i]=b[i][j];
50 p   fv   }
51 p   fv   }
52 p   }

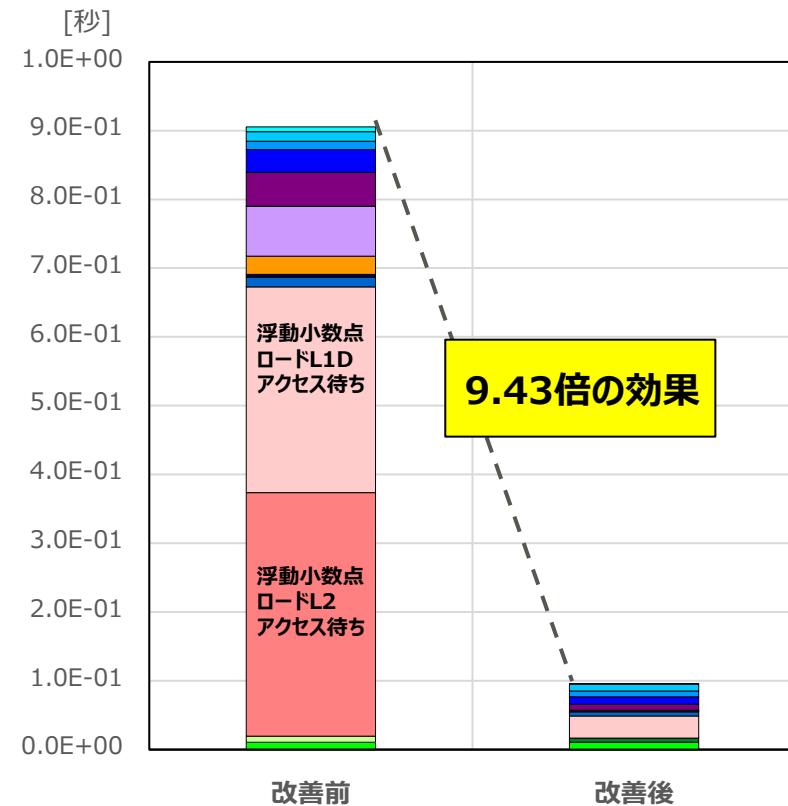
```

M=16
N=60000

配列宣言
double a[N][M],
b[N][M];

ループ交換を行い
外側で並列化

False Sharing
回避



False Sharing 回避により L1D ミス数が
削減され性能向上した

Cache	L1 miss rate (/Effective instruction)	Load-store instruction	L1D miss	L1D miss rate (/Load-store instruction)	L1D miss demand rate (%)	(/L1D miss)	L1D miss hardware prefetch rate (%)	(/L1D miss)	L2 miss	L2 miss rate (/Load-store instruction)	L2 miss demand rate (%)	(/L2 miss)	hardware prefetch rate (%)	(/L2 miss)	software prefetch rate (%)	(/L2 miss)
改善前	0.00	4.04E+09	7.28E+08	0.18	30.45%		69.61%		0.00%	3.94E+04	0.00	48.46%	58.66%	0.00%		
改善後	0.00	1.27E+09	4.81E+07	0.04	6.15%		93.86%		0.00%	1.93E+04	0.00	15.65%	89.19%	0.00%		

処理量が不規則なループ

- ・ 処理量が不規則なループ（改善前）
- ・ 処理量が不規則なループ（OpenMP チューニング）

処理量が不規則な場合、スケジュール方法が static のサイクリック分割ではロードインバランスが発生します。下の例ではロードインバランスのためバリア同期待ちが多くなっています。

改善前ソース

```

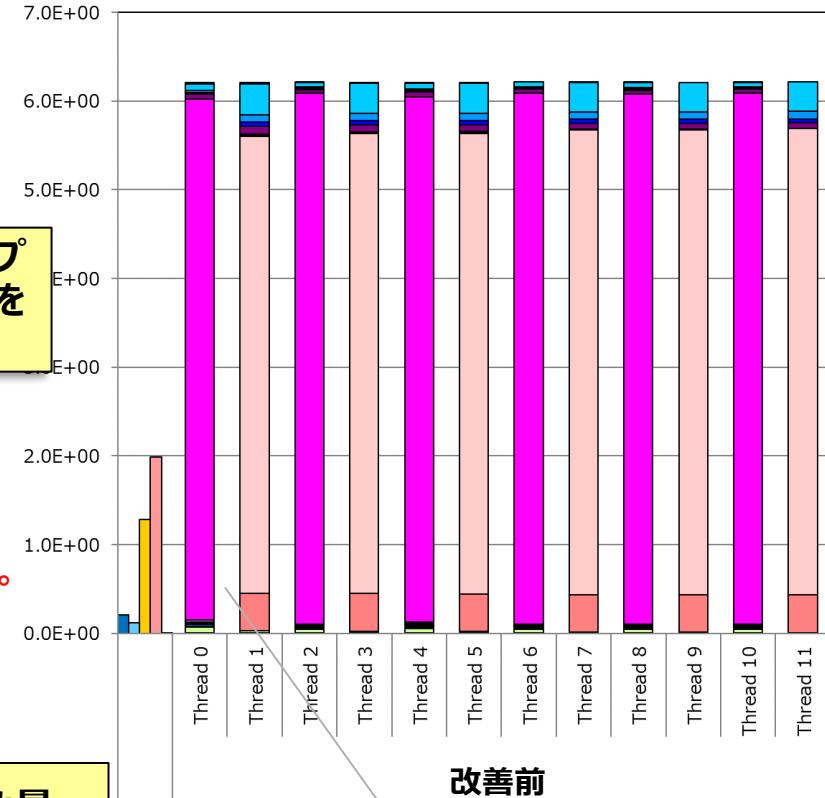
1 subroutine init(a,b,ie,n)
:
8   1     <<< Loop-information Start >>>
9   2     <<< [OPTIMIZATION]
10  2     <<< FUSED
11  2     <<< Loop-information End >>>
12  1       do i=1,n
13  1         if (mod(i,2).eq.0) then
14  2           ie(i)=100000
15  1         endif
16  1       enddo
:
17
18
19
20  1 p       subroutine sub(a,b,s,ie,n)
19  1 p         real a(n),b(n),s
20  1 p         integer ie(n)
21  1 p         !$omp parallel do schedule(static,1)
22  2 p           do j=1,n
23  2 p             <<< Loop-information Start >>>
24  2 p             <<< [OPTIMIZATION]
25  2 p             <<< SIMD(VL: 16)
26  2 p             <<< SOFTWARE PIPELINING(IPC: 3.50,
27  2 p               ITR: 448, MVE: 7, POL: S)
28  2 p             <<< Loop-information End >>>
29  1 p               do i=1,ie(j)
30  2 p                 a(i) = a(i)*b(i)*s
31  2 p               enddo
32  1 p             enddo
33  1 p         end subroutine sub
:
34
35       program main
36         parameter(n=1000000)
37         call init(a,b,ie,n)
38
39         call sub(a,b,2.0,ie,n)
40
41       end program main

```

偶数回転時のみ評価ループ
の終値である配列 ie に値を
設定している

評価ループ

制御変数 j が偶数時のみ最
内ループは100,000回転する。



改善前

ロードインバランスによる
バリア同期待ちが発生

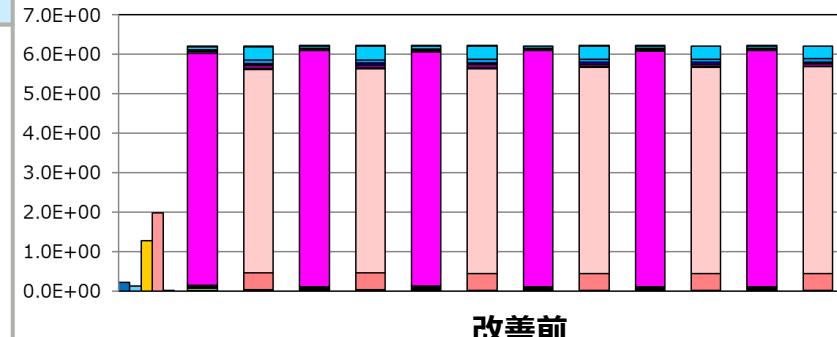
スケジュール方法をdynamicにすることで、先に処理が終了したスレッドが次の処理を行えるようになり、ロードインバランスが改善しました。

改善後ソース

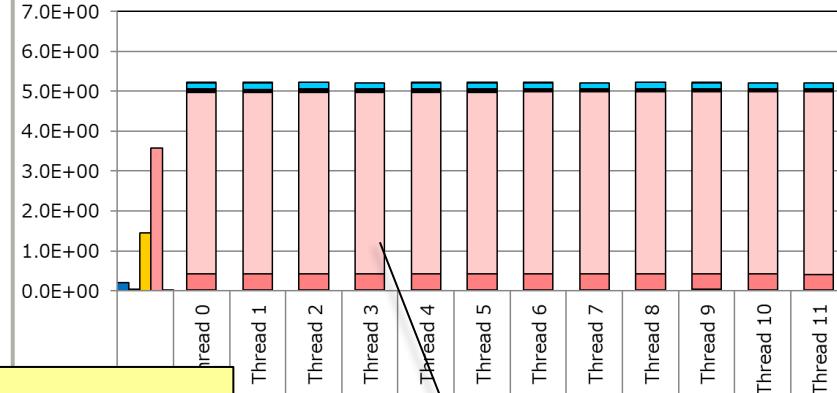
```

1      subroutine init(a,b,ie,n)
2      <<< Loop-information Start >>>
3      <<< [OPTIMIZATION]
4      <<< FUSED
5      <<< Loop-information End >>>
6          do i=1,n
7              if (mod(i,2).eq.0) then
8                  ie(i)=100000
9              endif
10             enddo
11
12
13
14
15
16      subroutine sub(a,b,s,ie,n)
17      real a(n),b(n),s
18      integer ie(n)
19      !$omp parallel do schedule(dynamic,1)
20      1 p
21      do j=1,n
22      <<< Loop-information Start >>>
23      <<< [OPTIMIZATION]
24      <<< SIMD(VL: 16)
25      <<< SOFTWARE PIPELINING(IPC: 3.50, ITR: 448,
26                               MVE: 7, POL: S)
27      <<< Loop-information End >>>
28      2 p 2v      do i=1,ie(j)
29      2 p 2v          a(i) = a(i)*b(i)*s
30      2 p 2v      enddo
31      1 p      enddo
32      end subroutine sub
33
34      program main
35      parameter(n=1000000)
36      call init(a,b,ie,n)
37
38      call sub(a,b,2.0,ie,n)

```



改善前



改善後

dynamicを指定することで
先に処理が終了したスレッドが次の処理を行うようになる。

ロードインバランスが改善された

処理量が不規則な場合、スケジュール方法が static のサイクリック分割ではロードインバランスが発生します。下の例ではロードインバランスのためバリア同期待ちが多くなっています。

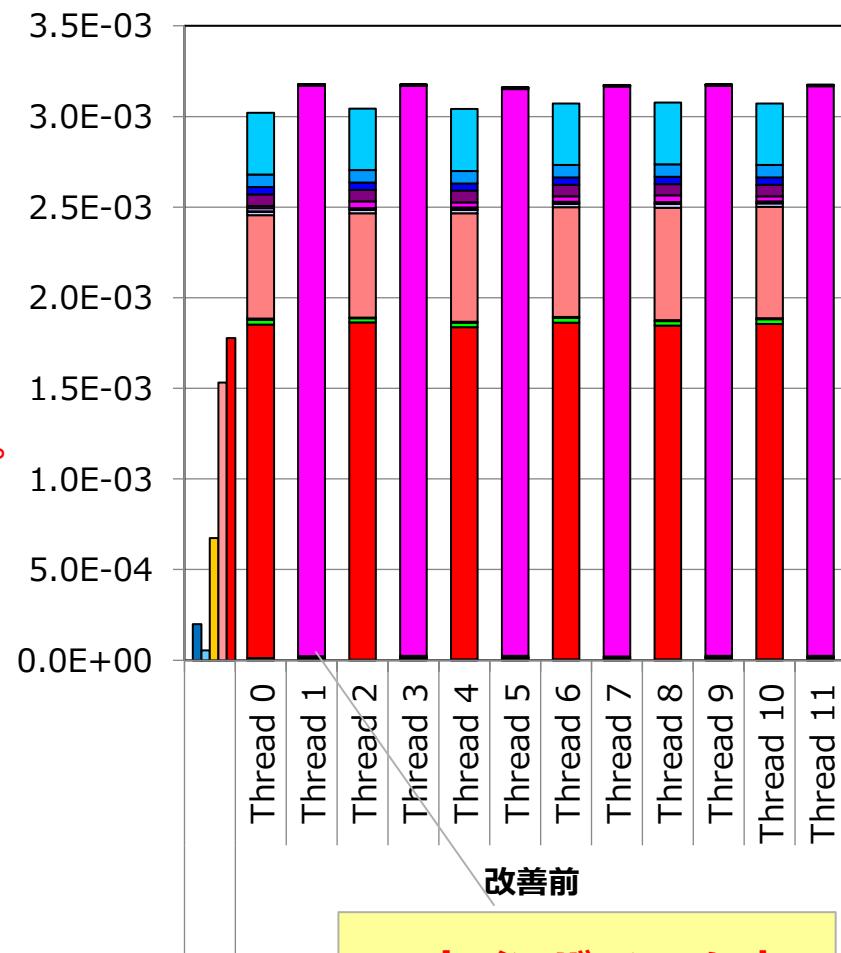
改善前ソース

```

28 void sub(int n, float (* restrict a)[n], float * restrict b,
           float s, int * restrict ie){
29     int i, j;
30
31     #pragma omp parallel for schedule(static, 1)
32     <<< Loop-information Start >>>
33     <<< [OPTIMIZATION]
34     <<< PREFETCH(HARD) Expected by compiler :
35     <<< (unknown)
36     <<< Loop-information End >>>
37     for (j = 0; j < n; j++){
38         <<< Loop-information Start >>>
39         <<< [OPTIMIZATION]
40         <<< SIMD(VL: 16)
41         <<< SOFTWARE PIPELINING(IPC: 3.50,
42             ITR: 448, MVE: 7, POL: S)
43         <<< PREFETCH(HARD) Expected by compiler :
44         <<< (unknown)
45         <<< Loop-information End >>>
46         for (i = 0; i < ie[j]; i++){
47             a[j][i] = a[j][i]*b[i]*s;
48         }
49     }
50 }
```

最内ループは制御変数 j が
奇数のときは1,000,000
回転、偶数ときは100,000
回転する。

評価ループ



ロードインバランスによる
バリア同期待ちが発生

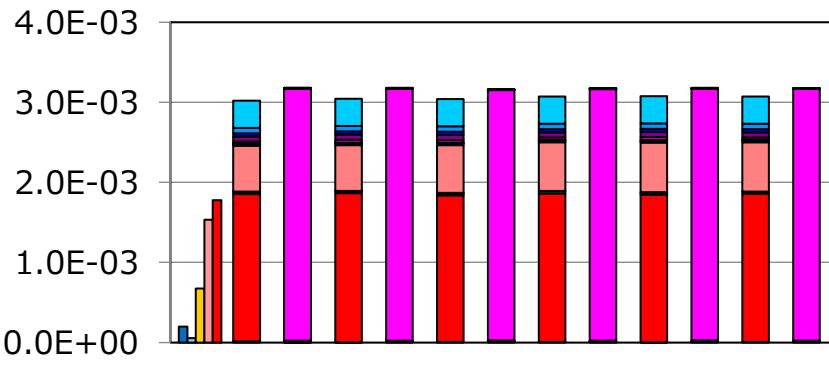
スケジュール方法をdynamicにすることで、先に処理が終了したスレッドが次の処理を行えるようになり、ロードインバランスが改善しました。

改善後ソース

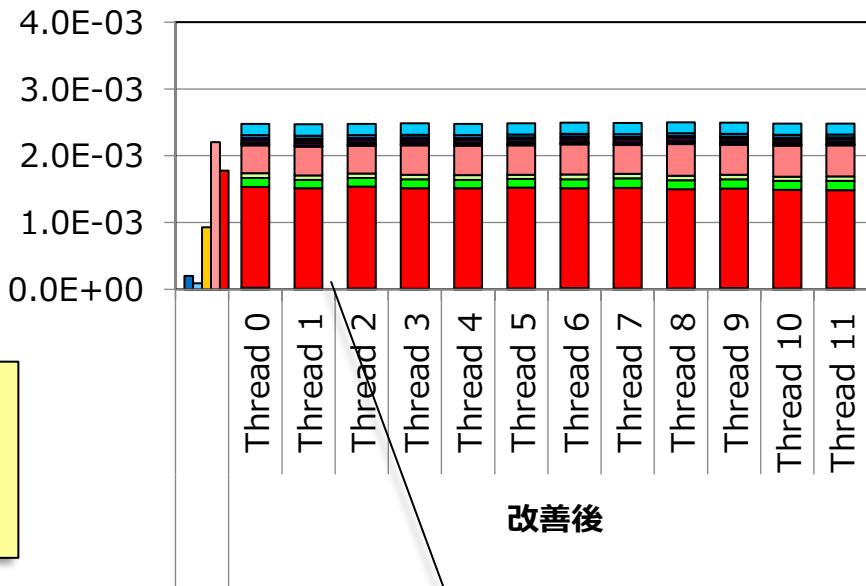
```

28 void sub(int n, float (* restrict a)[n],
           float * restrict b, float s, int * restrict ie){
29     int i, j;
30
31     #pragma omp parallel for schedule(dynamic, 1)
32     <<< Loop-information Start >>>
33     <<< [OPTIMIZATION]
34     <<< PREFETCH(HARD) Expected by compiler :
35     <<< (unknown)
36     <<< Loop-information End >>>
37     p   for (j = 0; j < n; j++){
38         <<< Loop-information Start >>>
39         <<< [OPTIMIZATION]
40         <<< SIMD(VL: 16)
41         <<< SOFTWARE PIPELINING(IPC: 3.50, ITR: 448,
42                               MVE: 7, POL: S)
43         <<< PREFETCH(HARD) Expected by compiler :
44         <<< (unknown)
45         <<< Loop-information End >>>
46     p   2v   for (i = 0; i < ie[j]; i++){
47     p   2v       a[j][i] = a[j][i]*b[i]*s;
48     p   2v   }
49     p   }
50 }
```

dynamicを指定することで
先に処理が終了したスレッドが次
の処理を行うようになる。



改善前



改善後

ロードインバランスが改善された

適切な並列化次元での並列化

- ・ 適切な並列化次元での並列化（改善前）
- ・ 適切な並列化次元での並列化（最適化制御行チューニング）
- ・ 適切な並列化次元での並列化（翻訳オプションチューニング）
- ・ 適切な並列化次元での並列化（OpenMPソースチューニング）

並列化次元のループ繰返し数が少なく翻訳時に繰返し数が不明な場合、繰返し数がスレッド並列数（例は12並列）未満となるケースではロードインバランスが発生します。そのため、バリア同期待ちが多くなっています。

改善前ソース

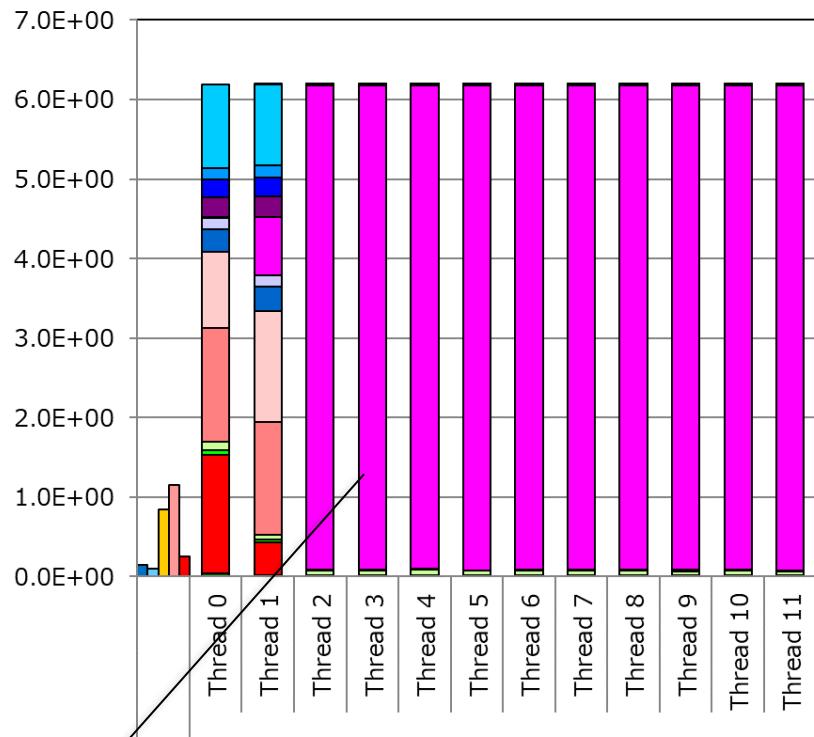
```

32   1 pp
      <<< Loop-information Start >>>
      <<< [PARALLELIZATION]
      <<< Standard iteration count: 2
      <<< Loop-information End >>>
          do k=1,l
      <<< Loop-information Start >>>
      <<< [OPTIMIZATION]
      <<< PREFETCH(HARD) Expected by compiler :
      <<< c, b, a
      <<< Loop-information End >>>
          do j=1,m
      <<< Loop-information Start >>>
      <<< [OPTIMIZATION]
      <<< SIMD(VL: 8)
      <<< SOFTWARE PIPELINING(IPC: 3.25,
                                ITR: 144, MVE: 4, POL: S)
      <<< PREFETCH(HARD) Expected by compiler :
      <<< c, b, a
      <<< Loop-information End >>>
34   3 p 2v      do i=1,n
35   3 p 2v      a(i,j,k)=b(i,j,k)+c(i,j,k)
36   3 p 2v      enddo
37   2 p         enddo
38   1 p         enddo

```

I=2
m=512
n=256

各スレッドのロードバランスが悪い



改善前

並列化次元の k の繰返し数が 2 回転で
あるためロードインバランスが発生

SERIAL指示子、**PARALLEL**指示子を指定することで、適切な次元で並列化されるようになり、ロードインバランスが改善しました。

改善後ソース

```

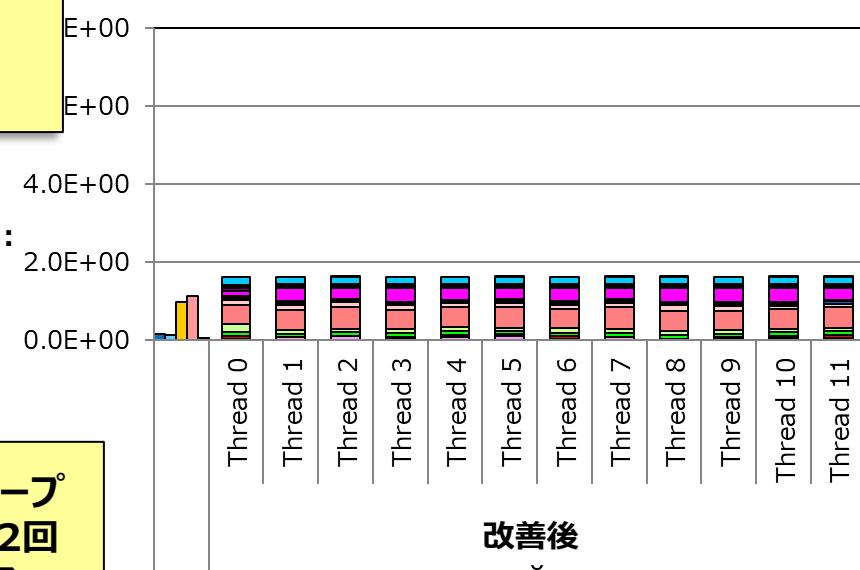
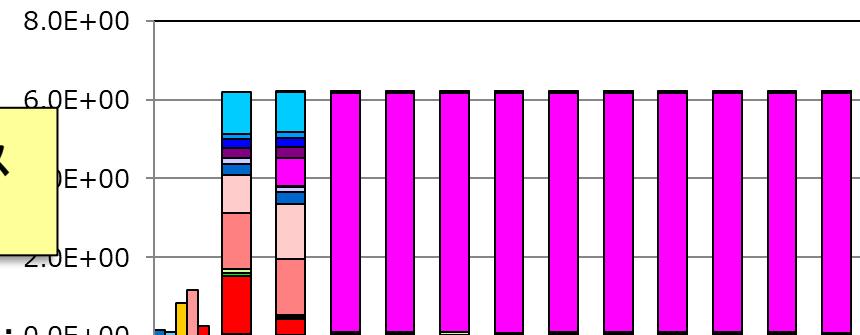
31      1      !ocl serial
32      1      do k=1,l
33      1      !ocl parallel
34      2 pp   <<< Loop-information Start >>>
35      3 p    <<< [PARALLELIZATION]
36      3 p    <<< Standard iteration count: 4
37      3 p    <<< [OPTIMIZATION]
38      2 p    <<< PREFETCH(HARD) Expected by compiler : 0.0E+00
39      1      <<< c, b, a
34      2 pp   <<< Loop-information End >>>
35      3 p    <<< c, b, a
36      3 p    do j=1,m
37      3 p    <<< Loop-information Start >>>
38      2 p    <<< [OPTIMIZATION]
39      1      <<< SIMD(VL: 8)
34      2 pp   <<< SOFTWARE PIPELINING(IPC: 3.25,
35      3 p    ITR: 144, MVE: 4, POL: S)
36      3 p    <<< PREFETCH(HARD) Expected by compiler :
37      3 p    <<< c, b, a
38      2 p    <<< Loop-information End >>>
39      1      do i=1,n
34      2 pp   a(i,j,k)=b(i,j,k)+c(i,j,k)
35      3 p    enddo
36      3 p    enddo
37      3 p    enddo
38      2 p    enddo
39      1      enddo

```

ループスライスを抑止する

I=2
m=512
n=256

並列化次元のループjの繰返し数512回で並列実行をする



翻訳オプション -Kdynamic_iteration を指定することで、適切な並列化次元が実行時に自動選択されるようになり、ロードインバランスが改善しました。

改善後ソース

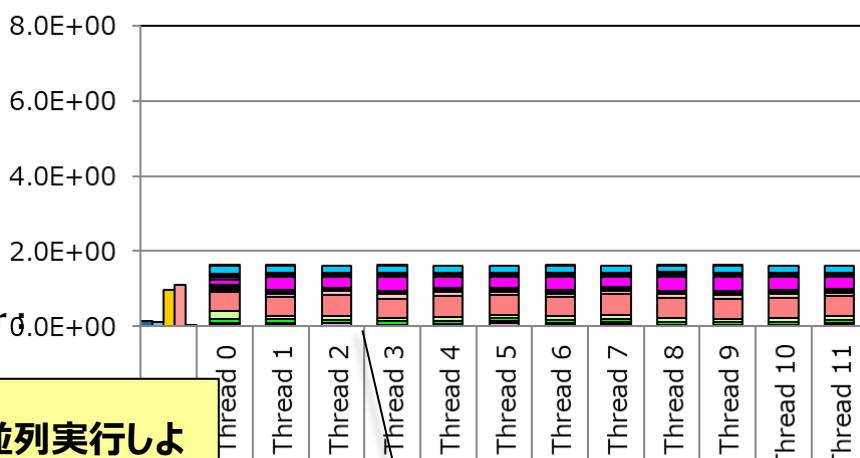
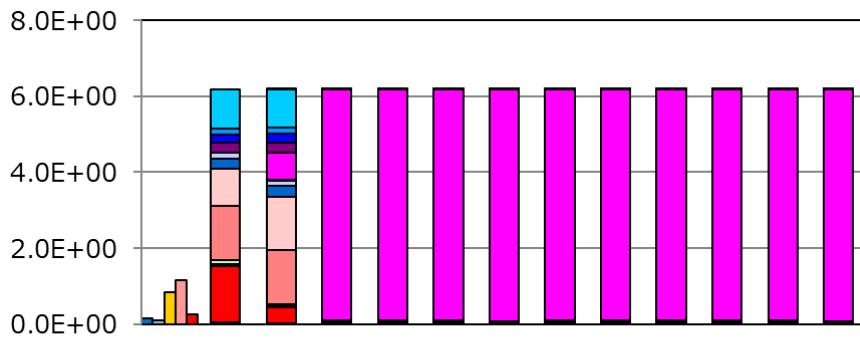
```

31   1 pp
      <<< Loop-information Start >>>
      <<< [PARALLELIZATION]
      <<< Standard iteration count: 2
      <<< Loop-information End >>>
          do k=1,l
      <<< Loop-information Start >>>
      <<< [PARALLELIZATION]
      <<< Standard iteration count: 4
      <<< [OPTIMIZATION]
      <<< PREFETCH(HARD) Expected by compiler :
      <<< c, b, a
      <<< Loop-information End
          do j=1,m
      <<< Loop-information Start
      <<< [PARALLELIZATION]
      <<< Standard iteration count: 728
      <<< [OPTIMIZATION]
      <<< SIMD(VL: 8)
      <<< SOFTWARE PIPELINING(IPC: 3.25,
                                ITR: 144, MVE: 4, POL: S)
      <<< PREFETCH(HARD) Expected by compiler
      <<< c, b, a
      <<< Loop-information End
          do i=1,n
              a(i,j,k)=b(i,j,k)
          enddo
      <<< 3 pp 2v
      <<< 3 p 2v
      <<< 3 p 2v
      <<< 2 p
      <<< 1 p

```

**I=2
m=512
n=256**

外側ループから並列実行しようとするが、ループ k では繰返し数が 2 回転と少ないため内側の 512 回転のループ j で並列実行される



ロードインバランスが改善された

OpenMPのCOLLAPSE指示節を指定することで、外側のループ[°]を一重化します。
その結果、ロードインバランスが改善しました。

改善後ソース

```

32      !$omp parallel do private(k,j,i) collapse(2)
<<< Loop-information Start >>>
<<< [OPTIMIZATION]
<<< PREFETCH(HARD) Expected by compiler :
<<< c, b, a
<<< Loop-information End >>>
33 1 p      do k=1,l
34 2 p      do j=1,m
<<< Loop-information Start >>>
<<< [OPTIMIZATION]
<<< SIMD(VL: 8)
<<< SOFTWARE PIPELINING(IPC: 3.25,
ITR: 144, MVE: 4, POL: S)
<<< PREFETCH(HARD) Expected by compiler :
<<< c, b, a
<<< Loop-information End >>>
35 3 p 2v      do i=1,n
36 3 p 2v      a(i,j,k)=b(i,j,k)+c(i,j,k)
37 3 p 2v      enddo
38 2 p      enddo
39 1 p      enddo
40 !$omp end par

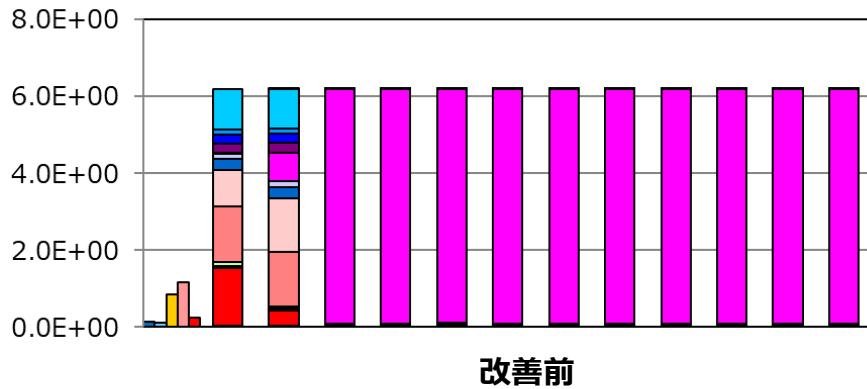
```

**I=2
m=512
n=256**

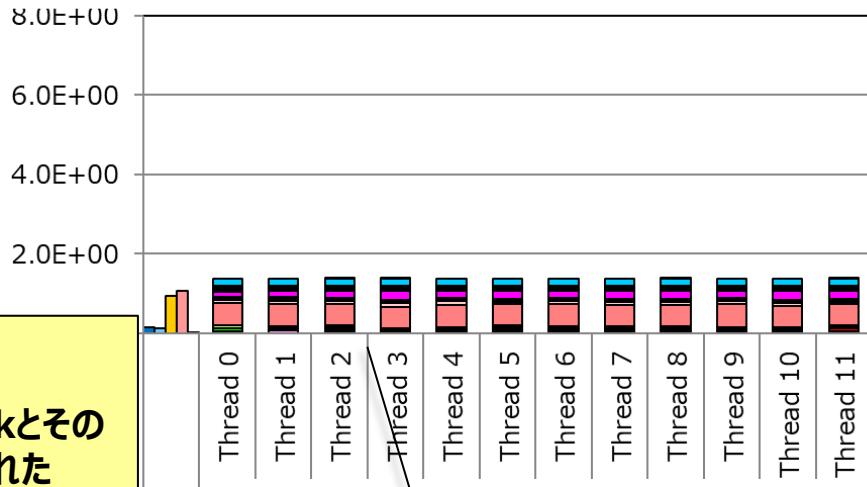
COLLAPSE指示節により
繰り返し数の少ないループkとその
内側のループjが一重化された

■ 注意事項

SIMD化軸(最内ループ)はCOLLAPSE
しないように注意してください。



改善前



改善後

ロードインバランスが改善された

並列化次元のループ繰返し数が少なく翻訳時に繰返し数が不明な場合、繰返し数がスレッド並列数（例は12並列）未満となるケースではロードインバランスが発生します。そのため、バリア同期待ちが多くなっています。

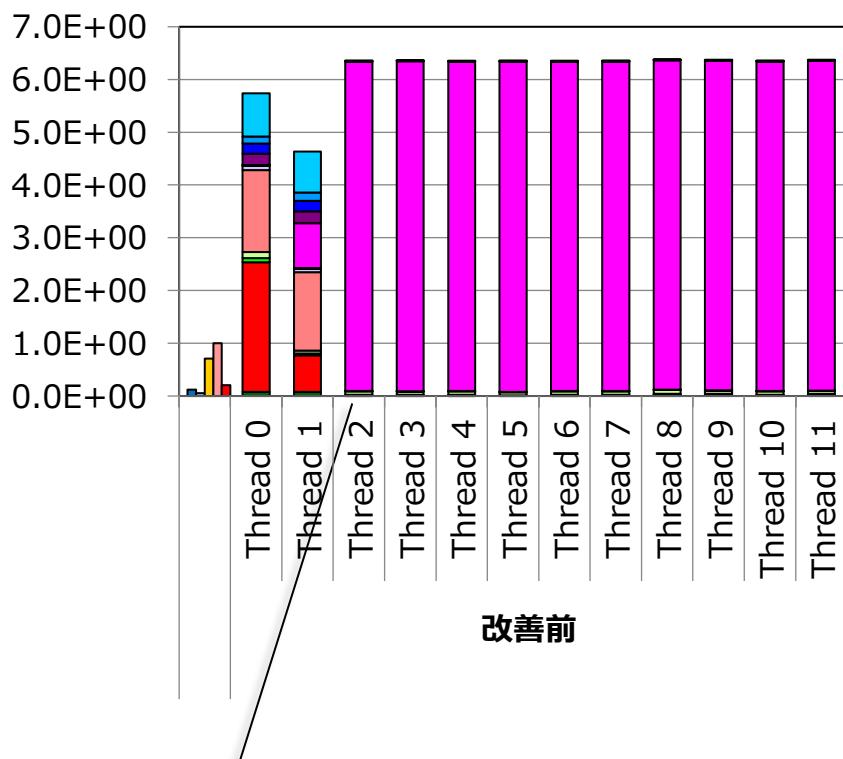
改善前ソース

```

42 p      #pragma omp parallel for
43   for (k = 0; k < l; k++) {
44     <<< Loop-information Start >>>
45     <<< [OPTIMIZATION]
46     <<< PREFETCH(HARD) Expected by compiler :
47       (unknown)
48     <<< Loop-information End >>>
49     for (j = 0; j < m; j++) {
50       <<< Loop-information Start >>>
51       <<< [OPTIMIZATION]
52       <<< SIMD(VL: 8)
53       <<< SOFTWARE PIPELINING(IPC: 3.00,
54           ITR: 144, MVE: 5, POL: S)
55       <<< PREFETCH(HARD) Expected by compiler :
56         (unknown)
57       <<< Loop-information End >>>
58     2v     for (i = 0; i < n; i++) {
59     2v       a[k][j][i] = b[k][j][i] + c[k][j][i];
60     2v     }
61   p     }
62   p   }
63 }
```

I=2
m=512
n=256

各スレッドのロードバランスが悪い



並列化次元の k の繰返し数が 2 回転であるためロードインバランスが発生

parallel for指示子やparallel指示子を指定することで、適切な次元で並列化されるようになり、ロードインバランスが改善しました。

改善後ソース

```

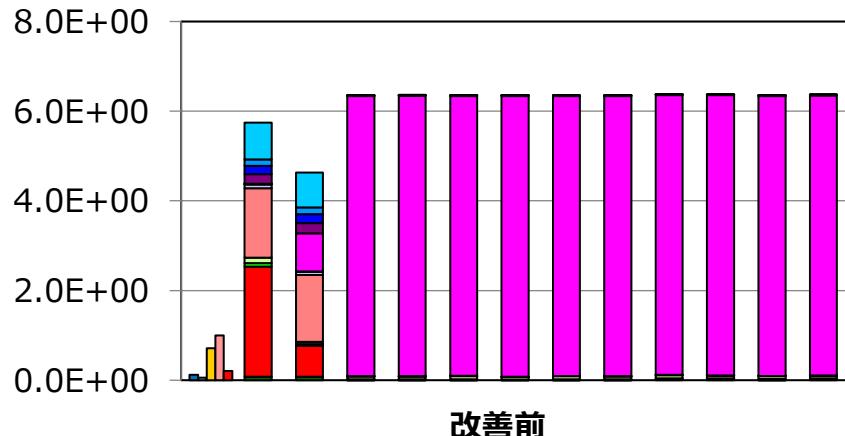
42      #pragma omp parallel private(i,j,k)
43      {
44          for (k = 0; k < l; k++){
45              #pragma omp for
46              <<< Loop-information Start >>>
47              <<< [OPTIMIZATION]
48              <<< PREFETCH(HARD) Expected by compiler :
49              <<< (unknown)
50              <<< Loop-information End >>>
51          }
52      }
53      p
54      2v      for (i = 0; i < n; i++){
55          2v          a[k][j][i] = b[k][j][i] + c[k][j][i];
56      }
57  }
```

l=2

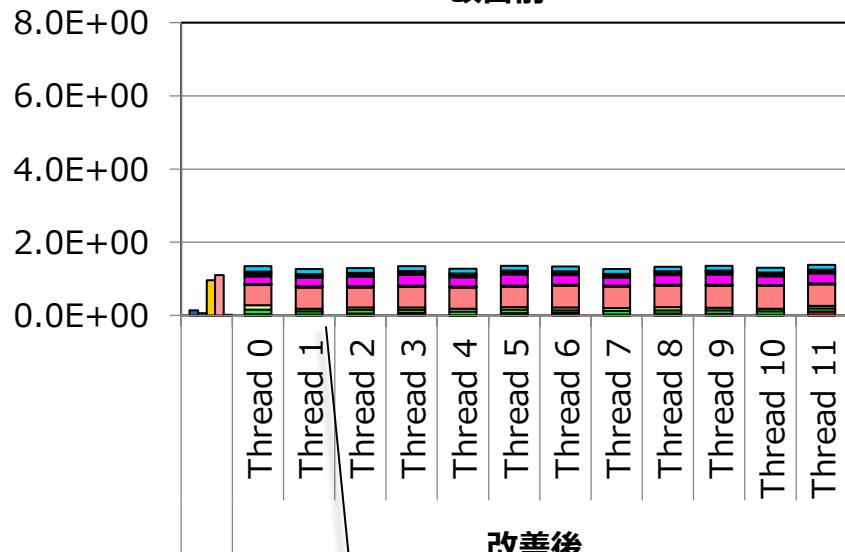
m=512
n=256

ループスライスを抑止する

並列化次元のループjの繰返し数512回で並列実行をする



改善前

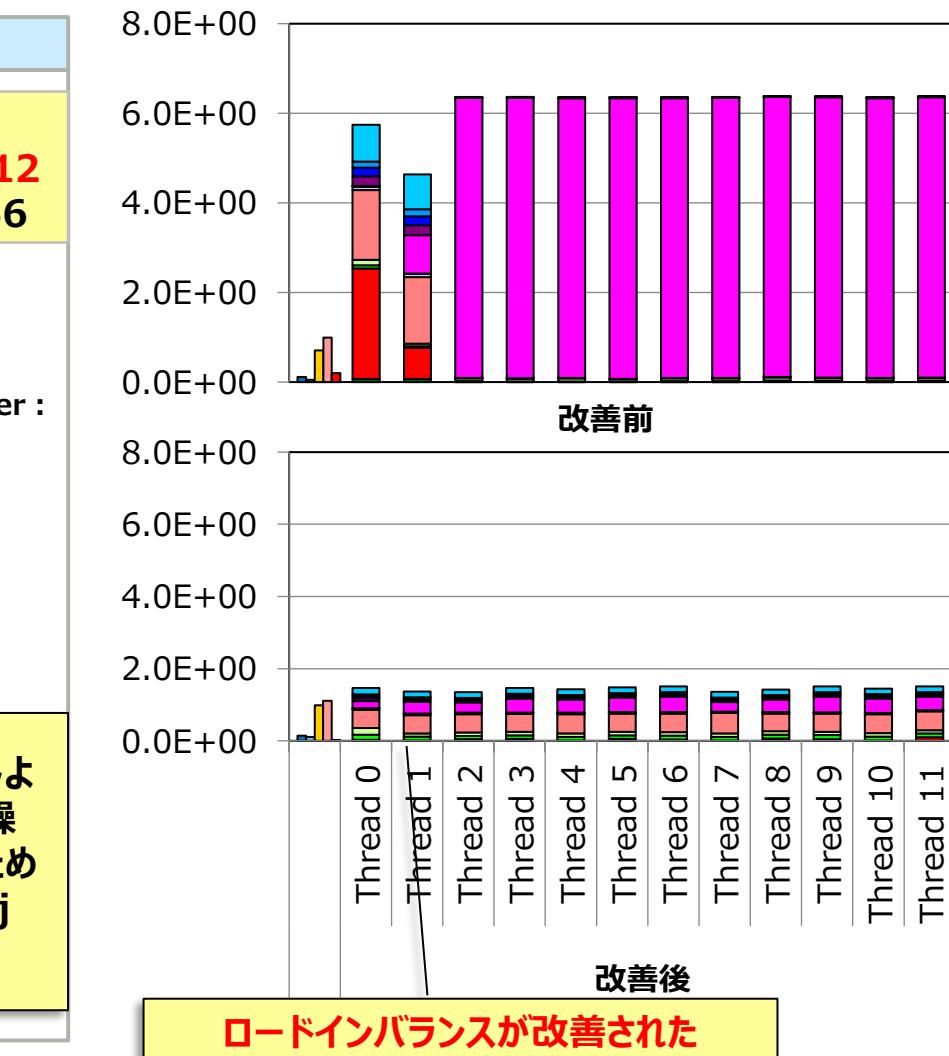


改善後

ロードインバランスが改善された

翻訳オプション **-Kdynamic_iteration** を指定することで、適切な並列化次元が実行時に自動選択されるようになり、ロードインバランスが改善しました。clangモードには-Kdynamic_iterationがないため、本手法によるチューニングはできません。

改善後ソース	
42 pp	<pre><<< Loop-information Start >>> <<< [PARALLELIZATION] <<< Standard iteration count: 3 <<< Loop-information End >>> for (k = 0; k < n3; k++){ <<< Loop-information Start >>> <<< [PARALLELIZATION] <<< Standard iteration count: 37 <<< [OPTIMIZATION] <<< PREFETCH(HARD) Expected by compiler : <<< (unknown) <<< Loop-information End >>> for (j = 0; j < n2; j++){ <<< Loop-information Start >>> <<< [PARALLELIZATION] <<< Standard iteration count: 552 <<< [OPTIMIZATION] <<< SIMD(VL: 8) <<< SOFTWARE PIPELINING(TPC-C 2.25 ITR: 1) <<< PREFETCH <<< (unknown) <<< Loop-information Start >>> 2v for (i = 0; 2v a[k][j][i] 2v } 2v }</pre>
43 pp	
44 pp	<p>外側ループから並列実行しようとするが、ループ k では繰返し数が 2 回転と少ないため内側の 512 回転のループ j で並列実行される</p>
45 p	
46 p	
47 p	
48 p	



OpenMPのcollapse指示節を指定することで、外側のループ[°]を一重化します。
その結果、ロードインバランスが改善しました。

改善後ソース

```

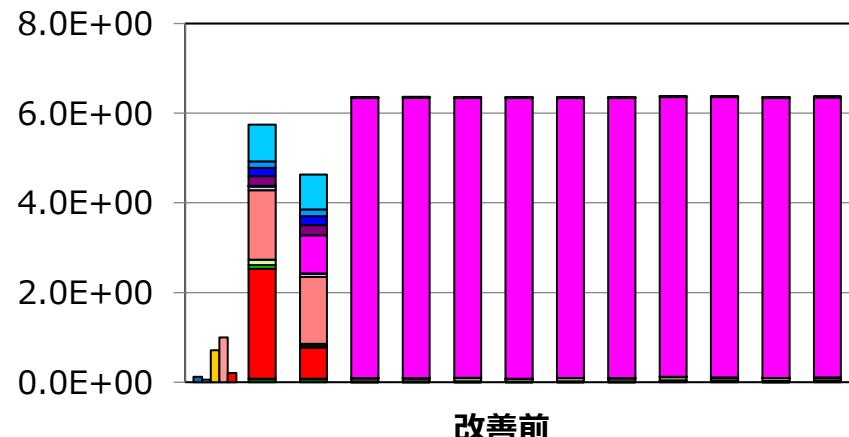
42      #pragma omp parallel for private(i,j,k) collapse(2)
43      <<< Loop-information Start >>>
44      <<< [OPTIMIZATION]
45      <<< PREFETCH(HARD) Expected by compiler :
46      <<< (unknown)
47      <<< Loop-information End >>>
48      p      for (k = 0; k < l; k++){
49      p          for (j = 0; j < m; j++){
50          <<< Loop-information Start >>>
51          <<< [OPTIMIZATION]
52          <<< SIMD(VL: 8)
53          <<< SOFTWARE PIPELINING(IPC: 3.00,
54              ITR: 144, MVE: 5, POL: S)
55          <<< PREFETCH(HARD) Expected by compiler :
56          <<< (unknown)
57          <<< Loop-information End >>>
58          p      2v      for (i = 0; i < n; i++){
59          p          a[k][j][i] = b[k][j][i] + c[k][j][i];
60          p      } } }
```

I=2
m=512
n=256

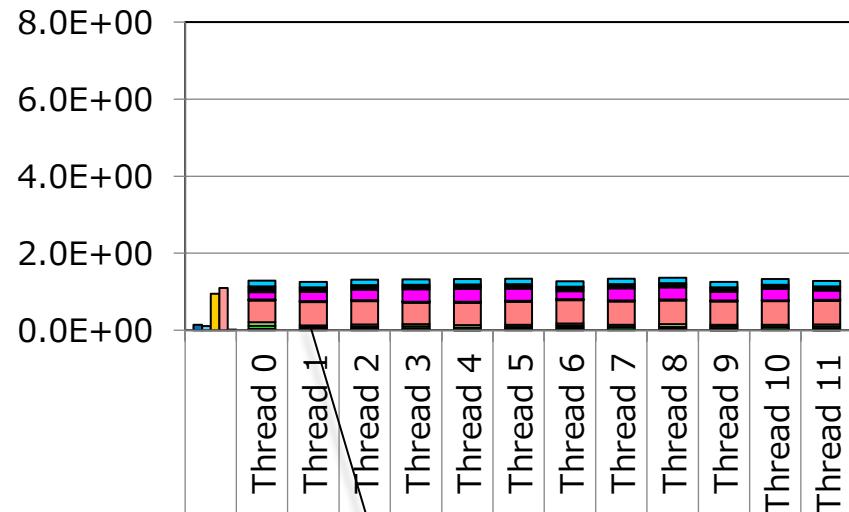
COLLAPSE指示節により
繰り返し数の少ないループ[°]kとその
内側のループ[°]jが一重化された

■ 注意事項

SIMD化軸(最内ループ[°])はCOLLAPSE
しないように注意してください。



改善前



改善後

ロードインバランスが改善された

ラージページの設定による実行効率の改善

- ・ ラージページページングポリシーの指定
- ・ ロックタイプの変更

ラージページページングポリシーの指定

- ・ ラージページページングポリシーについて
- ・ ラージページのページングポリシー(demand)の効果

■ XOS_MMM_L_PAGING_POLICY=prepage:demand:prepage

■ スレッド並列で複数CMGにて走行する時は、以下のメモリアロケーションとなる。

- ・プリページングでは、ロードモジュール起動時に、CMG0からデータが載る。
- ・デマンドページングでは、初回アクセス時に実行CMGにデータが載る。

■ 複数CMGをまたぐときはデマンドページングを推奨します。

環境変数名	指定値 (_付はdefault)	説明
XOS_MMM_L_PAGING_POLICY	[demand <u>prepage</u>]: [<u>demand</u> prepage]: [demand <u>prepage</u>]	各メモリ領域のページング方式(ページの割り当て契機)を選択する設定です。 「demand」はデマンドページング方式、「prepage」はプリページング方式を意味します。本変数はコロン(:)区切りで3つのメモリ領域のページング方式を指定します。 第1指定は、静的データの.bss領域です。(静的データの.data領域はページング方式指定の対象外で常にprepageとなります。) 第2指定は、スタック領域およびスレッドスタック領域です。 第3指定は、動的メモリ確保領域です。 指定値以外の値を指定した場合は、「prepage:demand:prepage」を指定したものとみなします。

ラージページのページングポリシー(demand)の効果

FUJITSU

プリページングではCMG0からデータが載るため、48スレッドのストリームの性能がでません。
デマンドページングに変更することで実行CMGにデータが載り、性能が大幅に向上了しました。

ソース

```
Subroutine sub(n,iter,x1,x2,y1)
  real(8) :: x1(n), x2(n), y1(n),c0
  integer n,i,k
  c0=2.0
  call fapp_start("sub",0,0)
  do k=1,iter
    !$omp parallel do
    do i=1,n
      y1(i) = x1(i) + c0 * x2(i)
    end do
  enddo
  :
  parameter(N=45000000,ITER=100)
  real*8 x1(N),x2(N),y1(N)
  call init(N,ITER,x1,x2,y1)
  call sub(N,ITER,x1,x2,y1)
```

	Memory throughput (GB/s)
prepage指定(default)	93GB/s
demand指定	804GB/s

翻訳オプション : -Kfast,openmp
-Kprefetch_sequential=soft -Kprefetch_line=9
-Kprefetch_line_L2=70 -Kzfill=18

ストリーム(データサイズは約1GB)

ロックタイプの変更

- ロックタイプ(XOS_MMM_L_arena_lock_type)とは
- ロックタイプの変更による効果

ロックタイプ(XOS_MMM_L_arena_lock_type)とは

■ XOS_MMM_L_arena_lock_type

- メモリ割り当てポリシーに関する設定です。
- 「0」はメモリ獲得性能優先。パラレルリージョン内でmallocしている場合に推奨します。
(メモリの獲得/開放がスレッド毎に独立したメモリプールから行われるようになり、デフォルト設定よりも排他制御のコストが減って性能改善する場合があります)
- 「1」はメモリ使用効率優先(デフォルト)。メモリ使用量が多い場合に推奨します。

ロックタイプの変更による効果

XOS_MMM_L_ARENA_LOCK_TYPE=0を指定することで、mallocの性能が向上しました。(実行時間は0.56秒から0.35秒に1.60倍の性能向上)

ソース

```
subroutine sub(n,m,iter,x1,x2,y2)
    integer(8) :: pZ1(iter)
    real(8) :: x1(n), x2(n), y2(n,m),c0
    c0=2.0
    !$omp parallel do shared(n,m,iter,x1,x2,c0,y2) private(pZ1,i,j,k) default(none)
    do k=1,m
        do j=1,iter
            pZ1(j) = malloc(8 * n)
        end do

        do i=1,n
            y2(i,k) = x1(i) + c0 * x2(i)
        end do

        do j=1,iter
            call free( pZ1(j))
        end do
    end do
end subroutine sub

program main
parameter(N=1048512,ITER=80)
real*8 x1(N),x2(N),y2(N,12)
call sub(N,12,ITER,x1,x2,y2)
end program main
```

メモリ使用量の削減

- OMP SINGLE からOMP MASTERへの書き換え

■ 書き換えによるメモリ使用量の削減

- `omp single`を`omp master`と`barrier`に書き換えすることで、実行スレッドを固定させ、メモリ領域の冗長使用を抑止します。

改善前ソース	改善後ソース
<pre>!\$omp parallel !\$omp single call loop(a,b,alpha,max); !\$omp end single !\$omp end parallel</pre>	<pre>!\$omp parallel !\$omp master call loop(a,b,alpha,max); !\$omp end master !\$omp barrier !\$omp end parallel</pre>

■ 更新履歴

版数	発行月	変更理由及び内容
初版(1.0版)	2020/9	・新規作成
1.3版	2021/3	・記事見直しによる誤字や表現の修正
1.4版	2021/8	・ソフトウェア版数アップによる差分修正、および、記事見直しによる誤字や表現の修正 ・「アンロールアンドジャムとは」ページの追加 ・「OMP SINGLE から OMP MASTER への書き換え」ページの追加
2.0版	2022/5	・C/C++のチューニングを追加 ・記事見直しによる誤字や表現の修正 ・フォーマットの変更
2.1版	2022/7	・記事見直しによる誤字や表現の修正
2.2版	2023/3	・記事見直しによる誤字や表現の修正

FUJITSU[∞]

