ULIBC ライブラリを用いた共有メモリ型並列アルゴリズムの高速化

安井 雄一郎^{1,2,a)} 藤澤 克樹^{1,2,b)} 竹内 聖悟^{3,4,c)} 湊 真一^{3,4,d)}

概要:現在,共有メモリ型並列計算機の主要なアーキテクチャとなっている NUMA (Non-Uniform Memory Access) アーキテクチャを有する計算機上には,各プロセッサと対となるローカルメモリと,他のプロセッ サと対となるリモートメモリが存在し,各プロセッサコアと各メインメモリ間の距離が異なるためアクセ スコストが均一でない. NUMA アーキテクチャ上での高速化にはコストの大きなリモートメモリへのアク セスを削減することが重要になるものの,並列実行時の各スレッドが計算機トポロジ上のどこに配置して いるかといった NUMA アーキテクチャを考慮した制御に必要な情報の取得は,既存ライブラリ群だけで は容易ではない. 本研究では CPU アフィニティやローカルメモリ確保などの機能をまとめたライブラリ ULIBC (Ubiquity Library for Intelligently Binding Cores) を開発し,いくつかの並列アルゴリズムに対し て, ULIBC を用いた NUMA アーキテクチャを考慮した高速化を示す.

キーワード:スレッド並列計算, CPU アフィニティ, グラフ処理, ZDD, 数理最適化問題

Fast implementation of shared-memory parallel algorithm using ULIBC (Ubiquity Library for Intelligently Binding Cores)

Yuichiro Yasui^{1,2,a)} Katsuki Fujisawa^{1,2,b)} Shogo Takeuchi^{3,4,c)} Shin-ichi Minato^{3,4,d)}

Abstract: Under NUMA (Non-uniform memory access) which is recent major shared memory multi-core architecture, each processor has its own local memory faster than non-local memory. Some speedup techniques which based on reducing access costs to the non-local memory larger than local memory considering NUMA architecture require development cost due to not enough to functions of some libraries for managing of memory access. Therefore, we implemented ULIBC (Ubiquity Library for Intelligently Binding Cores) for CPU affinity and local memory allocation. Finally, we show that some applications achieved performance improvements using ULIBC.

Keywords: thread parallel computing, CPU affinity, graph processing, ZDD, mathematical optimization

1. はじめに

現在, 共有メモリ型並列計算機の主要な NUMA (Non-Uniform Memory Access) アーキテクチャを有する計算機

- 3 北海道大学
- Hokkaido University
- 4 JST ERATO
- a) yasui@indsys.chuo-u.ac.jp
- ^{b)} fujisawa@indsys.chuo-u.ac.jp
 ^{c)} takeuchi@erato.ist.hokudai.ac.jp
- ^{d)} minato@ist.hokudai.ac.jp

上には、プロセッサと対となるローカルなメインメモリ (ローカルメモリ)と、他のプロセッサと対となるメインメ モリ (リモートメモリ)が存在する.プロセッサとローカ ルメモリの対は NUMA ノードと呼ばれ、一般的に NUMA ノード内のローカルメモリへのアクセスコストは、NUMA ノードを横断した他のプロセッサと対となるリモートメモ リへのアクセスコストよりも小さい.このように以前まで の UMA (Uniform Memory Access) アーキテクチャを有す る計算機での均一なメモリアクセスと比べて、NUMA アー キテクチャを有する計算機でのメモリアクセスは複雑かつ 不均一化している.そのため、並列数の増加に従う並列効

¹ 中央大学

Chuo University

 ² JST CREST
 ³ 北海道大学

率の低下を抑えることは難しく、特に演算量に比べてデー タ移動量が大きい並列アルゴリズムではその傾向が強い.

NUMA アーキテクチャを考慮した高速化に関する先行 研究 [1], [5], [6], [7], [9] では、コストの大きなリモートメ モリへのメモリアクセスを削減し、ローカルメモリへの メモリアクセスの局所性を高めて、性能改善を達成してい る. それらの制御には CPU アフィニティ設定やメモリ 領域のローカルメモリへの固定などが必要だが、Portable Hardware Locality (hwloc) [2], Likwid [3], Intel コンパ イラの Thread Affinity Interface [4], OpenUH コンパイ ラ [5], Linux 上での numactl コマンド, libnuma ライブラ リ, sched_{get,set}affinity や mbind などのシステム コール関数などが既に存在する. しかしながら、これらの 環境だけでは著者らが文献 [9] で提案したアルゴリズムの 実装に必要な「スレッド並列での実行中の各スレッドが何 番目の NUMA ノードの何番目のコア上に固定されるか」 といった情報の取得は容易に実現することができない.

そこで本研究では、NUMA アーキテクチャ計算機上で汎 用的な高速化を行うための、アフィニティ設定やローカルメ モリ設定などの機能をまとめたライブラリ ULIBC (Ubiquity Library for Intelligently Binding Cores)を開発した. ULIBC は計算機上の CPU ソケット、物理コア、論理コア などの CPU トポロジ情報を取得し、トポロジ情報からス レッドと論理コアの割当表を作成する.その後、割当表を 参照して、ローカルメモリへのメモリ確保や、スレッド並列 時に各スレッドを論理コア上への固定を行う.ULIBC は割 当表の参照によるシンプルな制御を行っており、以下のよ うな利点が挙げられる.

- 各スレッドの論理コアへの割当は事前に決定するため、 スレッド並列時における各スレッドが固定される位置の把握や性能予測が容易である
- 並列時のアフィニティ設定のためのオーバーヘッドは システムコールの発行コスト以外ほぼ発生しない
- 割当表を参照するための関数を挿入すれば良いため、
 実装上の手間を抑えられる

加えて本論文では, いくつかの性質の異なるいくつかの 並列アルゴリズムに対し, ULIBC を用いた NUMA アーキ テクチャを考慮した高速化を示す.対象は以下の4種類で いずれも高速化を達成し,計算機上のボトルネック箇所の 特定を行った.

- メモリバンド幅に対する STREAM ベンチマーク
- Graph500 ベンチマークの対象であるグラフに対する 並列幅優先探索 (BFS) [9]
- グラフ上の閉路を持たないパスを列挙する並列アルゴ リズム Edge-based approach (EBA) [13]
- 近年, 最も注目されている数理最適化問題の一つである 半正定値計画問題に対する高性能汎用ソルバ SDPA [14]

2. NUMA アーキテクチャ上でのメモリアクセス

NUMA アーキテクチャにおける並列アルゴリズムのメ モリアクセスに関して説明を行う.まず図 1 は NUMA アーキテクチャである Intel Xeon E5-4640 が 4 基搭載さ れた共有メモリ型計算機環境 4-way Sandybridge-EP であ る.この計算機環境において,まず 4 基ある CPU ソケッ トはそれぞれ 8 個の物理コアを持ち,各物理コアは個別の L1/L2 キャッシュを持ち,同時に 2 つのスレッド (論理コ ア)を実行することが可能である.また,ソケット上には 16 の論理コア (8 の物理コア) で共有する L3 キャッシュ が配置している.



図1 計算機環境 4-way Sandybridge-EP

続いて、4-way Sandybridge-EP 上での動的メモリ確 保に関して説明する.特殊な設定を行っていなければ (get_mempolicy 関数で得られるメモリ確保のポリシーが MPOL_DEFAULT ならば), C/C++ 言語の malloc/new 関数 を用いた動的メモリ確保は、関数呼び出しを行ったコアの ローカルメモリへ連続した領域を試みる. その際, 確保す る領域が1つのローカルメモリ上に収まりきらなければ、 図2のように他のメモリに対象を移し、先程と同様に連続 した領域を試みる. このように確保された領域はサイズに よっては複数のローカルメモリを横断してしまい、コスト が小さいローカルメモリへのアクセスと、コストが大きい リモートメモリへのアクセスが混在することになり, 負荷 分散が難しく並列効率が下がる原因となる. そのため, 並 列計算での各スレッド間でのメモリアクセス要求の衝突や、 各 NUMA ノード間を接続するインターコネクトに非常に 大きな負荷がかかることになる.しかしながら、実際には メモリ領域はページ単位でファーストタッチに基づいて確 保されるため, 確保されるメモリ領域がいずれの NUMA ノードに配置されるかはどのように初期化するかに依存 する.



図2 一般的なメモリ確保と複数ローカルメモリへのメモリアクセス

図3のようにアクセスするメモリ領域を各ローカルメモ リに分割することができ、メモリアクセスの局所性が高い アルゴリズム設計を行うことができるならば、各スレッド がアクセスコストの小さいローカルメモリへのアクセスす るように制御を行うことで、不均質なデータアクセスによ る性能低下やメモリアクセス要求の衝突による性能低下を 抑えることができる.



図3 NUMA を考慮したメモリ確保とローカルなメモリアクセス

3. ULIBC ライブラリ

本研究で実装したライブラリ ULIBC は, CPU アフィニ ティとローカルメモリへの動的メモリ確保などの制御を容 易に行うためのライブラリである.類似のライブラリ等は 既に存在するものの,既存のライブラリだけでは「スレッ ド並列計算時に各スレッドを実行する論理コアが何番目の CPU ソケットの何番目の物理コア上に固定されるか」と いった情報を取得することは容易ではない.

詳細なトポロジ情報の取得には Portable Hardware Locality (hwloc) [2], Likwid [3] が挙げられるが, アフィニ ティ設定を行った後に実際にどのコア上に固定されている かという機能は提供されていない. また, コンパイラによる アフィニティ設定として, GCC コンパイラの OpenMP ラ イブラリ, Intel コンパイラ Thread Affinity Interface [4], OpenUH コンパイラ [5] などが挙げられるが, GCC と ICC では環境変数によるアフィニティ設定のみで, OpenUH で はアフィニティ設定やローカルメモリへのメモリ確保の 両方を提供しているものの独自プラグマ文による指定と なり環境に依存してしまう. コマンドラインツールによる HPCS2014 2014/1/8

アフィニティ設定やローカルメモリへのメモリ確保設定 として、Linux 上での numactl コマンドが挙がられるが、 プログラムの実行者が事前に計算機トポロジ情報を把握 している必要がある. さらに一般的な Linux 上で利用可 能な libnuma ライブラリ、sched_{get,set}affinity や mbind などのシステムコール関数を用いると、自由な制御 ができるものの実装コストが大きく実装が煩雑になってし まう.

既存のライブラリに対し ULIBC は、一般的な Linux 環 境を対象とし、ULIBC を適用したアプリケーションが実行 すると、まず計算機トポロジ情報(図4の processor topology)の自動的取得を行い、各スレッドの論理コアへの割当 表(図4の mapping table)を作成する.アフィニティ設定 とローカルメモリへのメモリ確保はこの割当表を参照して 実現する.なお、ULIBC の利用に管理者権限は不要である. (1)計算機トポロジの取得.

- (2) スレッド並列時に参照するための, スレッド ID と論 理プロセッサ ID の割当表の作成.
- (3) (2) を参照しローカルメモリ上に固定した動的メモリ 領域確保 (mmap 関数と mbind 関数による制御).
- (4) スレッド並列時に, (2) を参照して各スレッドを論理 コア上に固定 (sched_setaffinity 関数による制御) し, (1) を参照して (3) で確保したメモリ領域にアクセ スするように制御.



図 4 スレッド ID と論理プロセッサ ID の割当表 (mapping table) と計算機トポロジ (processor topology)

3.1 計算機トポロジの取得

ULIBC が自動的に行う計算機トポロジの取得について 説明する.まず各論理コアを識別するための ID について 表1にまとめる.

表 1 計算機トポロジを表現するための ID (括弧内の数字は 4-way Sandybridge-EP 上で取り得る範囲を表している)

ID	説明
Processor ID	論理 CPU コアに対する ID (0,1,,63)
Package ID	CPU ソケットに対する ID (0,1,2,3)
	(NUMA ノードに対応)
Core ID	CPU ソケット内の物理コアに対する ID (0,1,,7)
SMT ID	物理コア上での論理コア番号 (0,1)

計算機毎に CPU ソケット数, ソケット数内の物理コア

数, 物理コア上の論理コア数は異なり, プロセッサコア ID, CPU ソケット ID, ソケット内の ID の順序は異なる. 一 般的な Linux 環境を想定すると, 汎用的な計算機トポロジ の取得には以下のような方法を挙げることができる. いず れの取得方法でも同様の情報を取得することが可能で, 計 算機トポロジの取得に要する実行時間はプログラムの性能 に影響を及ぼさないほど小さい.

- (a) デバイスファイル /proc/cpuinfo
 Processor ID, Package ID, Core ID と, processor,
 physical id, core id がそれぞれ対応.
- (b) ディレクトリ構成 /sys/devices/system/ [9]
 Processor ID, Package ID, Core ID と, cpu, node
 / physical_package_id, core_id がそれぞれ対応.
 物理コア上に動作する Processor ID をまとめた
 thread_siblings_list が提供される.
- (c) APICID のサブセット情報 [10]
 各論理コアに一意の ID となる APIC ID のサブセット 情報により, Processor ID, Package ID, Core ID, SMT ID を取得できる.

表 2 は, Hyper-threading 機構を有効とした論理コアが 64 となる 4-way Sandybridge-EP 上の計算機トポロジを 示した Processor ID, Package ID, Core ID の対応表であ る. これらはアプリケーションが実行時に ULIBC が自動 的に取得するものである.

表	2	4-way	Sandybridge-EP	上の計算機ト	ポロジ
---	----------	-------	----------------	--------	-----

	Core ID							
	0	1	2	3	4	5	6	7
Package ID			I	Proces	sor II)		
0	0	1	2	3	4	5	6	7
0	32	33	34	35	36	37	38	39
1	8	9	10	11	12	13	14	15
1	40	41	42	43	44	45	46	47
9	16	17	18	19	20	21	22	23
2	48	49	50	51	52	53	54	55
9	24	25	26	27	28	29	30	31
	56	57	58	59	60	61	62	63

3.2 割当表を用いたアフィニティ設定

ULIBC におけるアフィニティ設定について説明を行う. Linux 環境での各スレッドの論理コアへの固定には sched_setaffinity 関数を用いて Processor ID (論理 CPU コア ID) で指定し,動的確保したメモリ領域を特定 のローカルメモリ上に固定するための関数 mbind を用いて Package ID (NUMA ノード ID) で指定する. ULIBC では, スレッド並列計算時の各スレッドに与えた一意のスレッド ID と Processor ID との割当表を作成し制御を行う.

ULIBC ではスレッドと論理コアを割り付けるアフィニ ティ設定の方針を3種類 Scatter, Compact, Compact+ を 用意している. Scatter では使用するスレッドを NUMA ノードに分散するように、Compact では使用するスレッド を NUMA ノードに集約するように、それぞれスレッドを 各論理コアに割り当てる.また Compact+は、Compact と 似たアフィニティ設定だが、物理コア上へのスレッドの重 複 (Hyper-Threading 機構の利用) は回避する割当となる. 表 3.2 と表 3.2 に 4-way Sandybridge-EP 上のアフィニ ティ設定毎のスレッド ID と Processor ID の割当表、ア フィニティ設定毎のスレッド数と使用する NUMA ノード 数、NUMA ノード内の物理コア数、物理コア上で同時実行 するスレッド数をまとめる.

表 3 4-way Sandybridge-EP 上のアフィニティ設定毎のスレッド ID と Processor ID の割当表

(a)	Septton	77	,		,
(a)	Scatter	11	1	-7	ィ

	(-)							
スレッド ID			I	Proces	sor Il)		
00 15	00	08	16	24	01	09	17	25
00 - 15	02	10	18	26	03	11	19	27
16 21	04	12	20	28	05	13	21	29
10 - 31	06	14	22	30	07	15	23	31
20 47	32	40	48	56	33	41	49	57
52 - 41	34	42	50	58	35	43	51	59
48 - 63	36	44	52	60	37	45	53	61
	38	46	54	62	39	47	55	63

(b) Compact アフィニティ

スレッド ID			I	Proces	sor II)		
00 15	00	01	02	03	04	05	06	07
00 - 15	32	33	34	35	36	37	38	39
16 91	08	09	10	11	12	13	14	15
10 - 51	40	41	42	43	44	45	46	47
20 47	16	17	18	19	20	21	22	23
52 - 41	48	49	50	51	52	53	54	55
48 - 63	24	25	26	27	28	29	30	31
	56	57	58	59	60	61	62	63

(c) Compact+ アフィニティ

	. ,							
スレッド ID			Η	Proces	sor II)		
00 15	0	1	2	3	4	5	6	7
00 - 15	8	9	10	11	12	13	14	15
16 21	16	17	18	19	20	21	22	23
10 - 31	24	25	26	27	28	29	30	31
20 47	32	33	34	35	36	37	38	39
32 - 41	40	41	42	43	44	45	46	47
48 - 63	48	49	50	51	52	53	54	55
	56	57	58	59	60	61	62	63

3.3 ULIBC を用いてアフィニティ設定を行った STREAM ベンチマーク

STREAM は大きさが n の実数ベクトル $a, b, c \in \mathbb{R}^n$ に 対する 4 種類の演算 Copy, Scale, Add, Triad に要する実 行時間を用いたメモリ帯域幅に対するベンチマークである. なおいずれのベクトルも倍精度浮動小数点型 (double 型) の大きさ n の配列として, scalar は実数変数 (3.0) として 実装されている.

表 4 4-way Sandybridge-EP 上のアフィニティ設定毎のスレッド 数と使用する NUMA ノード数, NUMA ノード内の物理コア 数,物理コア上で同時実行するスレッド数

(a) Scatter	・アフィニティ
-------------	---------

スレッド数	1	2	4	8	16	32	64
NUMA ノード数	1	2	4	4	4	4	4
NUMA ノード内の物理コア数	1	1	1	2	4	8	8
物理コア上のスレッド数	1	1	1	1	1	1	2
(b) Compa	ctア	・フィ	ニテ	イ			
スレッド数	1	2	4	8	16	32	64
NUMA ノード数	1	1	1	1	1	2	4
NUMA ノード内の物理コア数	1	2	4	8	8	8	8
物理コア上のスレッド数	1	1	1	1	2	2	2
(c) Compac	t+ 7	アフィ	ィニラ	- 1			
スレッド数	1	2	4	8	16	32	64
NUMA ノード数	1	1	1	1	2	4	4
NUMA ノード内の物理コア数	1	2	4	8	8	8	8
物理コア上のスレッド数	1	1	1	1	1	1	2
~			<i>(</i> .	~ ~	-		>

Copy	$c_i =$	a_i ,	(i = 0, 1,, n - 1)
Scale	$b_i =$	scalar $\cdot c_i$,	$\left(i=0,1,,n-1\right)$
Add	$c_i =$	$a_i + b_i$,	$\left(i=0,1,,n-1\right)$
Triad	$a_i =$	$b_i + \operatorname{scalar} \cdot c_i,$	(i = 0, 1,, n - 1)

まず,図5にC/C++ 言語上で malloc/new 関数を用い てメモリ確保を行った Triad 演算の実装例を示す. この場 合,図2で示したように,特定のローカルメモリ上に固定 するような制御は行っていないため, スレッドに応じてア クセスコストの異なる不均質なメモリアクセスが多発する ことが予測される.しかしながらメモリ確保はファースト タッチに基づいて確保されるため, 初期化時のスレッドと 同じスレッドで演算を行うことができればローカルなメモ リアクセスとなる.

```
#pragma omp parallel for
 for (long i = 0; i < n; ++i) {
   A[i] = B[i] + scalar * C[i];
```

図5 OpenMP を用いたスレッド並列 Triad 演算 (アフィニティ設 定なし w/o)

続いて,図6はULIBCの関数(図7)を用いて,アフィ ニティ設定とローカルメモリへのメモリ確保設定を行った Triad 演算である. まず (1) set_affinity_policy 関数で アフィニティ設定 {COMPACT, COMPACT_PLUS, SCATTER} のいずれかを指定し, (2) get_online_numa_nodes 関数で 取得したスレッド並列時に実行される NUMA ノード数 分の固定されたローカルメモリ領域を (8) lmalloc 関数 を用いて動的メモリ確保する. その後, スレッド並列領 域の先頭で OpenMP の omp_get_thread_num 関数を用 いてスレッド ID を取得し, アフィニティ設定で割り当て られる Processor ID, Package ID, NUMA ノード内の各

論理コアに割当表での割り当て順に与えた通し ID をそ れぞれ(3) get_numa_procid 関数,(4) get_numa_nodeid 関数, (7) get_numa_ordcoreid 関数を用いて取得する. 各スレッドは得られた Processor ID で指定した論理コ ア上に (9) set_affinity 関数を用いて固定する. その 後,各 NUMA ノード上に配置されたスレッド数を (2) get_online_numa_cores 関数を用いて取得し, 各スレッ ドが担当する配列の範囲を計算して、各スレッドは割り当 てられた Triad 演算を行う. また, スレッド計算後には (10) clear_affinity 関数を用いて、 プログラムの実行開 始時に取得したデフォルトのアフィニティ設定で現在の アフィニティ設定を上書きし、以後の処理に現在のアフィ ニティ設定が影響が出ないように初期化する.また、この 例では使用していない (5) get_numa_coreid 関数や (6) get_numa_smtid 関数は、それぞれ スレッド ID に対応す る物理コア ID (Core ID) や, 物理コア上で実行されるス レッド ID (SMT ID) を取得することができる.

```
long N[MAX_NODES];
double *A[MAX_NODES], *B[MAX_NODES], *C[MAX_NODES];
set_affinity_policy(SCATTER);
for (int k = 0; k < get_online_numa_nodes(); ++k) {</pre>
 A[k] = lmalloc(sizeof(double) * N[k], k);
 B[k] = lmalloc(sizeof(double) * N[k], k);
  C[k] = lmalloc(sizeof(double) * N[k], k);
#pragma omp parallel
  int id = omp_get_thread_num();
  int procid = get_numa_procid(id);
  set_affinity(procid);
  int nodeid = get_numa_nodeid(id);
  int coreid = get_numa_ordcoreid(id);
  int lnp = get_numa_online_cores(nodeid);
  double *1A = A[nodeid], *1B = B[nodeid], *1C = C[nodeid];
  const long q = N[nodeid] / lnp, r = N[nodeid] % lnp;
  const long ls = q * (coreid+0) + min(r, coreid+0);
  const long le = q * (coreid+1) + min(r, coreid+1);
  for (long j = ls; j < le; ++j) {
   1A[i] = 1B[i] + scalar * 1C[i];
 7
  clear_affinity();
```

図 6 ULIBC による Scatter アフィニティ設定を適用した STREAM ベンチマークの Triad 演算

図8に、以上のように実装した STREAM ベンチマーク の Triad 演算に対する 4 種類アフィニティ設定 w/o (ア フィニティ設定なし), Scatter, Compact, Compact+毎の 実行結果 (メモリ帯域幅 GB/s) をまとめる. 要素数は 2²⁷ とし、用いた計算機環境は 4-way Sandybridge-EP である. Scattar では, 1, 2, 4 とスレッド数の増加に従い使用する CPU ソケットも 1, 2, 4 と増加する. 続いて 4, 8, 16, 32 とスレッド数を変化させると、各ソケットに 1, 2, 4, 8 コ

2014年ハイパフォーマンスコンピューティングと計算科学シンポジウム High Performance Computing Symposium 2014

(1) int set_affinity_policy(int policy) (2) int get_online_numa_nodes(void) (3) int get_numa_procid(int thread_id) (4) int get_numa_nodeid(int thread_id) (5) int get_numa_coreid(int thread_id) (6) int get_numa_smtid(int thread_id) (7) int get_numa_ordcoreid(int thread_id) (8) void *lmalloc(size_t sz, int node_id) (9) int set_affinity(int processor_id) (10)void clear_affinity(void) 図7 ULIBC で提供する関数群

アずつ使用した 4 ソケット実行となる.一方, Compact, Compact+ では、1、2、4、8 とスレッド数の増加に従い、1、 2,4,8コアを使用した1ソケット実行となる.その後, Compact では 16 スレッド時に各物理コア上に 2 スレッ ドずつ固定した配置となり、32、64 とスレッド数の増加に 従い, 2, 4 ソケット実行となる. また, Compact+ では 16, 32 スレッド時には各物理コアに1つずつ配置した 2,4 ソ ケット実行となる. いずれも 64 スレッド時には,各物理コ ア上に2スレッドずつ固定した配置 (Hyper-Threading 機 構)と同じ配置となる.以上より,使用するソケットが最大 になるようにコア配置を行う Scatter が最も高いメモリ帯 域幅を示し、 続いて、 Compact+, Compact となった. この ように、メモリバンド幅に対しては有効とはならなかった ものの、各スレッド間に頻繁な通信が必要となる場合には、 Compact もしくは Compact+ による性能向上に期待でき る. また, アフィニティ設定を行わなった場合は, 32, 64 ス レッド時で性能低下が起きるが, ULIBC によるアフィニ ティ設定を行うことで、メモリアクセスの不均一化による 性能低下を抑えていることが確認できる.



図 8 NUMA を考慮した STREAM ベンチマーク

HPCS2014 2014/1/8

ULIBC を用いた NUMA を考慮した共有 4. メモリ型並列アルゴリズムの高速化

4.1 幅優先探索

近年、グラフを用いた解析手法は様々な分野で盛んに 研究されており、中でも幅優先探索 (Breath-first search; BFS) は基本的かつ重要なグラフ処理といえる. 並列 BFS アルゴリズムである Beamer's algorithm [8] は, 探索済の 点集合から未探索の隣接点集合を発見するという従来手法 である Top-down 探索と, 未探索の隣接点集合から探索済 みの点集合を発見する Bottom-up 探索を組み合わせ, 点 数 2²⁸, 枝数 2³² からなる Kronecker graph に対し 4-way Intel Xeon E7-8870 上で高い性能 5.1 GTEPS を達成して いる. ここで GTEPS は 1 秒あたり 10⁹ 枝を探索する性 能を表す. さらに著者らは Scatter アフィニティを基本と した NUMA を考慮した高速化を行い,同等の計算機環境 4-way Sandybridge-EP 上で 2.2 倍となる 11.15 GTEPS を達成した [9]. なお本論文では、文献 [9] では取り扱わな かったアフィニティ毎での性能,特性について議論する.

図 9に、SCALE=26、edgefactor=16 とした点数 $n = 2^{26}$ 、 枝数 $m = 2^{30}$ からなる Kronecker graph に対する BFS 計算性能 (GTEPS) をまとめたものである. 我々のアルゴ リズムは実行時間の大部分を各 NUMA ノードに分割した 領域へのローカルなメモリアクセスで計算できるように 局所性の改善を行っており、使用するスレッド数の増加に 従う性能向上が得られる.また、16スレッド時のアフィニ ティ設定毎の NUMA ノード数は Compact : Compact+ : Scatter = 1:2:4 であるが、その際の性能差は Compact : Compact+: Scatter = 1.00: 1.47: 1.87 となることか ら,同一のスレッド数であれば NUMA ノード数が多い割 当の方が高い性能となることが確認される.このように、 STREAM ベンチマークの Triad 演算に似た性能特性を確 認できるものの、参照や書出し対象のデータ量が大きくか つ複雑アクセスパターンとなるため、ファーストタッチで のメモリ確保が有効とならず、アフィニティ設定を行わな い w/o が極端に低速となる.

4.2 閉路を持たないパスの数え上げ

Knuth の Simpath アルゴリズム [12] は, 与えられたグラ フG = (V, E)上の2点 $s, t \in V$ 間の閉路を持たないパス を効率的に列挙する. ここで, 枝集合 E は {e₁, e₂, ..., e_m}, 点集合 V は $\{s, v_1, v_2, ..., v_\ell, t\}$ と与えられ, 各枝の添字は 始点 s からの幅優先探索で決定するものとする. 図 10 (a) のようなグラフ上で s-t 間の閉路を持たないパスの列挙を 行う際, Simpath アルゴリズムは各枝 ei をパスに含めるか 否かを表現した図 10 (c) の DD (Decision Diagram; 決定 グラフ)を構築する. ここで, DD 内の各節点 ei が持つ実線 矢印と破線矢印はそれぞれ 1-arc と 0-arc と呼び, ei をパス





に含めるか否かを表す. 各 e_i の選択が G における正しい パスを構成する場合は DD の端点が 1 (1-terminal), そうで なければ 0 (0-terminal) となるため, 閉路を持たないパス の総数は端点が 1 となる組合せを数えれば良い. 図 10 (a) の例では, 端点が 1 となるのは $\{e_1, e_4\}$ と $\{e_2, e_3, e_4\}$ とな り, 閉路を持たないパスの総数は 2 となる.



図 10 ZDD と Simpath アルゴリズム

Simpath アルゴリズムは始点 *s* からトップダウンに幅優 先探索で DD を構築していく際に, (1) 等価な節点を共有 する, (2) 節点を含めると条件を満たさない場合その節点 を削除して 0-terminal に接続する, といった ZDD (zerosuppressed binary decision diagram) [11] の簡約化規則を 用いる. 最終的には図 10 (d) のような規約な ZDD に変換 することができる. Simpath は効率的ではあるものの逐次 アルゴリズムであるため, 現在のマルチコアプロセッサの 性能を引き出すことは難しい.

すでに著者らは Simpath に対する 3 種類並列アルゴリ ズム Node-based, Range-based, Edge-based を提案し, 中 でもロックフリーアルゴリズムの Edge-based approach (EBA) は高速かつ並列効率が高く, 32 スレッド並列実行 時の EBA は 逐次アルゴリズムの 7 倍の性能となる [13]. EBA におけるスレッド数を k とした並列計算では. ス レッド j は各レベル $i \mod k$ ($0 \le i \le m$) に対応した FIFO キュー N_i の情報を保持する. スレッド j は FIFO キュー N_i から節点 n_i をデキューし, 枝 e_i を選択するか 否かに対応する新たな節点 $n_i^0 \ge n_i^1$ を生成し, それぞれ $l = (i+1) \mod k$ となるスレッド l に移動する. なお, ス レッド l における $n_i^0 \ge n_i^1$ の重複確認は, スレッドローカ ルなハッシュテーブルを参照しロックフリーで処理される. 続いて, 本研究により ULIBC を用いた NUMA アーキテ クチャを考慮した高速化について説明を行う. EBA におけ るスレッド j は自身が保持している FIFO キュー $N_i \mod k$ へのデータ参照と, スレッド l が保持している FIFO キュー $N_{(i+1) \mod k}$ へのデータ書込みが頻繁に行われる (図 11). そのため隣接のスレッドの距離が短くなるように, 論理コ アへ割り当てることが重要となる.

HPCS2014

2014/1/8



図 11 Edge-based parallel Simpath のメモリアクセスパターン

図 12 は, 16×16 の格子グラフ (点数 289, 枝数 544) 上 で EBA を用いて経路数の列挙に要した時間を各アフィ ニティ設定 w/o (アフィニティ設定なし), Scatter, Compact, Compact+毎にまとめる. なお, 計算機環境は 4-way Sandybridge-EP である. まず w/o は並列数の増加に従い 性能向上しており、アフィニティ設定を行わずとも並列効 果が高いことが確認できる.一方で,1スレッド時に w/o は 1478 秒要したが, Scatter, Compact, Compact+ はそれ ぞれ 752 秒となり、アフィニティ設定による安定性の向上 が確認できる. Scatter は,1スレッド時と比べて 2 スレッ ド時では隣接スレッドの距離の増大による一時的に性能低 下するが, さらなる並列数の増加に従い性能向上し, 最も 高速な 64 スレッド時には w/o よりも高い性能となる. 一 方, Compact と Compact+ は 1, 2, 4, 8 と並列数の増加 に従い性能向上するものの、Compact では 32 スレッド時 に, Compact+ では 16 スレッド時に, 一時的に性能低下す る. このとき、それまで割り当てが1ソケット内に限定し ていたのに対して,2ソケットに横断した割当となってい る. 16 スレッド時に Compact と Compact+ の性能差は ほぼ倍となることから、各物理コア上での2つ論理コアの 同時実行 (Hyper-Threading 機構) よりも CPU ソケット の横断によるメモリアクセスの不均一化によるオーバヘッ ドが大きいことが確認できる. さらに 32 スレッド時での Compact と Compact+ はほぼ同等であるため、物理コア を 16 ずつ使用した 2 ソケットの性能と, 論理コアを 32 使 用した1ソケットの性能が同等であることも推測できる. Compact, Compact+の 64 スレッド時には w/o (85.8 秒) と比べて, 49.2 秒と 1.74 倍高速化に成功した.

続いて図 13 に, 格子グラフの大きさを 10 × 10 から 18 × 18 と変化させた際の, アフィニティ設定毎の 64 ス レッド時の実行時間をまとめる. いずれのサイズに対して も Compact は w/o よりも同等かより高い性能を示して いる.



図 12 16×16-格子グラフに対するアフィニティ設定毎の実行時間



図 13 格子グラフに対するアフィニティ設定毎の EBA の実行時間

4.3 SDPA

半正定値計画問題 (Semidefinite programming; SDP) は 組合せ最適化,システムと制御,データ科学,金融工学,量 子化学など非常に幅広い応用を持ち,現在最適化の研究 分野で最も注目されている最適化問題の一つとなってい る.また今後のエネルギー供給計画(スマートグリッド 等)では非線形の複雑な最適化問題を扱う必要があり,こ れらの問題に対して強力な緩和値を算出できる SDP の高 速計算技術の確立が急務とされている.SDP に対しては 高速かつ安定した反復解法である内点法アルゴリズムが存 在しているが,巨大な線形方程式系の計算が大きなボトル ネックとなっている.具体的には線形方程式系の行列要素 の計算(以下 ELEMENTS)と行列の Cholesky 分解(以下 CHOLESKY)である.著者らのグループでは内点法アル

ゴリズムを記述したソフトウェアの開発・評価・公開を15 年以上行っており、疎性の追求、計算量やデータ移動量な どによる計算方法の自動選択などの技術を他に先駆けて実 現し、大規模な並列計算等によって上記のボトルネックの 高速化と世界最大規模の SDP を高速に解くことに成功し ている [14]. 特に CHOLESKY に関しては東京工業大学の スーパーコンピュータ TSUBAME 2.0 上において, 多数 GPU の活用や計算と通信のオーバーラップ技術を応用す ることによって、制約式の数が148万以上となる世界最大 規模の巨大 SDP を解き SDP の世界記録の更新及び最大 で 533 TFlops (Cholesky 分解: 4080 個の NVIDIA Tesla M2050)の性能を達成した [15]. その後, TSUBAME 2.5 において制約式の数が233万以上となる巨大 SDP を解き、 更なる SDP の世界記録の更新及び最大で 1.713 PFLOPS (Cholesky 分解: 4080 個の NVIDIA Tesla K20X) の性能 を達成した [16]. 一方, ELEMENTS に関しては入力問題 の疎性の活用や行列要素の並列計算によって大幅な高速化 に成功しているが、メモリバンド幅に律速されているため、 浮動小数点演算律速の CHOLESKY のように GPU 等に よる加速は容易ではない. 文献 [16] においては分散メモリ 環境上でのアフィニティ設定とメモリインターリーブ設定 による ELEMENTS を提案しているが、今回は著者らが開 発した 1 ノード版の SDP 用のソフトウェア SDPA [14] を 用いて、図 14 に示す ELEMENTS 計算の性能向上に関す る検証実験を行う.

<pre>set_affinity_policy(policy)</pre>
$B = \emptyset$
for $l = 1, 2, \cdots, h$ do
for $j \in \{j \mid \boldsymbol{F}_j^l \neq \emptyset\}$ parallel(thread) do
$\mathbf{set_affinity}(\mathtt{omp_get_thread_num}())$
for $i \in \{i \mid \boldsymbol{F}_i^l eq \emptyset\}$ do
Selects ${\mathcal F}_1, {\mathcal F}_2$ or ${\mathcal F}_3$ and compute B_{ij}^l
$B_{ij} = B_{ij} + B_{ij}^l$
end
\mathbf{end}
end
$clear_affinity()$

図 14 線形方程式系の行列要素 **B**_{ij} の計算 (ELEMENTS)

図 15, 16 に, Intel Xeon E7-4870 が 4 基搭載された共 有メモリ型計算機環境 4-way Westmere-EX 上の 2 種類の SDP 問題 mater-6, tensor4441 の実行時間を示す. mater-6, tensor4441 ともに疎性を有し, ELEMENTS がボトル ネックとなっている. まず tensor4441 はスレッド数の増加 に従い一定の性能向上が得られているが, w/o や Scatter では 2, 4, 8 スレッド時が同等でそれ以降は性能低下して しまうのに対して, Compact では 8 スレッドまで性能向上 を確認することができる. 一方, mater-6 はスレッド数の増 加に従い実行時間が増加してしまい,最も性能が高いのは スレッド数を1とした逐次実行時である.また,スレッド数 の増加による性能低下の影響を受けにくいのは Compact であることも確認できる.この問題に関しては今後とも注 意深く観察する必要がある.以上のように,計算機実験に おいて特性異なるどちらの問題に対しても,ソケットを横 断するデータアクセスがオーバヘッドになっていることが 確認された.



図 15 SDP 問題 tensor4441 に対するアフィニティ毎の実行時間



図 16 SDP 問題 mater-6 に対するアフィニティ毎の実行時間

5. **まとめ**

本研究では CPU アフィニティ設定やローカルメモリ上 に固定したメモリ領域確保を行うためのライブラリ ULIBC を開発し, 性質の異なる 4 種類のアプリケーションに対す る高速化を行い, いずれに対しても性能改善を達成した. ULIBC は 3 種類のアフィニティ設定 Scatter, Compact, Compact+ に対応し, スレッド並列領域においては計算機 トポロジを参照することで, 各スレッドが配置している論 理コアが計算機上のどの箇所に配置されているか容易に把 握することができる. これらの情報は NUMA アーキテク チャに対する高速化を行う上で非常に重要で、他のライブ ラリでは提供されていない機能である. 例えば, 我々が文 献 [9] で行ったように NUMA を意識しメモリアクセスの 局所性を高めたアルゴリズムを設計することができた場 合には、既存のライブラリや OS が提供する関数群を用い ることで実現することは難しくない. しかしながら論文中 でも行っているように Simpath や SDPA の ELEMENTS 演算などのより複雑なメモリアクセスを行うアルゴリズム に対する高速化においては、アルゴリズム上の性質上、どの アフィニティ設定が有効なのか、またその際にどれほどの ボトルネックが存在しているかなどの解析が必要となる. そういった解析にはアフィニティ設定やメモリ領域確保の 挙動の容易に把握できる環境が必要である. 今後も引き続 き ULIBC をベースとした高速なアルゴリズムの開発を行 う予定である.

謝辞 本研究の一部は,科学技術振興機構 CREST「ポ ストペタスケール高性能計算に資するシステムソフトウェ ア技術の創出」,ならびに科学技術振興機構 ERATO 湊離 散構造処理系プロジェクトの助成を受けたものである.

参考文献

- V. Agarwal, F. Petrini, D. Pasetto, and D. A. Bader: Scalable graph exploration on multicore processors, Proc. ACM/IEEE Int. Conf. High Performance Computing, Networking, Storage and Analysis (SC10), IEEE Computer Society, 2010.
- [2] F. Broquedis, J. Clet-Ortega, S. Moreaud, N. Furmento, B. Goglin, G. Mercier, S. Thibault, and R. Namyst: hwloc: a Generic Framework for Managing Hardware Affinities in HPC Applications. In Proceedings of the 18th Euromicro International Conference on Parallel, Distributed and Network-Based Processing (PDP2010), 2010.
- [3] J. Treibig, G. Hager, G. Wellein: LIKWID: A lightweight performance-oriented tool suite for x86 multicore environments. PSTI2010, 2010.
- [4] Intel(R) C++ Compiler XE 13.1 User and Reference Guide, Intel Thread Affinity Interface.
- [5] L. Huang, H. Jin, L. Yi, and B. Chapman: Enabling locality-aware computations in OpenMP, Journal Scientific Programming - Exploring Languages for Expressing Medium to Massive On-Chip Parallelism archive Vol. 18, Issue 3-4, pp. 169–181, 2010.
- [6] Y. Yasui, K. Fujisawa, K. Goto, N. Kamiyama, and M. Takamatsu: NETAL: High-performance implementation of network analysis library considering computer memory hierarchy, J. Oper. Res. Soc. Japan, vol. 54, no. 4, pp. 259–280, 2011.
- [7] M. Frasca, K. Madduri, and P. Raghavan: NUMA-Aware Graph Mining Techniques for Performance and Energy Efficiency, Proc. ACM/IEEE Int. Conf. High Performance Computing, Networking, Storage and Analysis (SC12), IEEE Computer Society, 2012.
- [8] S. Beamer, K. Asanović, and D. A. Patterson: Directionoptimizing breadth-first search, Proc. ACM/IEEE Int. Conf. High Performance Computing, Networking, Stor-

age and Analysis (SC12), IEEE Computer Society, 2012.

- [9] Y. Yasui, K. Fujisawa, and K. Goto: NUMA-optimized Parallel Breadth-first Search on Multicore Single-node System, The proceedings of the IEEE BigData 2013, Santa Clara, USA, IEEE Computer Society, 2013.
- [10] S. Kuo: Intel(R) 64 Architecture Processor Topology Enumeration, Intel Corporation White paper, 2012.
- [11] S. Minato: Zero-suppressed BDDs for set manipulation in combinatorial problems, DAC '93: Proc. 30th Int. Conf. Design automation, 1993.
- D. E. Knuth: The Art of Computer Programming, Vol. 4, Fascicle 1, Addison-Wesley, pp.121–123, 2009.
- [13] S. Takeuchi, J. Kawahara, A. Kishimoto, and S. Minato: Shared-Memory Parallel Frontier-Based Search, Algorithms and Computation Lecture Notes in Computer Science Volume 7748 (WALCOM 2013), 2013.
- [14] K. Fujisawa, K. Nakata, M. Yamashita, M. Fukuda: SDPA Project: Solving large-scale semidefinite programs. J. Oper. Res. Soc. Japan, 50, 278–298, (2007)
- [15] K. Fujisawa, T. Endo, H. Sato, M. Yamashita, S. Matsuoka and M. Nakata: High-Performance General Solver for Extremely Large-scale Semidefinite Programming Problems, Proc. ACM/IEEE Int. Conf. High Performance Computing, Networking, Storage and Analysis (SC12), IEEE Computer Society, 2012.
- [16] K. Fujisawa, T. Endo, Y. Yasui, H. Sato, N. Matsuzawa, S. Matsuoka, and H. Waki: Peta-scale General Solver for Semidefinite Programming Problems with over Two Million Constraints, Proc. ACM/IEEE Int. Conf. Parallel & Distributed Processing Symposium (IPDPS 2014), IEEE Computer Society, 2014.