グラフ処理における 同時実行制御の効率化 Efficient concurrency control in graph processing

伊藤 竜一♡ 藤森 俊匡♡ 鬼塚 真◇

Ryuichi ITO Toshimasa FUJIMORI Makoto ONIZUKA

近年,ネットワークやストレージ,センサなどの進歩により,大 規模なグラフ構造のデータが一般的になり、それを如何に効率的に 処理するかが技術課題となっている.特に,グラフ処理の基盤とな るエンジンの高速化は重要な要素である.しかしながら,並列度が 増加すると,同時実行制御がボトルネックとなり,実行効率が悪化 することが知られている.そこで,本稿では,並列環境として単一 マシンメニーコア環境に焦点を当て、同時実行制御に Compareand-Swap 命令や Hardware Transactional Memory を用 いて最適化を図ることで並列化のオーバーヘッドを減らし,高速な 処理を実現した.ベースとして用いたグラフ処理エンジンである GraphChi と比較して,同時実行制御に Compare-and-Swap 命令を利用した制御や Hardware Transactional Memory を用いた手法に置き換えることで,処理全体にかかる時間を最大 約19%,同時実行制御にかかる時間を最大約76%短縮すること に成功した. また, Hardware Transactional Memory の利 用最適化を行うことで,更に最大約8%の高速化に成功した。

In recent years, fast graph processing becomes more important as the size of graph data increases. It is a well-known issue that the concurrency control becomes a performance bottleneck in graph processing, in particular on massively parallel environments. In this paper, we propose an efficient design for graph processing on single PC with manycores. We introduce a notion of "virtual node" that achieves multi-stage aggregation based on Gather-Apply-Scatter computation model. We also accelerate the processing by employing a lock-free technique for concurrency control, which is suitable for workloads with low conflict rate. The result of experiments shows that our approach reduced both the elapsed time to 19% and the time required for concurrency control to 76% at most compared with the approach that uses Mutex.

1 はじめに

身の回りの様々な事象がデータとして蓄積されることが一般的 になってきたことで,結果として得られるソーシャルグラフや Web

♡ 学生会員 大阪大学情報科学研究科博士前期課程 <u>{ito.ryuichi, fujimori.toshimasa</u>}@ist.osaka-u.ac.jp グラフといったグラフ構造のデータが大規模化している. 例えば, 2016年12月時点のFacebookの月間アクティブユーザ数は18 億6000万人にも及び[1],その中でのユーザ同士の繋がりやポス トとそれに対するコメントなどが,大規模且つ複雑なネットワー クを形成している.このような莫大なグラフデータ源から,有益な 情報を如何に迅速に取り出すかが技術課題となっている.これま ではプロセッサの動作周波数向上により処理速度を上げることが 可能であったが,その成長が鈍化しつつあり,増加するコアを活か した並列処理による高速化の重要性が増している.

並列処理へのアプローチとして、単一マシン、限られたコア数、 共有メモリ環境であることを利用する方法がある.実行環境のセッ トアップや管理が簡単であり、データのやり取りにネットワークを 介することがないのでレイテンシが少ないといった利点がある一 方、CPU の数が比較的少数で計算能力に制限があり、巨大なデー タを一度に全てメモリ上にロードすることができないといった欠 点が存在する.そのため、限られた計算資源で、Memory-based な システムと比較すると非常に低速である Disk-based なシステム となる.単一マシン用グラフ処理エンジンの既存の研究として、 GraphChi [2] や TurboGraph [3], VENUS [4], GridGraph [5], NXGraph [6] が挙げられる.いずれの研究も、アクセスの局所性 を生かすデータ構造とアルゴリズムを設計することで、コアの利 用率および IO 操作量の削減して高速化を図っている.

グラフ処理高速化の課題の一つとして、更新処理時の一貫性を 保つためにコストがかかるという問題が挙げられる. 一般的に、現 実世界で見られるようなグラフデータの多くには、少数ながら、多 数の辺が集中する頂点が含まれることが知られている [7]. このよ うな頂点は、隣接頂点の多さから並列処理する際に更新操作の資 源として参照される回数が非常に多くなるため、一貫性を保つた めのコストが高く、ボトルネックになりやすい、このような頂点に 関する計算タイミングを分散させることで,競合を解決するコス トを抑えることができるが、 グラフデータの局所性を十分に活か すことができなくなる.一方,計算をシリアル実行にすることで局 所性を最大限に活かすことができるが、並列環境を活かすことが できなくなる、というトレードオフが発生する.本稿の提案手法で は同時実行制御に Lock-free である Compare-and-Swap 操作や Hardware Transactional Memory [8] を適用することで並列環 境と局所性の両方を維持し、グラフ処理の高速化を図った. 並列計 算モデルとして,同時実行制御を必要とするがデータの競合発生率 は低い、という Lock-free な手法で用いられる投機的実行と相性の 良い特徴を持つ Gather-Apply-Scatter Computation Model [9] を利用する. 更に, Gather-Apply-Scatter Computation Model を構成する演算の一つである Sum 演算に関する部分が可換モノ イドであるという性質に注目し、Sum 演算を多段階に分解するこ とで投機的実行の失敗を抑制する最適化を行った.

本稿では Gather-Apply-Scatter Computation Model を利用 できる単ーマシングラフ処理エンジンとして, GraphChi をベー スに提案手法の実装し, 評価実験を行った. 同時実行制御に従来の Mutex を用いた手法を Compare-and-Swap を用いた手法に置き 換えることで処理全体にかかる時間を最大約 19%, 同時実行制御 にかかる時間を最大約 76%, Hardware Transactional Memory を用いた手法に置き換えることで, 処理全体にかかる時間を最大 約 10%, 同時実行制御にかかる時間を最大約 53% 短縮すること に成功した. また, Hardware Transactional Memory の利用最 適化を行うことで処理全体にかかる時間を更に最大約 8% 短縮す ることに成功した.

本稿の構成は以下の通りである.2章 で本稿の前提となる知識 について述べる.3章 で提案手法を詳説し,4章 で提案手法の実 験とその評価・分析を行う.5章 で関連研究を紹介し,6章 で本稿 全体を俯瞰しまとめる.

[◇] 正会員 大阪大学情報科学研究科教授 onizuka@ist.osaka-u.ac.jp

2 前提知識

本章では、グラフ処理における同時実行制御で必要となる前提 知識と、提案手法のベースに用いたグラフ処理エンジンと計算モ デルについて述べる. 2.1 節 では同時並列的処理におけるデータ の一貫性を保つためのモデルについて、2.2 節 では 2.1 節 を実現 するために用いられる同期プリミティブについて、2.3 節 ではグ ラフ処理で用いられる計算モデルについて、2.4 節 ではベースに 用いたグラフ処理エンジンである GraphChi について述べる.

2.1 Data Consistency Model

従来はシングルマシンーシングルコアーシングルスレッドでの計 算が主流であったが,近年,コア単体の動作周波数上昇による性能 向上が必要とされる処理の大規模化に追いつかず,マルチマシンー マルチコアーマルチスレッドでの計算へと移行してきている.大 規模グラフ処理のように,一度に巨大な共有資源を扱う処理をマル チマシンーマルチコアーマルチスレッドで,データの一貫性を保 ちながら計算するためには,同時実行制御を行う必要がある.そ のような処理では,一度の更新処理に必要となる共有資源は共有 資源全体からするとごく一部であることが多いため,同時実行制 御の対象範囲を局所化することで効率化を図ることができる.こ の同時実行制御の粒度や方法を定めたものを Data Consistency Model と呼ぶ.

グラフ処理で利用される同期実行で制御を行うための具体的 な Data Consistency Model として, Bulk Synchronous Parallel(BSP) [9] がある. これは, 処理に Superstep という同期単位 を定め, Superstep ごとに設けられた Barrier のタイミングで全 ワーカーのデータを一斉に同期させるというモデルである. 値を 参照する場合は前回の Superstep で得られた値を利用すること で,他のワーカーによるデータの変更の影響を受けず,計算精度と 一貫性を保つことができる. しかしながら, Barrier のタイミング はすべてのワーカーが処理を終えたときと決められているため, 各 ワーカーへ分散されたタスクに偏りが生じていると処理時間の長 いワーカーに全体の処理時間が左右されてしまう問題がある.

一方,分散グラフ処理エンジンである GraphLab などで用いられ ている非同期実行型の Data Consistency Model として, Vertex Consistency / Edge Consistency / Full Consistency [10] があ る. Vertex Consistency は注目頂点の値のみ, Edge Consistency は注目頂点と隣接する辺のみ, Full Consistency は注目頂点と隣 接する辺と頂点のみに同時実行制御を行う.当然ながら,同時実行 制御対象外の辺や頂点の値は任意のタイミングで他のワーカーか ら書き換えられてしまう可能性があるため,計算精度が下がるが, 制御対象が少ないほど同時実行制御によるオーバーヘッドが抑え られ, また, 各ワーカーへのタスク分散に全体の実行時間が影響を 受けなくなり, 実行時間を短縮することができる.

2.2 同時実行制御と同期プリミティブ

Data Consistency Model を導入してデータの一貫性を保つため には、内部で共有資源へのアクセスを管理する同期プリミティブを 利用して同時実行制御をする必要がある.同時実行制御には多く のアプローチがあるが、Lock-based と Lock-free に大別される. それぞれ 2.2.1 節 と 2.2.2 節で述べる.

2.2.1 Lock-based

Lock-based な手法として,同期プリミティブに Lock 命令を 使ったものが知られている.特定の共有資源を利用する前にロッ クを取得することで,他のワーカーが同時に同じ共有資源にアク セスしてしまうことを防ぐ排他的な方法である.コストがかかり やすいという問題点があるが,仕組みが単純なため用いられるこ とが多い.

2.2.2 Lock-free

Lock-based な手法とは異なり、ロックを用いない、つまり、他の ワーカーに影響を与えない Lock-free な手法が知られている.こ の手法は同期プリミティブに専用の機械語命令を用いることで、大 きく分けて2つのアプローチから利用されている.

1 つ目は特定ドメインに特化した Lock-free データ構造を作 るというアプローチである.同期プリミティブに Compare-and-Swap(CAS) という機械語命令を用いる.ロックを取ることなく値 の変更を検知することでデータの一貫性を保つことが可能だが,複 数の値に跨がる一貫性が保てないという問題がある.

2 つ目のアプローチとして, Transactional Memory(TM) [8] が挙げられる. TM とは, データベース分野におけるトランザク ションをメモリに応用したものである.一連の処理を行う際に、ト ランザクション内のすべての処理が完了するまでメモリの共有部 分に反映させず、トランザクション内のすべての処理が終わった 時点でトランザクション内で使われたデータとメモリの共有部分 のデータをチェックし,T) 一貫性が保てる場合はトランザクショ ンの結果をメモリに反映 (Commit) し、トランザクションが成功 したことを通知する; F) 一貫性が保てない場合はトランザクショ ンの結果を破棄 (Abort) し、トランザクションが失敗したことを 通知する.このように一連のデータの読み込み・書き込みを1つ のトランザクションとして扱うことで原子性・一貫性を,各トラ ンザクションを個別に管理することで独立性を保証する.¹ TM を 利用することで, 読み込み側責任で悲観的となる Lock-based な 手法とは違い,書き込み側責任で楽観的に処理できるため,トラ ンザクションの Abort が発生しない限りでは同時実行制御のス ループット向上を望むことができる.TM の実装として,Software Transactional Memory と Hardware Transactional Memory がある. Software Transactional Memory(STM) とは、TM をソ フトウェアレベルで実装したものである.同期プリミティブとし て前述の CAS 命令が用いられ,1 つ目に挙げたアプローチの一 種でもある. 一方, Hardware Transactional Memory(HTM) と は、TM をハードウェアレベルで実装したものである. 現時点では, Intel Haswell [11], IBM Blue Gene/Q [12] · POWER8 [13] · zEnterprise EC12 [14] などの一部のプロセッサに, それぞれ独 自の命令として実装されている.STM と比較して,*)トランザ クションの途中経過はすべて高速なキャッシュ上で管理される;*) ユーザランドにオブジェクトを生成しない;等,ハードウェア資源 を有効に利用できるためスループットが高い.

2.3 Gather-Apply-Scatter Computation Model

ユーザ定義のグラフ処理アルゴリズムを記述する方法は処理エン ジンにより様々であるが、前述の Data Consistency Model を効 率的に実現するため、アルゴリズムを合成可能な形に分解して定義 するインターフェースが利用されることが多い. その具体的な計算 モデルとして, PowerGraph [9] で提案され, 他のグラフ処理エン ジンでも類似モデルが使われることが多い Gather-Apply-Scatter Computation Model(GAS) が知られている. GAS では, 辺から 値を取得し、集約するための Gather と Sum, 頂点の値を更新す るための Apply, 頂点の値から隣接する辺の値を更新するための Scatter と呼ばれる 4 つの演算で 1 つのアルゴリズムを定義す る. GAS の処理フローの擬似コードを Algorithm 1 に示す. 1~5 行目の foreach ループで Gather と Sum, Scatter 演算を行い, 全ての辺に対して値の割り当てと取り出しを完了させてから、6~ 8 行目の foreach ループで一括して Apply を実行し,頂点デー タの更新をする. なお, これを非同期に計算する場合は, Sum 演算 のみ同時実行制御の対象とすれば全体のデータの一貫性を保つこ とができ、並列処理の恩恵を受けやすいという特徴を持つ.

¹揮発性メモリ上で行われるため,耐久性については保証されない

表 1: Algorithm で使用されている主な記号の定義

$G_i(V, E)$	iイテレーション目の頂点集合 V,
	辺集合 E からなるグラフ G
get_vertex(i)	頂点番号 i から頂点データを取得する関数
e.cur_val	BSP における現イテレーションの辺 e の値
e.prev_val	BSP における前イテレーションの辺 e の値
v.acc	頂点 v の Sum 演算におけるアキュムレータ

Algorithm 1 Gather-Apply-Scatter Computation Model

Input: $G_n(V, E)$

Output: $G_{n+1}(V, E)$

1: foreach $e \in E$ do

- 2: $v_{dest} \leftarrow get_vertex(e.out)$
- 3: $e.cur_val \leftarrow Scatter(e)$
- 4: lock
- 5: $v_{dest}.acc \leftarrow Sum(v_{dest}.acc, Gather(e.prev_val))$
- 6: unlock
- 7: **end for**
- 8: foreach $v \in V$ do
- 9: $v.val \leftarrow Apply(v.acc)$
- 10: **end for**

2.4 GraphChi

本稿では, 提案手法を GraphChi [2] に組み込むことで, メニー コア環境に対応した高速なグラフ処理エンジンを提案する. そこ で本節では、ベースとなるグラフ処理エンジンである GraphChi について概説する. グラフ処理を行うユーザプログラムのために, GAS に基づくインターフェースと,1つの更新関数を定義して 利用するインターフェースを提供している.前者の場合は,Data Consistency Model を選択することができ, Sum 演算は Mutex を使って同時実行制御が行われる.後者の場合は,基本的に非同期 で処理が行われるが, 競合が起こり得るデータに関しては全てシ リアル実行することで同時実行制御は行わない. ディスクアクセ スには Parallel Sliding Windows(PSW) と呼ばれる, 実行時 にランダムディスクアクセスを発生させないアルゴリズムを利用 してグラフデータの読み書きを行っている. PSW は,全てのグラ フデータを一度にメモリに載せることが出来ず, Disk-based とな ることを前提としている. グラフ G = (V,E) に対して, 1~3 から なるプリプロセッシングを行う.1) Vを P 個のインターバルに分 割する (後述する Shard がメモリに載るように調整する); 2) 各 インターバルに含まれる頂点を終点とする全ての辺をそれぞれの インターバルに **Shard** として関連付ける. つまり, 全ての辺は必 ず唯一の Shard に含まれる; 3) 各 Shard に含まれる辺を始点で ソートする; 実行時には, i~iv からなる, Shard に含まれる辺が ソートされていることを利用した高速なシーケンシャルアクセス で処理を行う.i) 注目インターバルに関連付けられた Shard を Memory-Shard としてメモリにロードする; ii) 注目インターバ ルに含まれる頂点を始点とする辺をディスク上の Shard からメモ リにロードする. 始点でソートされているため, インターバルごと に、各 Shard (Parallel) の対象となる辺集合 (Window) が前から 後ろへとスライドしていく (Sliding); iii) 注目インターバルに含 まれる各頂点に対し、ユーザ定義の更新処理を並列実行する; iv) i

~iii を残りのインターバルに対して繰り返す; このように実行す ることで,1インターバルにつき P 回のシーケンシャルアクセス, つまり,全体で P² 回のシーケンシャルアクセスで処理を完了する ことができる.Pは グラフサイズ/メモリサイズ 程度となるため, 現実的な値と言える.

3 提案手法

事前実験として, CAS 操作や HTM を含む, 同期プリミティ ブの特性評価を行った. 類似した先行ベンチマーク調査として, STAMP [15] [16] などが挙げられるが, 具体的なアプリケーショ ンを対象としたものであり, 競合率や計算負荷といった基本的な指 標に基づくベンチマークはなされていなかった. そこで, いくつか の指標に基づく同期プリミティブに関するベンチマークを作成し た. 結果として, 図 1 から分かるように, ほぼすべての条件におい て CAS 操作が高速であることが, また, 並列度が高い / データ競 合率が低い / 同時アクセス要素数がある程度少ない / トランザク ションサイズが小さい, といった条件下で HTM の性能が最大限 発揮されることが確認された. ただし, CAS 命令を使った操作は, 直接的には複数値に対する同時実行制御ができないという問題が ある.

本稿の提案手法では、GAS に基づく単一マシン用グラフ処理エ ンジンの同時実行制御のコストを下げることでグラフ処理の高速 化を実現する. GAS は, 構成する 4 つの演算のうち, Sum の実行 時のみ共有資源のデータの一貫性を考慮することで全体の一貫性 を保つことができる計算モデルである. GAS による更新処理のう ち, 並列セクションにクリティカルセクションが占める時間は多く のアルゴリズムで 20~30%, 意図的に負荷を偏重させても 75% 程 度と事前の計測により分かっている.多数のデータからごく一部の データを非同期に取り出して扱っていることに注意すると,Sum は同時実行制御を必要とするが, データの競合が起こる確率は非常 に低い演算ということが分かる. つまり, Lock-based な悲観的同 期実行制御を用いると, lock/unlock によるオーバーヘッドが大き い. 一方, Lock-free は楽観的同時実行制御であるため, データの競 合が起こる確率が低い場合においては低コストで実行可能である. そこで,提案手法では,Sum 演算で既定の単一アキュムレータ³ のみにアクセスする場合は CAS 操作を, それ以外のオブジェクト にもアクセスする場合は CAS では制御出来ないため, HTM を適 用する. Algorithm 2 は同時実行制御に CAS 操作を, Algorithm 3 は HTM を適用した GAS による更新処理の擬似コードを示す. 前者は 5~12 行目で CAS 操作を利用することで共有資源である v_{dest}.acc の一貫性を保つ. 後者は 5~7 行を HTM のトランザク ションとすることで,他ワーカーからトランザクション内すべての 共有資源の一貫性を保つ. どちらも Data Consistency Model に BSP を利用している.

3.1 Hardware Transactional Memory 適 用の最適化

ソーシャルグラフのような現実世界によるグラフデータでは、そ のデータが大規模化すると隣接辺が極端に多い頂点が発生し得る. このような頂点が含まれると、GASの更新処理が集中し、データ の競合が非常に発生しやすくなる.特に、HTMにおいては、デー タの競合が発生した時点でトランザクションがAbortされるため、 急激に同時実行制御のコストが増加する.そこで、本節では、隣接 辺が多い頂点が含まれる場合でも、仮想ノードという実際の1つ の頂点に対して複数の中間結果を保持する概念を導入することで、 共有資源に対するアクセスを分散させてデータの競合発生率を抑 える手法を提案する.

³集約用の変数. 例として PageRank であれば, 各隣接頂点から流入す る PageRank 値を集約するために用いられる.

²CAS の同時アクセス数が 2 以上の場合と No Contrl は正確な同時実 行制御が行われておらず参考値である



(a) 同期プリミティブと並列度の関係 (b) 同期プリミティブと計算負荷の関係



(d) 同期プリミティブとアクセス対象 要素サイズの関係

般論文

ンあたりの同時アクセス要素数の関係 図 1: 同時実行制御ベンチマーク²

(e) 同期プリミティブと1クリティカルセクショ (f) 同期プリミティブと並列セクション中に占め

るクリティカルセクションの割合の関係

GAS 処理のうち, Sum 演算は, 注目頂点の値を更新する際に必 要となる隣接辺データを集約するための演算である.隣接辺は複 数存在する可能性があり、Sum 演算による集約演算は隣接辺の数 だけ呼び出される.この際,集約結果が処理対象となる辺の順序に 依存しないようにするため, 定義1に示すように, Sum 演算を二 項演算とする隣接辺データ集合が可換モノイドであることが必要 条件である.

- 完美 1	Sum 演算 に お	1+2	る可換エノイド	
尼我 I	Bulli (肉弁にの	11		
隣接	接辺データ集合		S	
	二項演算 Sum		$\odot: S \times S \to S$	
	$\forall s_1 \forall s_2 \forall s_3 \in S$;	$(s_1 \odot s_2) \odot s_3 = s_1 \odot (s_2 \odot s_3)$	
	$\exists z \forall s \in S$;	$z \odot s = s \odot z = s$	
	$\forall s_1 \forall s_2 \in S$;	$s_1 \odot s_2 = s_2 \odot s_1$	

そこで、組 (S, Sum) が可換モノイドであることを利用して、 Sum 演算による集約を多段階に分割する. 一旦隣接辺データを仮 想ノードに集約し、その後、頂点単位で対応する仮想ノードを再度 集約することで最終的な値を算出する.このように仮想ノードを 挟んで部分集約をすることによって,辺が集中する頂点の演算が 発生してもデータの競合発生率が上がらず、トランザクションの Abort 率が抑制されることで HTM の特性を活かした処理を高速 化することができる.

GraphChi への実装 3.2

2.4節で述べた通り, GraphChi はユーザプログラムのために2つ のインターフェースで提供しているが、本稿での提案手法は GAS を対象とした手法であるため、GAS に基づくインターフェースに 対して実装を行った.なお、提案手法では、単一マシンでスレッド プールを利用した並列処理を行うため、ロードバランスを保つこと

は容易なので Data Consistency Model には BSP に基づいたモ デルを採用した. 隣接辺データを集約する Sum 演算は,メモリか ら Shard をロードする際に, 辺単位で呼び出される. そこで, トラ ンザクションの範囲を狭めるため,共有資源である演算の中間結果 に Sum 演算でアクセスする部分のみに CAS 操作 と HTM を適 用した.また,各頂点に,スレッドごとに利用する中間結果を分け て持たせることで仮想ノードを実現した.これは、仮想ノードがス レッドローカルになっているため、同時実行制御を必要としないた め高速な実行が可能である.一方,仮想ノードの個数が多くなり, 中間結果を保持するためのメモリの使用量とのトレードオフが生 じる.そのため,各頂点ごとに,流入辺数が閾値以上であり,データ の競合が発生しやすい頂点のみに仮想ノードを実装し, それ未満 の頂点に関してはデータの競合率は低いと見なして仮想ノードを 用いずに HTM での同時実行制御を行う実装とした.

実験・評価 4

提案手法の性能を評価するため, GraphChi における既存の Mutex による実装をベースラインとして検証を行った. グラフ データやアルゴリズムの特徴による差異を確認するため,表2に 示すグラフデータで、複数のグラフ処理アルゴリズムを用いて実 験を行った.

本実験で用いたグラフデータを概説する. Pokec [17] [18] は、SNS である Pokec 内でのユーザの繋がりを能わすグラフ データである. Power law exponent⁴ が高いという特徴を持つ. Flickr [19] [20] は Flickr という画像投稿コミュニティ内でのユー ザの繋がりを表すグラフデータである. Power law exponent が 低いという特徴を持つ. LiveJournal [21] [22] は, ブログサービ スである LiveJournal でのユーザの繋がりを表すグラフデータで ある. Wikipedia, English [23] [24] は 英語版 Wikipedia 記事間 のハイパーリンクを表すグラフデータである.

⁴グラフデータが従う冪乗則の指数.数値が高いと辺が集中する頂点の 割合が下がる.

表 2: データセット							
Name	Nodes	Edges	Power law exp.	Shard 数			
Pokec [17] [18]	1,632,803	30,622,564	3.0810	2			
Flickr [19] [20]	2,302,925	33,140,017	1.7110	2			
LiveJournal [21] [22]	4,847,571	68,993,773	2.6510	3			
Wikipedia, English(dbpedia) [23] [24]	18,268,992	136,537,566	2.3710	5			

Algorithm	2	GAS	with	CAS	operation
-----------	---	-----	------	-----	-----------

Input: $G_n(V, E)$

Output: $G_{n+1}(V, E)$

1: parallel foreach $e \in E$ do

 $v_{dest} \leftarrow get_vertex(e.out)$ 2:

 $e.cur_val \leftarrow Scatter(e)$ 3:

- $fetched \leftarrow Gather(e.prev_val)$ 4:
- while true do 5:
- $loaded \leftarrow v_{dest}.acc.load()$ 6:
- intermediate $\leftarrow Sum(v_{dest}.acc, fetched)$ 7:
- **if** *acc*_{*dest*}.*acc*.*compare_and_swap*(8:
- loaded, intermediate) **then** 9:
- break 10:
- end if 11:
- end while 12:
- 13: end for
- 14: parallel foreach $v \in V$ do
- $v.val \leftarrow Apply(v.acc)$ 15:
- 16: end for

本実験で利用したグラフ処理アルゴリズムを概説する. PageRank は有向グラフの各頂点の重要度を決定するためのアルゴリズ ムである. 更新時の計算量が少なく, データの更新回数が多いとい う計算的特徴を持つ. 単一始点最短経路問題 (SSSP) は, 更新時の 計算量が少ないという計算的特徴を持つ.予め開始頂点を指定し, 収束するまで計算を行う.加えて,更新時のクリティカルセクショ ン内計算量を意図的に増やした検証用アルゴリズム (Heavy) を利 用した.

本実験で利用した環境について述べる. CPU は TSX をサポー トしている Intel Xeon E7-8890v3, 動作周波数は 2.5 GHz で物 理 18 コアである.メモリは 2TB, 補助記憶装置として書き込み速 度が約1.3 GB/sec, 読み込み速度が約200 MB/sec である HDD を利用した. ソフトウェア環境として, OS に Red Hat Enterprise Linux7(kernel-3.10.0), コンパイラに g++5.3.0 を利用した. コン パイル時最適化として,-O3 オプションを用いている. 並列度は物理 コア数とした.また,グラフデータの環境別 Shard 数は GraphChi の機能により、表2のように定められた.なお、同様の実験をコン シューマ向けの性能を持つデスクトップパソコンでも行ったが,同 様な結果が得られたため説明は省略する. CPU は TSX をサポー トしている Intel Core i5-6600, 動作周波数は 3.3GHz で 物理 4 コアである.メモリは 16GB, 補助記憶装置として書き込み速度が 約 350 MB/sec, 読み込み速度が 約 210 MB/sec である HDD を利 用した. ソフトウェア環境として, OS に Fedora23(kernel-4.2.3), コンパイラに g++5.3.1 を利用した.

Algorithm 3 GAS with HTM

Input: $G_n(V, E)$

Output: $G_{n+1}(V, E)$

1: parallel foreach $e \in E$ do

- 2: $v_{dest} \leftarrow get_vertex(e.out)$
- $e.cur_val \leftarrow Scatter(e)$ 3:
- $fetched \leftarrow Gather(e.prev_val)$ 4.
- transaction start 5:
- $v_{dest}.acc \leftarrow Sum(v_{dest}.acc, fetched)$ 6:
- transaction end 7.
- 8: end for

```
9: parallel foreach v \in V do
```

- $v.val \leftarrow Apply(v.acc)$ 10:
- 11: end for

(A) 様々な環境における性能特性, (B) 並列度による性能変化, (C)HTM に仮想ノードを用いた最適化を行った手法についての実 験をそれぞれ行った.

(A) 同時実行制御方法の処理全体の時間の比較を 図2 に示す. また,既存手法の Mutex によるオーバーヘッドが処理全体の時間 に占める割合を表3に,提案手法の同時実行制御にかかる時間の Mutex との比率を表4に示す. 図2,表3,表4より, CAS 操作 を利用した場合は全て, HTM もほとんどの場合において既存の Mutex を使った実装に対して高速化され,処理全体の時間が最大 19%,同時実行制御にかかる時間が最大76%短縮されたことが分 かる. 一部の条件において, HTM を利用した場合の方が Mutex を利用した場合に比べて処理全体の時間が増加する結果となった. 表4と照らし合わせると、元々の Mutex による同時実行制御コス トが処理全体に対して小さい計算の場合は性能の向上があまり望 めず、場合によっては性能が劣化する傾向が見受けられた.また、 横軸が CPU 使用率を表すヒストグラムを 図3に示す. 図3(a)と 図 3(b) を比較すると, 提案手法が既存手法に比べて最も CPU 使 用率が高く理想的な状態を表す階級である Ideal の占める割合が 増加したことから,並列環境を活かすことができたと分かる.

(B) データセットを Wikipedia, アルゴリズムを PageRank に 固定した環境で、コア数を制限することで並列度を変化させた時の 性能比較を 図 4 に示す. 図 4 より, 常に HTM が Mutex に対し て優勢であることが分かった.その中での傾向として,並列度が上 がることで HTM の性能は確かに向上しているものの, HTM と Mutex の差異の減少が見受けられる. これは, HTM は並列度が高 くなるとデータの競合が発生してしまい, 投機実行に失敗する確 率が増加するためと考えられる.図1(a)の結果を踏まえると,並 列度を更に上げるとやがて実行性能が悪化する可能性があるため, 高並列環境においては 3.1 節 で提案したような最適化の必要性が 再認された.



(D) SSSP 図 2: (A) 同時実行制御方法のデータセットごとの性能比較

表 3: (A) Mutex にかかる時間が全体に占める割合 (%)

アルゴリズム	Wikipedia	Flickr	LiveJournal	Pokec
PageRank	14.948	19.839	17.039	19.521
SSSP	11.623	13.761	11.222	12.970
Heavy	9.710	12.186	11.153	12.252

表 4: (A) Mutex 比の同時実行制御にかかる時間 (%)

アルゴリズム	手法	Wikipedia	Flickr	LiveJournal	Pokec
PageRank	HTM	94.008	73.696	46.680	52.050
	CAS	65.779	59.451	58.294	43.142
SSSP	HTM	103.755	57.188	51.183	60.815
	CAS	81.217	39.549	64.143	54.082
Heavy	HTM	190.068	84.201	66.535	55.249
	CAS	97.568	60.363	94.153	75.727

(C) データセットを Wikipedia, アルゴリズムを PageRank に 固定した環境で,既存手法と実験(A)での提案手法に加え,いくつ かの閾値における仮想ノードを用いた手法の性能比較と仮想ノー ドの使用メモリ量を図5に示す. 横軸ラベルの仮想ノードに続く 数値が閾値を示す.図5の結果より,HTMにスレッドローカル仮 想ノードを利用した最適化を行うことで,処理全体の時間を更に最 大約8%短縮することに成功した.閾値を下げる,つまり,スレッ ドローカル仮想ノードによる実行の割合を増やすと,同時実行制御 のコストを抑えることができ、実行性能を上げることができるが、 多量のメモリを必要とする. PageRank や SSSP のような Sum 演算の集約対象が単一値であるアルゴリズムでは影響は少ないが, 多値を取るアルゴリズムの場合は必要メモリ量が非常に多くなり, 1 Shard あたりのサイズを減少させ、処理全体のパフォーマンス を下げる可能性がある.一方, 閾値を上げる, つまり, HTM での実 行の割合を増やすと、スレッドローカル仮想ノードによる高速化 分を初期化処理のオーバーヘッドが上回り実行性能が低下してし まう. 図5からこのトレードオフが見受けられ、アルゴリズムや データの特徴に合わせて適切な閾値を設定する必要があることが 分かった.

5 関連研究

般論文

5.1 単一マシン用グラフ処理エンジン

GraphChi [2] を参考に設計された単一マシン用グラフ処理エン ジンとして, TurboGraph [3] や VENUS [4], GridGraph [5],



図 3: (A) CPU 使用率ヒストグラムの一例 ⁵ 灰: Idle, 赤: Poor, 黄: Ok, 緑: Ideal



図 4: (B) 同時実行制御方法ごとの図 5: (C) HTM と仮想ノードを用いコア数と性能の関係た手法の性能とメモリ使用量

NXGraph [6] などが挙げられる.

TurboGraph はデータの保存単位を細分化し, Pin-and-Slide という方法でデータのロードを行うことで, IO と CPU のオーバー ラップ効率を上げている. Pin-and-Slide では,ページと呼ばれる 単位でメモリにロード (Pin) し,次のページのロードをしている間 にロード済みのページ内のデータに関する更新処理を行う. ロード が終わり次第更新処理に移り (Slide),別スレッドで次のページの ロードを行う. このように,処理を細分化することで,ディスク IO 待ち / 更新処理待ち を減らしている. また,ディスクとして SSD を使うことで,ディスク IO のレイテンシを減らしている.

VENUSは、多くのグラフ処理アルゴリズムにおいて、隣接頂点から伝搬させる値を計算する際に、伝搬値を一時的に対象頂点間の辺の値として保持する必要が無いことに着目している.対象頂点の更新処理に隣接頂点の値を直接利用することで、中間データをメモリ上に置く必要性やディスクに書き出す必要性をなくして処理時間の短縮を行っている.しかしながら、辺に中間データを保

 $^{^5 {\}rm Intel}$ VT une Amplifier 2016: https://software.intel.com/enus/intel-vt une-amplifier-xe

持させる必要性があるアルゴリズム,例えば2ホップ以上先の頂 点の値を利用するようなアルゴリズム,を実行することができな いという制限がある.辺データの読み出しと更新処理をオーバー ラップさせることで,IO と CPU を効率的に利用している.

GridGraphは、頂点データを1次元に、辺データを始点と終点 に基づく2次元に分割して管理することで、従来のVertex-centric や Edge-centric な手法と比べ、1イテレーション当たりの辺デー タのロード回数を1回限りに抑えている.また、グラフ処理のス ケジューリングをユーザ定義のアルゴリズムに合わせて選択でき るようにすることで、不必要な辺データのロードを防ぎ、高速化を 行っている.

NXGraph は、辺の始点と終点を考慮する Destination-Sorted Sub-Shard という手法を用いることで、限られたメモリ空間でグ ラフの局所性を利用し、更新処理時に発生する書き込みによるデー タの競合発生率を抑えている.また、グラフデータのサイズと利用 可能なメモリサイズに基いて更新処理に Single-Phase Updating と Double-Phase Updating, Mixed-Phase Updating の3 種類 を使い分けることでディスクアクセスを減らしている.

5.2 Hardware Transactional Memory の 応用

HTM は古くから提案されてきたが [8], 近年ようやく商用のプロ セッサにも実装され、普及し始めたばかりの若い技術である.し かしながら、同時実行制御は分散処理やデータベース等、広い分野 で重要な技術であるため、その有用性は注目を集めており、本稿で 提案した単一マシングラフ処理エンジンだけでなく、既に分散処 理におけるメッセージパッシング [25] やインメモリデータベー ス [26] などにも応用されている.これらの既に応用化が進んでい る分野の特徴として、大部分は並列的に行われるものの、同時実行 制御が真に必要となるような複数アクセスが発生する割合が低い 処理を扱っているということが挙げられる.データの競合が少な い限り、同時実行制御としてのオーバーヘッドが非常に少ない、と いう HTM の特性を活かすことが重要であると分かる.

5.3 MapReduce との類似点

本稿で提案した仮想ノードへの部分集約に類似する概念として,分 散処理モデルである MapReduce [27] における Combiner によ る部分集約が挙げられる. MapReduce は基本的なフローは,値を Mapper で処理してからネットワーク越しで Reducer に渡し,集 約するというものである.そこに Combiner を導入し, Reducer での集約の前に Mapper のノード上で部分的に集約をすることで 通信量を削減して高速化を図ることができる [28].提案手法では 同時実行制御を減らすため, MapReduce では通信量を減らすため にそれぞれ部分集約を利用しており,目的は異なるがどちらも分 割統治法の集約ストラテジを最適化することでボトルネックにな り得る処理を減らしている.

6 まとめ

本稿では、単一マシンメニーコア環境でのグラフ処理エンジンに おいて、同時実行制御方法を従来の Mutex から CAS 操作や HTM に置き換えることで高速化を行った.更に、GAS の特性を考慮し、 HTM とスレッドローカル仮想ノードを組み合わせる高速な処理 手法を提案した.また、提案手法を GraphChi に実装し、性能評価 を行った.同時実行制御に Mutex を用いた実装と比較して、CAS 操作に置き換えることで処理全体にかかる時間を最大約 19%、同 時実行制御にかかる時間を最大約 76%、HTM を用いた手法に置 き換えることで、処理全体にかかる時間を最大約 10%、同時実行 制御にかかる時間を最大約 53% 短縮することに成功した.また、 HTM とスレッドローカル仮想ノードを用いた最適化を行うこと により、処理全体にかかる時間を更に最大約 8% 短縮することに 成功した.

謝辞

本研究の実験環境を提供して頂いたヒューレット・パッカード 社に深謝申し上げます.

[文献]

- [1] Facebook. Facebook company info. http://newsroom. fb.com/company-info/.
- [2] Aapo Kyrola, Guy E Blelloch, and Carlos Guestrin. GraphChi: Large-scale graph computation on just a pc. In *Proceedings of OSDI*, Vol. 12, pp. 31–46, October 2012.
- [3] Wook-Shin Han, Sangyeon Lee, Kyungyeol Park, Jeong-Hoon Lee, Min-Soo Kim, Jinha Kim, and Hwanjo Yu. TurboGraph: a fast parallel graph engine handling billion-scale graphs in a single pc. In *Proceedings of SIGKDD*, pp. 77–85. ACM, August 2013.
- [4] Jiefeng Cheng, Qin Liu, Zhenguo Li, Wei Fan, John CS Lui, and Cheng He. VENUS: Vertex-centric streamlined graph computation on a single pc. *Proceedings of ICDE*, April 2015.
- [5] Xiaowei Zhu, Wentao Han, and Wenguang Chen. Grid-Graph: Large-scale graph processing on a single machine using 2-level hierarchical partitioning. In Proceedings of the USENIX ATC, pp. 375–386, July 2015.
- [6] Yuze Chi, Guohao Dai, Yu Wang, Guangyu Sun, Guoliang Li, and Huazhong Yang. NXgraph: An efficient graph processing system on a single machine. *Proceed*ings of ICDE, May 2016.
- [7] Jure Leskovec, Deepayan Chakrabarti, Jon Kleinberg, Christos Faloutsos, and Zoubin Ghahramani. Kronecker graphs: An approach to modeling networks. *Proceedings of WAPA*, Vol. 11, pp. 985–1042, September 2010.
- [8] Maurice Herlihy and J Eliot B Moss. Transactional memory: Architectural support for lock-free data structures, Vol. 21. ACM, 1993.
- [9] Joseph E. Gonzalez, Yucheng Low, Haijie Gu, Danny Bickson, and Carlos Guestrin. PowerGraph: Distributed graph-parallel computation on natural graphs. In *Proceedings of OSDI*, pp. 17–30, Hollywood, CA, October 2012. USENIX.
- [10] Yucheng Low, Joseph E Gonzalez, Aapo Kyrola, Danny Bickson, Carlos E Guestrin, and Joseph Hellerstein. GraphLab: A new framework for parallel machine learning. *Proceedings of UAI*, July 2014.
- [11] Per Hammarlund, Rajesh Kumar, Randy B Osborne, Ravi Rajwar, Ronak Singhal, Reynold D'Sa, Robert Chappell, Shiv Kaushik, Srinivas Chennupaty, Stephan Jourdan, et al. Haswell: The fourth-generation intel core processor. *Proceedings of IEEE*, No. 2, pp. 6–20, December 2014.
- [12] Ruud Haring, Martin Ohmacht, Thomas W Fox, Michael K Gschwind, David L Satterfield, Krishnan Sugavanam, Paul W Coteus, Philip Heidelberger, Matthias A Blumrich, Robert W Wisniewski, et al. The IBM Blue Gene/Q compute chip. *Proceedings of IEEE*, Vol. 32, No. 2, pp. 48–60, December 2012.
- [13] B Sinharoy, JA Van Norstrand, RJ Eickemeyer, HQ Le, J Leenstra, DQ Nguyen, B Konigsburg, K Ward, MD Brown, JE Moreira, et al. IBM POWER8 processor core microarchitecture. *IBM Journal of Research and Development*, Vol. 59, No. 1, pp. 2–1, January 2015.
- [14] C Kevin Shum, Fadi Busaba, and Christian Jacobi. IBM zEC12: The third-generation high-frequency mainframe microprocessor. *Proceedings of IEEE*, No. 2, pp. 38–47, June 2013.

- [15] Chi Cao Minh, JaeWoong Chung, Christos Kozyrakis, and Kunle Olukotun. STAMP: Stanford transactional applications for multi-processing. In *Proceedings of IISWC*, pp. 35–46. IEEE, September 2008.
- [16] Takuya Nakaike, Rei Odaira, Matthew Gaudet, Maged M Michael, and Hisanobu Tomari. Quantitative comparison of hardware transactional memory for Blue Gene/Q, zEnterprise EC12, Intel Core, and POWER8. In Proceedings of ISCA, pp. 144–157. ACM, June 2015.
- [17] Pokec network dataset KONECT, May 2015
- [18] Lubos Takac and Michal Zabovsky. Data analysis in public social networks. In International Scientific Conference and International Workshop Present DTI, May 2012.
- [19] Flickr network dataset KONECT, May 2015.
- [20] Alan Mislove, Hema Swetha Koppula, Krishna P. Gummadi, Peter Druschel, and Bobby Bhattacharjee. Growth of the Flickr social network. In *Proceedings of COSN*, pp. 25–30, March 2008.
- [21] LiveJournal network dataset KONECT, May 2015.
- [22] Jure Leskovec, Kevin J. Lang, Anirban Dasgupta, and Michael W. Mahoney. Statistical properties of community structure in large social and information networks. In *Proceedings of WWW*, pp. 695–704, April 2008.
- [23] Wikipedia, English network dataset KONECT, May 2015.
- [24] Sren Auer, Christian Bizer, Georgi Kobilarov, Jens Lehmann, Richard Cyganiak, and Zachary Ives. DBpedia: A nucleus for a web of open data. In *Proceedings* of *ISWC*, pp. 722–735, October 2008.
- [25] M. Besta and T. Hoefler. Accelerating Irregular Computations with Hardware Transactional Memory and Active Messages. In Proceedings of the 24th Symposium on High-Performance Parallel and Distributed Computing (HPDC'15), pp. 161–172. ACM, 06 2015.
- [26] Viktor Leis, Alfons Kemper, and Tobias Neumann. Exploiting hardware transactional memory in mainmemory databases. In *Proceedings of ICDE*, pp. 580– 591. IEEE, March 2014.
- [27] Jeffrey Dean and Sanjay Ghemawat. MapReduce: simplified data processing on large clusters. *Proceedings of* ACM, Vol. 51, No. 1, pp. 107–113, March 2008.
- [28] Jimmy Lin and Chris Dyer. Data-intensive text processing with mapreduce. Synthesis Lectures on Human Language Technologies, Vol. 3, No. 1, pp. 1–177, 2010.

伊藤 竜一 Ryuichi ITO

大阪大学情報科学研究科博士前期課程在学中. 2016 年 大阪大学 工学部電子情報工学科卒業.マルチコア環境でのグラフ処理エン ジンの高速化に関する研究に従事.

藤森 俊匡 Toshimasa FUJIMORI

大阪大学情報科学研究科博士前期課程在学中. 2015 年 大阪大学 工学部電子情報工学科卒業.分散環境でのグラフ処理エンジンの 高速化に関する研究に従事.

鬼塚 真 Makoto ONIZUKA

大阪大学大学院情報科学研究科教授. 1991 年 東京工業大学情報工 学科卒.同年,NTT 入社. 2000-2001 年ワシントン大学客員研究 員,2013-2014 年電気通信大学客員教授,2012-2014 年 NTT 特 別研究員などを経て現職に至る.博士(工学).2004 年 情報処理学 会山下記念賞.2008 年 日本データベース学会上林奨励賞.2013 年 電子情報通信学会論文賞,情報処理学会論文賞.2014 年 電子 情報通信学会 I-Scover チャレンジ最優秀賞,日本データベース学 会論文賞.2015 年 電子情報通信学会論文賞.2017 年 日本デー タベース学会論文賞.ACM,電子情報通信学会,情報処理学会各 会員.