

投機的コヒーレンス制御にともなう ネットワークトラフィックの評価

乗 貞 由 華[†] 鈴 木 圭 介[†] 多 田 野 陽 介[†]
古 川 文 人[†] 大 津 金 光[†]
横 田 隆 史[†] 馬 場 敬 信[†]

DSMシステムにおいて、リモートメモリへのアクセスレイテンシの問題と、ネットワークトラフィックの問題の解決は重要な課題である。我々はハードウェアによってデータの一意性を保つDSMシステムにおいて read および write アクセス両方のコヒーレンスオーバーヘッドの削減を目的とした投機的コヒーレンス制御機構 SCCM を提案してきた。本機構を適用することにより、投機的コヒーレンス処理が成功した場合には、ネットワークトラフィックの削減が期待できる。本論文ではその投機的コヒーレンス処理にともなうネットワークトラフィックに関して評価した。その結果、メッセージ受信回数が Barnes で 0.40% ほど減少し、Raytrace で 0.97% 増加することがわかった。

Evaluation of Network Traffic Caused by Speculative Coherence Control

YUKA NORISADA,[†] KEISUKE SUZUKI,[†] YOUSUKE TADANO,[†]
FUMIHITO FURUKAWA,[†] KANEMITSU OOTSU,[†]
TAKASHI YOKOTA[†] and TAKANOBU BABA[†]

It is important to reduce the overhead of coherence control caused by memory access and network traffic. In order to reduce the overhead caused by the memory access of both read and write, we proposed the speculative coherence control mechanism for Cache Coherent DSM systems which maintains the consistency of data by hardware. If the speculative coherence control transaction is successful, we predict that the speculative coherence control mechanism reduces network traffic. In this paper, we evaluate network traffic caused by speculative coherence transaction. The results reveal 0.40% reduction of messages in the Barnes program and 0.97% increment in Raytrace.

1. はじめに

分散共有メモリ (以下、DSM) システム¹⁾ は、相互結合網を介した通信処理で実現される。故に DSM システムではリモートメモリへのアクセス (以下、リモートアクセス) が発行されてから完了するまでのレイテンシ (以下、リモートレイテンシ) がローカルメモリのレイテンシに比べ長い²⁾³⁾ ため、問題となっている。また、リモートアクセスによってネットワークが混雑すると、それがボトルネックとなり性能向上を妨げることが問題となっている。

そこで、この問題の解決方法として、ハードウェアによる Cache Coherent DSM システム (以下、CC-DSM システム) における投機的コヒーレンス処理に関する研究が行われている。投機的コヒーレンス処理は、データの共有状態を動的に予測し、その予測情報をもとに

行われる。予測のための処理と投機的コヒーレンス処理はユーザに完全に透過な形でハードウェアが行う。したがって、ユーザに提供する開発環境への影響はない。

我々は CC-DSM システムにおいて read および write アクセスのリモートレイテンシを短縮することを目的とした新たな投機的コヒーレンス制御機構 (Speculative Coherence Control Mechanism (以下、SCCM))⁴⁾ を提案してきた。本機構では投機を行うことによって投機的コヒーレンス処理が成功した場合、リモートレイテンシが削減でき、ネットワークトラフィックを削減できることが期待される。本論文では、SPLASH-2⁵⁾ を用いて本機構のシミュレーションを行いネットワークトラフィックにどのような影響を与えるかを評価する。

2. 投機的コヒーレンス制御機構

本論文では、以下に示す DSM システムとキャッシュコヒーレンスプロトコルを前提として議論する。本章では、2.1 節でモデルとなる DSM システムについて述

[†] 宇都宮大学工学部情報工学科
Department of Information Science, Faculty of Engineering, Utsunomiya University

べ、2.2節でキャッシュコヒーレンスプロトコルについて述べる。

2.1 DSMシステムとノード内の構成

DSMシステムとノード内の構成を図1に示す。このDSMシステムは、複数のノードがネットワークで結合された物であり、2Dメッシュ構造である。そのネットワークは、全二重で東西南北にReplyとRequest用に1つずつリンクを持っているので1ノードは合わせて16本のリンクを持っている。リンク幅は1本あたり16ビットである。また、2ノード間のメッセージの到着順序を保証される。各ノードは、1台のプロセッサコア、キャッシュ、DSMコントローラ、ディレクトリおよびメモリ、ネットワークインターフェース(以下、NI)およびSCCMから構成される。また、キャッシュはライトバック型、ディレクトリはフルマップ形式である。

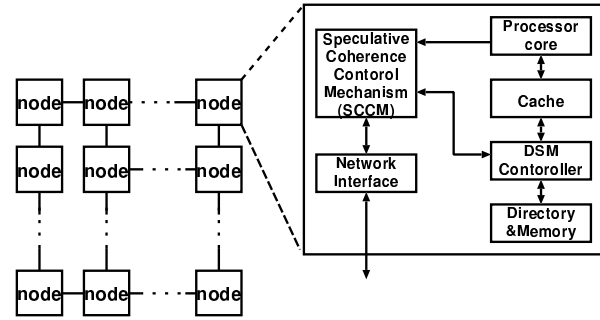


図1 DSMシステムとノード内の構成

2.2 キャッシュコヒーレンスプロトコル

キャッシュコヒーレンスプロトコルは無効化型である。また、キャッシュラインとメモリブロックは1対1に対応しており、それらの状態遷移はそれぞれ図2と図3に示される。なお、図2と図3でのX/Yは、Xが要求、Yが遷移後の応答を意味する。

キャッシュラインの状態は、Modified、Exclusive、Shared、Invalid、(以下、それぞれM、E、Sc、I)の4状態をとる。Mは、メモリブロックのコピーを専有し、その内容がメモリと一致していないことを示す。Eは、メモリブロックのコピーを専有し、その内容がメモリと一致していることを示す。Scは、他のノードに同一内容のラインがあることを示す。Iは、そのラインの内容が無効であることを示す。

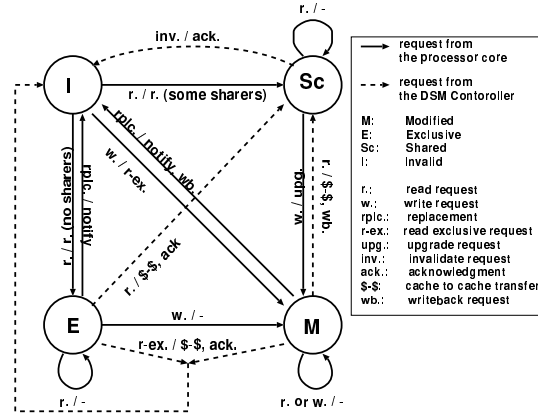


図2 キャッシュラインの状態遷移

メモリブロックの状態は、Uncached、Shared、Private、Busy-Shared、Busy-Private(以下、それぞれU、Sm、P、BS、BP)の5状態をとる。Uは、システム内にそのメモリブロックのコピーが無いことを示す。Smは、1個以上のノードのキャッシュがそのメモリブロックのコピーを保持し、その内容がメモリと一致していることを示す。Pは、1つのノードのキャッシュがそのメモリブロックのコピーを保持し、その内容がメモリと一致しない可能性があることを示す。BS、BPは、それぞれS、Pへ遷移するための処理を行っている最中であることを示す。BS、BPのメモリブロックへの要求があった場合、それを管理するDSMコントローラはその要求を全て保持し、その状態がそれぞれS、Pへ遷移した後で要求を受けた順に処理する。

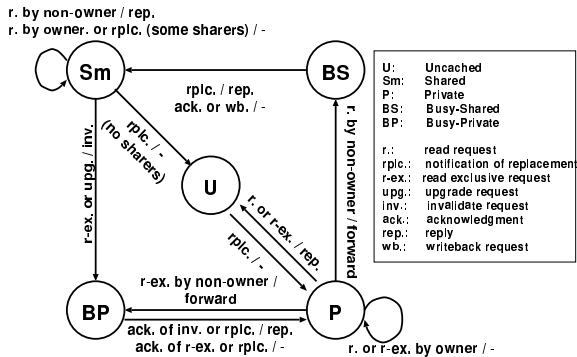


図3 メモリブロックの状態遷移

3. 投機的コヒーレンス処理とその適用

本章では、SCCMで行われる投機的コヒーレンス処理について述べる。

3.1 投機的コヒーレンス処理

SCCMは以下の5つの投機的コヒーレンス処理を行う。

(1) speculative self-downgrade

(以下、spec.self-dwg.)

Producer-Consumer共有データ(以下、PC共有デー

タ)へのreadアクセスによって発生するコヒーレンスオーバーヘッドを軽減するための処理である。spec.self-dwg.は、EまたはMの状態のコピーをすでに保持するノードが自分のキャッシュに対して発行し、そのコピーの状態をScに遷移させる。その後、コピーの状態がScに遷移したことをホームに通知するためにメッセージ(SPEC_SELF_DWG)を送信する。また保持しているコピーがMの状態の場合は、その内容をホームに書き戻す。

(2) speculative self-invalidation

(以下、spec.self-inv.)

PC共有データへのwriteアクセスおよびMigratory共

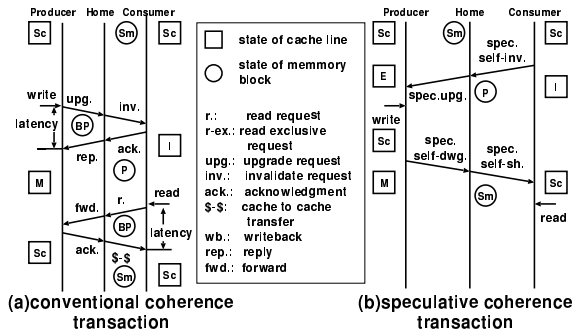


図4 PC共有状態における投機的コヒーレンス処理の例

有データ (M共有データ) への read と write 両方のメモリアクセスによって発生するコヒーレンスオーバーヘッドを軽減するための処理である。spec.self-inv. は、コピーをすでに保持するノードが自分のキャッシュに対して発行し、そのコピーを無効化する。その後コピーを無効化したことをホームに通知するためにメッセージ (SPEC_SELF_INV) を送信する。また、保持しているコピーの状態が M のときは、その内容をホームに書き戻す。

(3) speculative send block in shared state (以下、spec.send-sh.)

PC共有データへの read アクセスによって発生するコヒーレンスオーバーヘッドを軽減するための処理である。spec.send-sh. はホームにおいて Sm の状態のメモリブロックに対して read 要求を発行するノードを予測し、そのノードに対して Sc の状態のコピーを保持したメッセージ (REPLY_SPEC_SEND_SH) を送信する。それと同時に、予測されたノードが Sc の状態のコピーを保持していることを DSM コントローラに対して通知する必要がある。

(4) speculative send block in exclusive state (以下、spec.send-ex.)

M共有データへの read および write アクセスによって発生するコヒーレンスオーバーヘッドを軽減するための処理である。spec.send-ex. は、ホームにおいて U の状態のメモリブロックに対して read 要求に引き続く read-exclusive 要求を発行するノードを予測し、そのノードに対して E 状態のコピーを保持したメッセージ (REPLY_SPEC_SEND_EX) を送信する。それと同時に、予測されたノードがコピーを専有し、その内容がメモリブロックと一致しない可能性があることを DSM コントローラに対して通知する必要がある。

(5) speculative upgrade (以下、spec.upg.)

PC共有データおよび M共有データへの write アクセスによって発生するコヒーレンスオーバーヘッドを軽減するための処理である。spec.upg. は、ホームにおいて Sm の状態のメモリブロックに対して upgrade 要求を発行するノードを Sc の状態のコピーを持つノードの中から予測し、そのノードに対して E の状態への遷移を許可したことを通知するためにメッセージ (RE-

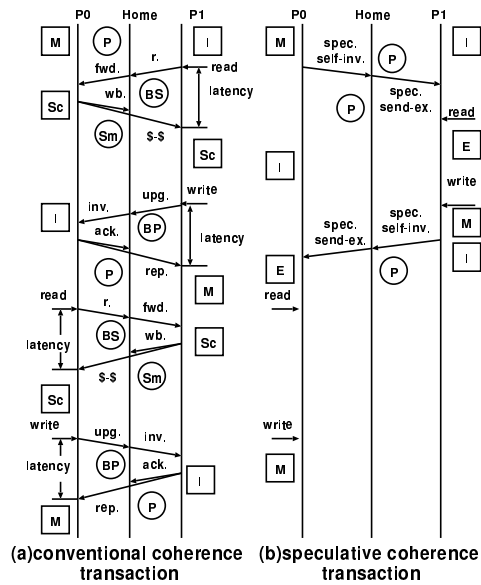


図5 M共有状態における投機的コヒーレンス処理の例

PLY SPEC_UPG) を送信する。それと同時に、予測されたノードがコピーを専有し、その内容がメモリブロックと一致しない可能性があることを DSM コントローラに対して通知する。

3.2 PC共有状態における処理

図4は、Producer-Consumer共有状態(以下、PC共有状態)におけるコヒーレンス処理の流れを示したものである。図4(a)はPC共有データへのメモリアクセスを投機を行わないコヒーレンス処理で解決する場合の例である。図4(b)SCCMによる投機的コヒーレンス処理によってコヒーレンスオーバーヘッドを完全に除去できたときの例である。

図4(a)では、ProducerのSc状態のデータへのwriteアクセスによって、Homeにupgrade要求(REQ_UPGRADE)が送信される。次にHomeからSc状態のデータのコピーを持っているConsumerへinvalidate要求(COHE_INVL)が送信される。Consumerはコピーを無効化し、その旨をHomeに通知するためのメッセージ(COHE_REPLY_INVL)を送信する。その後、Homeはupgrade要求を許可するためのメッセージ(REPLY_UPGRADE)をProducerに送信する。次にConsumerのM状態のデータへのreadアクセスによって、Homeにread要求(REQ_READ_SH)が送信される。次にHomeはProducerへ最新のデータを要求するためのメッセージ(COHE_COPYBACK)を送信する。ProducerはHomeへ最新のデータ(COHE_REPLY_COPYBACK)を書き戻す。HomeはConsumerに要求されたデータ(REPLY_SH)を送信する。

図4(b)では、spec.self-inv.によってSc状態のデータのコピーを持っているConsumerが予めコピーを無効化し、その後コピーを無効化したことをHomeに通知するためのメッセージ(SPEC_SELF_INV)を送信する。Homeはspec.upg.によってSm状態のメモリブ

ロックに対して upgrade 要求を発行するノードを Sc 状態のコピーを持つノードの中から予測し、そのノード (Producer) に対して E 状態への遷移を許可したことを通知するためのメッセージ (REPLY_SPEC_UPG) を送信する。その後 Producer からの write アクセスが発生しているため、リモートレイテンシが削減されている。次に、今度は E または M 状態のコピーをすでに持っている Producer が spec.self-dwg. によってコピーの状態を Sc に遷移させる。その後、コピーの状態が Sc に遷移したことを Home に通知するためのメッセージ (SPEC_SELF_DWG) を送信する。そして、Home は Consumer に予め Sc 状態のコピー (REPLY_SPEC_SEND_SH) を送信する。その後、Consumer からの read アクセスが発生しているため、リモートレイテンシが削減されている。

3.3 M 共有状態における処理

図5は、Migratory 共有状態 (以下、M 共有状態) におけるコヒーレンス処理の流れを示したものである。図5(a)はM共有データへのメモリアccessを投機を行わないコヒーレンス処理で解決する場合の例である。図5(b)はSCCMによる投機のコヒーレンス処理によってコヒーレンスオーバーヘッドを完全に除去できたときの例である。

図5(a)では、P1のM状態のデータへのreadアクセスによってP1からHomeへread要求 (REQ_READ) が送信される。次にHomeからP0へ最新のデータを要求するためのメッセージ (COHE_COPYBACK) が送信される。その後、P0はHomeに最新のデータを書き戻すためのメッセージ (COHE_COPYBACK) を送信し、同時にP0からP1に要求されたデータを持つメッセージ (REPLY_SH) が送信される。次に、P1からSc状態のデータへwriteアクセスが発生すると、P1からHomeへupgrade要求メッセージ (REQ_UPGRADE) が送信される。その後、HomeからSc状態のデータのコピーを持っているP0にinvalidate要求メッセージ (COHE_INVL) が送信される。P0はコピーを無効化し、その後、Homeへ無効化したことを通知するためのメッセージ (COHE_REPLY_INVL) を送信する。それと同時にP1へ要求されたデータを持つメッセージ (REPLY_EXCLDY) を送信する。以下、P0のM状態のデータへのreadアクセスが発生した場合も同様に処理される。

図5(b)では、P0が spec.self-inv. によって自分が保持している Sc 状態のデータのコピーを無効化する。そして、コピーを無効化したことを Home に通知するためのメッセージ (SPEC_SELF_INV) を送信する。Home は U 状態のメモリブロックに対して read-exclusive 要求を発行するノードを予測し、そのノード (P1) に対して E 状態のコピーを持つメッセージ (REPLY_SPEC_SELF_EX) を送信する。その後 Producer からの read および write アクセスが発生しているため、リモートレイテンシが削減されている。次に P1 が spec.self-inv. によって自分の保持している Sc 状態のデータのコピーを予め無効化し、その後コピーを無効化したことを Home に通知するためのメッセージ (SPEC_SELF_INV)

を送信する。次に Home は U 状態のメモリブロックに対して read-exclusive 要求を発行するノードを予測し、そのノード (P0) に対して E 状態のコピー (REPLY_SELF_SEND_EX) を送信する。その後 P0 からの read および write アクセスが発生しているため、リモートレイテンシが削減されている。

3.4 投機が失敗した場合の処理

投機が失敗した場合の処理を以下に示す。spec.self-inv. および spec.self-dwg. では、投機が失敗した場合、本来必要ない Request、Reply メッセージが流れる。

spec.send-sh.、spec.send-ex. および spec.upg. では、投機が失敗した場合、投機が失敗したことをディレクトリコントローラに通知するためのメッセージ (COHE_REPLY_SPEC_SEND_EX、COHE_REPLY_SPEC_UPG) を送信する。

4. 評価

本章では、投機を行わない従来のコヒーレンス制御機構 (以下、non-spec) を適用した DSM システム、read アクセスのみを投機の対象とする投機のコヒーレンス制御機構 (以下、spec-read) を適用した DSM システム、read と write アクセスの両方を投機の対象とする SCCM を適用した DSM システムの相互結合網のメッセージトラフィックがどのようになっているかを実験により明らかにする。

4.1 評価方法

評価は、表1に示す DSM システムのモデル上で SPLASH-2 の PC 共有データモデルのアプリケーションである Barnes と M 共有データモデルのアプリケーションである Raytrace を CC-NUMA 型並列計算機シミュレータ RSIM⁶⁾ でシミュレーションすることにより行った。RSIM は実験のために SCCM と Spec-read の機構をのせ、メッセージの情報を計測できるように改良したものを使用した。Barnes、Raytrace それぞれのシミュレーションで、表2に示すメッセージの受信回数、メッセージの受信回数×平均メッセージサイズ、平均レイテンシ、平均ブロック時間を計測した。以降はメッセージの表記は表中括弧内の略称を使用する。ここで、メッセージの受信回数×平均メッセージサイズとは、ネットワークを流れたメッセージの合計のサイズであり、平均レイテンシとは、送信されたメッセージが送信元のノードの NI を出てから送信先のノードの NI に到着するまでにかかったサイクル数の平均であり、平均ブロック時間とは、1つのメッセージが送信元から送信されて送信先に届くまでにネットワークの混雑のためにネットワークを使わず、使えるまで待機していたサイクル数の平均である。

4.2 評価結果

図6、図7に Barnes と Raytrace のそれぞれのメッセージ受信回数、図8、図9に受信回数×平均メッセージサイズ、図10、図11に平均レイテンシ、図12、図13に平均ブロック時間を示す。それぞれの図のグラフの並び方は、図6と同じ順序で並んでいる。

この結果から、Barnes においてメッセージ受信回

表1 シミュレーションパラメータ

ノード数		16
プロセッサコア		4-way スーパースカラ
命令キャッシュ		完全ヒット
データキャッシュ(ノンブロッキング)		
ラインサイズ		64バイト
最大未解決ミス数		4
サイズ	L1	64Kバイト
	L2	8Mバイト
アクセス レイテンシ	L1	1サイクル
	L2	タグ データ
		3サイクル
		4サイクル
メモリアクセスレイテンシ		
ローカルメモリ		80 サイクル
ネットワーク		
トポロジ		2D-メッシュ
メッセージヘッダサイズ		16バイト
フロー制御方法		ワームホール
レイテンシ /1hop	データなし (2hops以降)	128サイクル (+48サイクル/1hop)
	データ付き (2hops以降)	300サイクル (+48サイクル/1hop)

表2 評価メッセージ一覧

非投機メッセージ
REQ_READ_SH(R.R.S.)
REQ_READ_EX (R.R.E.)
REQ_UPGRADE(R.U)
REPLY SH(R.S.)
REPLY EXCL(R.E.)
REPLY_UPGRADE(R.U.)
REPLY_EXCLDY(R.ED.)
COHE_COPYBACK(C.C.)
COHE_COPYBACK_INV(C.C.I.)
COHE_INVL(C.I.)
COHE_REPLY_COPYBACK(C.R.C.)
COHE_REPLY_COPYBACK_INV(C.R.C.I.)
COHE_REPLY_INVL(C.R.I.)
COHE_REPLY_WRB(C.R.W.)
COHE_REPLY_REPL(C.R.R.)
投機メッセージ
SPEC_SELF_DWG(S.S.D)
SPEC_SELF_INV(S.S.I.)
REPLY_SPEC_SEND_SH(R.S.S.S)
REPLY_SPEC_SEND_EX(R.S.S.E)
REPLY_SPEC_UPG(R.S.U.)
COHE_REPLY_SPEC_SEND_EX(C.R.S.S.E)
COHE_REPLY_SPEC_SEND_UPG(C.R.S.S.U)

数ではSCCMでは全体で0.40%減少している。特にspec-read、SCCMともにR.S.S.Sが発生していることからR.R.S.とR.S.が減少している。またSCCMではS.S.I.が発生していることから、C.I.とC.R.I.が減少している。このことから投機によってメッセージの受信回数を減少させることがわかる。また、これにともないメッセージ受信回数×平均メッセージサイズもSCCMでは全体で1.13%減少している。特にR.R.S.とR.S.が減少している。しかし、平均レイテンシはSC-

CMでは全体で34%増加している。特にR.S、E.ED.、C.R.CとC.R.C.I.が増加しており、平均ブロック回数は、SCCMでは全体で1.54%増加している。特にR.S.、R.ED.、C.I.、C.R.C.、C.R.C.I.、C.R.I.が増加している。このことからネットワークトラフィックは削減されていないことがわかる。次にRaytraceにおいて、メッセージ受信回数ではSCCMでは全体で0.97%増加している。特にspec-read、SCCMともにR.R.S.、R.S.、R.E.、C.C、C.R.C.が増加している。これにともない、メッセージ受信回数×平均メッセージサイズは、SCCMでは全体では0.65%増加している。特にR.S.、R.E.が増加している。平均レイテンシは、SCCMでは全体で1.39%増加しているが、R.S.、R.E.、R.ED.、C.R.I.が減少しており、平均ブロック時間では、全体で7.35%減少しており。特にR.S.、R.E.、R.ED.、C.R.C.I.、C.R.I.が減少している。このことからネットワークトラフィックは削減されていることがわかる。

今回使用したBarnesとRaytraceでは投機メッセージの受信回数が少ないことから、投機をほとんど必要としなかったことが考えられる。そのためネットワークトラフィックがあまり削減されなかったのではないかと考える。

5. まとめ

本論文では、SPLASH-2を用いて投機的コヒーレンス制御機構のメッセージトラフィックについて評価した。その結果、メッセージ受信回数は最大で0.40%、メッセージ受信回数×平均メッセージサイズは最大で1.13%、平均レイテンシは最大で1.39%、平均ブロック時間は最大で7.35%減少することがわかった。

今後は、SPLASH-2の他のアプリケーションを用いて評価を行う予定である。

謝辞 本研究は、一部日本学術振興会科学研究費補助金(若手研究(B)14780186、及び基盤研究(B)14380135、(C)14580362)の援助による。

参考文献

- 1) Keleher, P. J.; Distributed Shared Memory Home Pages. <http://www.cs.umd.edu/keleher/dsm.html>.
- 2) 細見岳生、加納健、中村真章、広瀬哲也、中田登志之; 並列計算機Cenju-4の分散共有メモリ機構、並列処理シンポジウムJSP'99論文集、pp. 15-22(1999).
- 3) Landon, J. and Lenoski, D.; The SGI Origin: A cc-NUMA Highly Scalable Server, *Proceedings of the 24th Annual International Symposium on Computer Architecture* (1997).
- 4) 古川文人、多田野陽介、乗貞由華、大津金光、馬場敬信; DSMシステムにおける投機的コヒーレンス制御機構の提案と評価、情報処理学会研究報告2001-ARC-142 2001-HPC-85、pp. 169-

174(2001).

- 5) Woo, S.C., Ohara, M., Torrie, E., Singh, J. P. and Gupta, A.; The SPLASH-2 Program; Characterization and Methodological Considerations, *Proceedings of the 22th Annual International Symposium on Computer Architecture*, pp. 179-190(1998).
- 6) Pai, V. S., Ranganathan, P. and Adve, S. V.; RSIM; An Execution-Driven Simulator for ILP-Based Shared-Memory Multiprocessors and Uniprocessors, *Proceedings of the Third Workshop on Computer Architecture Education*(1997).

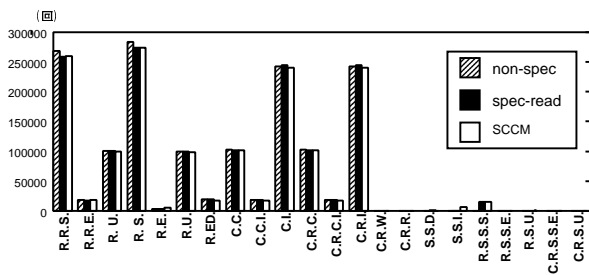


図6 Barnesのメッセージ受信回数

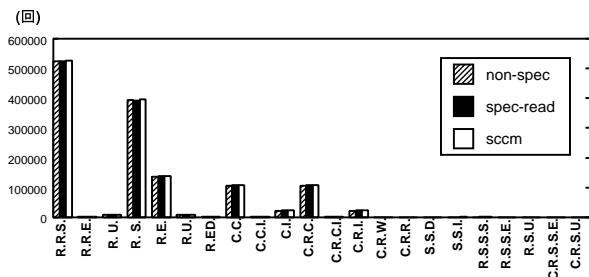


図7 Raytraceのメッセージ受信回数

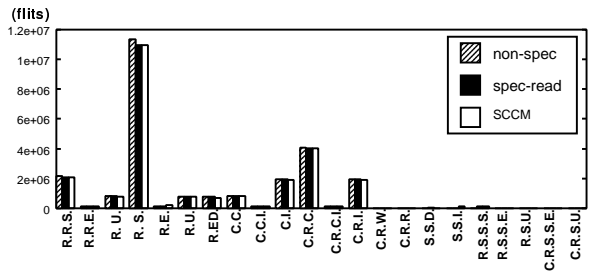


図8 Barnesのメッセージ受信回数×平均メッセージサイズ

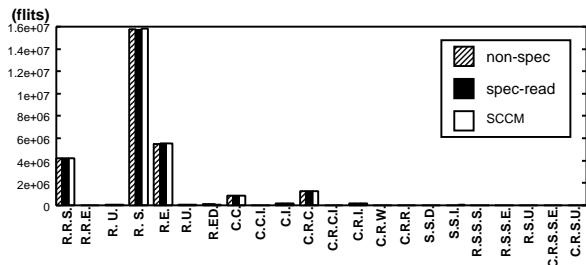


図9 Raytraceのメッセージ受信回数×平均メッセージサイズ

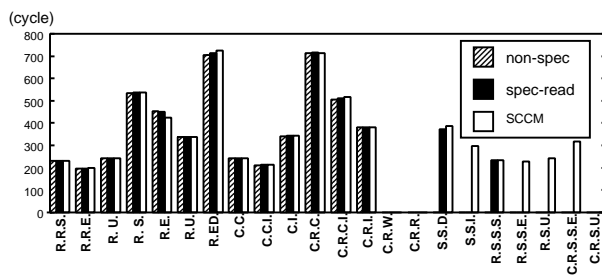


図10 Barnesの平均レイテンシ

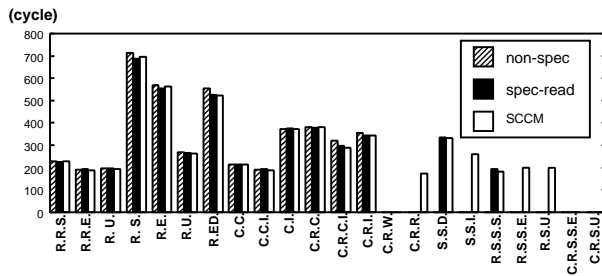


図11 Raytraceの平均レイテンシ

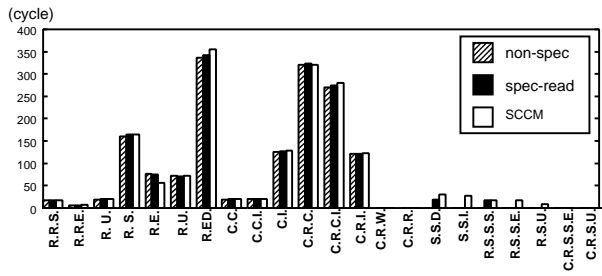


図12 Barnesの平均ブロック時間

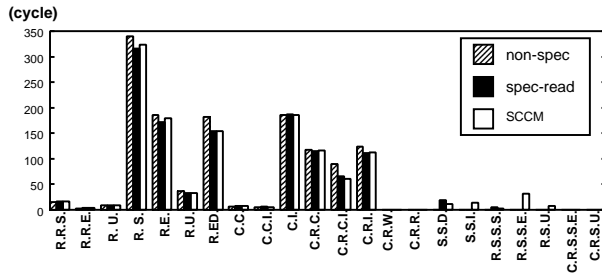


図13 Raytraceの平均ブロック時間