

# CAPRI: Prediction of Compaction-Adequacy for Handling Control-Divergence in GPGPU Architectures

著者： Minsoo Rhu, Mattan Erez

出典： *Proceedings of the 39th International Symposium on Computer Architecture (ISCA '12)* pp.61-71,2012.

発表者： 本多・近藤研究室 1253016 松本 洋平

## 1 研究背景

GPUは大量のSIMD型演算器を搭載する。SIMD型演算器は1命令で複数のデータに対して並列演算ができるため、高い演算能力を持つのに加え制御回路を小さくできる利点がある。しかし、制御系の命令を効率良く実行できないという問題がある。従来の制御命令の実行効率を改善する方法は性能低下を起こす可能性があった。本論文では、それを解決する手法を提案する。

## 2 GPUのスレッド実行とその問題点

### 2.1 warp 実行

GPUではwarp単位でメモリアクセスと演算が実行される。warpは単一命令で同じ処理を実行する32スレッドの集合である。warpを実行する時は8個または16個の演算器をグループにして演算をする。例えば16個の演算器が搭載されているGPUでwarpを実行する場合は、1warpを2回に分けて演算を実行する(図1)。なお、GPU上で96スレッドが生成された場合は32スレッドごとに3つのwarpに分割される(図2以降では紙面の都合上1warpを4スレッドで記載する)。

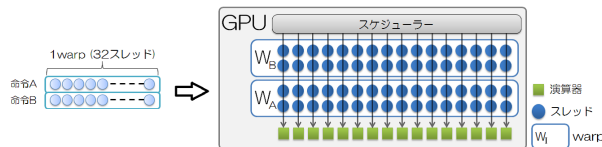


図1: GPUでのwarp実行

### 2.2 分岐実行における問題点

CPUにおいて分岐命令を実行するとき、各スレッドはIF文のTHEN節とELSE節の片方のみを選択して実行する。一方で、GPUでは分岐命令の実行時にwarp中の各スレッドはTHEN節とELSE節の双方にそれぞれ分岐してしまう場合があるがwarp単位で同じ命令を実行するためその場合は双方のパスを実行する必要がある。その時、THENを実行したスレッドはELSE節では休んでいることになる(図2)。これは非効率でありSIMD型アーキテクチャの構造上の欠点とされている。

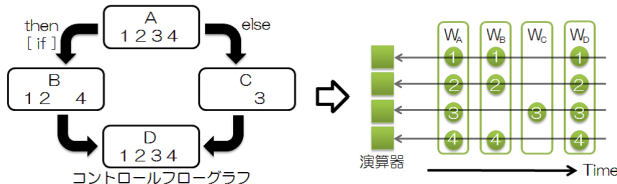


図2: GPUでの分岐命令実行順序

## 3 既存手法

提案手法のベースとなる2つの分岐実行の効率化手法について以下に述べる。

### 3.1 Warp-wide PDOM (WPDOM) [1]

分岐によってはwarp中の全スレッドがTHEN節とELSE節の片方のパスのみで実行されることがある。この場合は事前に片方のパスだけが実行されることがわかればもう片方のパスの命令実行を回避することができる。WPDOMはwarpごとにスタックを用意して分岐実行を管理する。このスタックでは分岐の際にそれぞれのパスで実行すべきwarp中のスレッドの番号を保持している。図3の例ではTHEN節である命令Bで全てのスレッド‘0123’が実行される。その後ELSE節である命令Cの各スレッドを実行する順番になるが実行するスレッドがないのでスタックは‘xxxx’となっている。この‘xxxx’の場合は命令を発行しなければ無駄な演算の実行を回避できる。

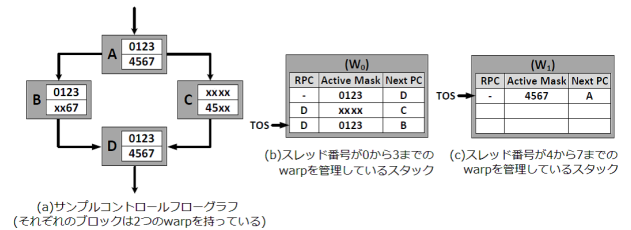


図3: WPDOMの実行例

### 3.2 Thread Block Compaction (TBC) [2]

TBCは同じパスを実行する別々のwarp内のスレッド同士をまとめて実行(コンパクション)する手法である。図4の命令Bでは‘0xxx’と‘xx67’の2つのwarpがある。これらをまとめて‘0x67’の1つのwarpにコンパクションして実行すれば無駄な演算が削減される。なお、TBCではコンパクションを行うために、複数ワーブの実行管理を図4のCooperating Thread Array (CTA) という1つのスタックで管理する。このため、複数ワーブの命令実行の進行が同期化されてしまい命令実行の並列性が失われて性能低下の原因となる場合がある。

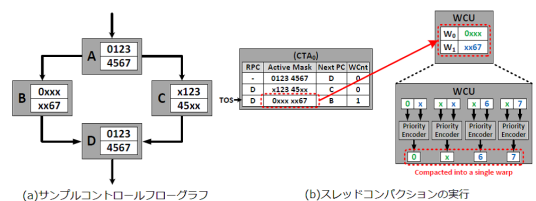


図4: TBCの実行例

## 4 提案手法：CAPRI

TBCは実行する全てのwarpに対してスレッドコンパクションを行う。この際、個々のスレッドに空きがなくコンパクションしてもwarp数が減らないあるいは分岐の片方のパスだけが実行される命令間など、スレッドコンパクションの効果がないwarp間でコンパクションを実行してしまうとメモリアクセスの不連続性や命令実行の進行が同期されwarp間の並列性が活用できずに性能低下を引き起こす可能性がある。例えば図5(b)の例ではTBCは使用しているCTAにより複数命令間の実行を同期化してしまっている。これにより命令A中のあるwarpで起きたキャッシュミスが後続の命令Bの実行を遅らせている。

上記のようなコンパクションは同じようなパターンで繰り返されることがしばしばある。そこでCAPRIではスレッドコンパクションの効果がない分岐命令を記録することで次回以降でその分岐命令を実行する際にその部分でのコンパクションをしないようにする。具体的にはコンパクション前の命令中のwarp数とコンパクション後のwarp数を比較して数が一致してしまう場合、すなわちコンパクションの効果がない場合は次回以降のコンパクションを実行しないようテーブルに登録する。このようにすることで例えば図5(c)のようにコンパクションの効果が得られない命令Aと命令Bのコンパクションをキャンセルすることでキャッシュミスのレイテンシを隠蔽することも可能になる。

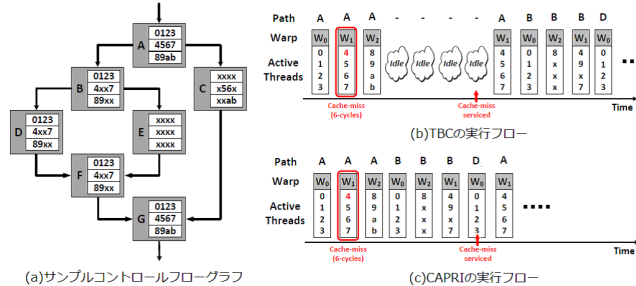


図 5: CAPRI の実行例

## 5 実験手法：GPGPU-sim[3]

GPGPU-simはTor Aamodtらによって開発されたGPUアーキテクチャを研究するためのツールである。

本実験ではNVIDIA's Quadro FX5800をモデル化したものにWPDOM, TBC, TCB+, CAPRI(1bS), CAPRI(1bL)を実装して各種法を評価した。なお、TCB+は片方のパスだけ実行されるヒントをコンパイル時に挿入する手法、CAPRI(1bS)は一度コンパクションが不要と判断されるとそれ以降のコンパクションをしない手法、CAPRI(1bL)は2bitカウンターを使用してコンパクション実行の可否を決定する手法である。

また、実行するベンチマークはParboil, Rodinia, CUDA SDK, GPGPU-simに付随のベンチマークの中から性能検証の精度が高いものを選択する。

## 6 実験結果

### 6.1 IPC

分岐のあるベンチマークではCAPRI(1bL)はWPDOMより12.6%の性能向上、TCB+より7.2%の性能向上が見られた。また、CAPRI(1bS)はCAPRI(1bL)ほどの性能向上は見られなかったがWPDOM, TCB+よりも高い性能が見られた。分岐のないベンチマークのではTBCはWPDOMに対して平均で10.1%、TCB+は平均で8.9%の性能低下が見られた。一方、CAPRIはWPDOMに対して±1%の性能が変化した。

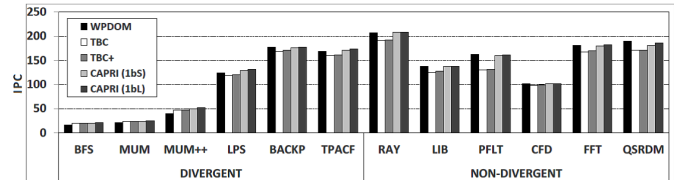


図 6: 各手法のIPC比較結果

### 6.2 実装オーバーヘッド

プロセッサ設計では性能と消費電力と回路面積の3つはトレードオフの関係にあり、性能向上は達成したがその分だけ消費電力と回路面積が増えるということが多くある。CAPRIの場合は最大で10.8%の性能向上を達成した。また、消費電力は演算実行数が減るため削減される。一方で8160bitのレジスタ量だけ回路面積増加(1mm<sup>2</sup>程度)してしまうが、これはチップ全体の1/300程度である。よって、ほんの僅かな回路面積のオーバーヘッドで性能向上と消費電力の削減を達成した。

## 7 結論

本論文ではGPU上で分岐命令を効率良く実行する提案した。分岐命令がある場合は、平均で7.6%性能向上が見られた。分岐命令がない場合は、±1%性能が変化した。また、ほんの僅かな回路面積で提案手法が実装できることを示した。

## 参考文献

- [1] W.W. Fung et al. Warp Formation and Scheduling for Efficient GPU Control Flow. MICRO-40 pp.407-420, 2007.
- [2] W. W. Fung and T. M. Aamodt. Thread Block Compaction for Efficient SIMT Control Flow. HPCA-17 pp25-36, 2011.
- [3] Ali Bakhoda et al. Analyzing CUDA workloads using a detailed GPU simulator. ISPASS 2009 pp163-174, 2009