GPUにおける細粒度パワーゲーティング向け スレッド発行制御手法の検討

松本 洋平 $^{1,a)}$ 近藤 正章 1 和田 康孝 1 本多 弘樹 1

概要:近年では GPU の消費電力の増加が問題となっている.本稿では GPU に搭載された SIMD 演算器 のリーク電力の削減手法として細粒度パワーゲーティングを適応することを考え,その際のリーク電力削 減効果を向上させるためのスレッド発行制御手法を提案する.前提とする細粒度パワーゲーティングでは GPU の SIMD の各演算器単位で電源の供給を制御する.提案手法では,電源の ON/OFF によるスリー プモードとアクティブモードの移行時に生じる電力的なオーバーヘッド抑えるために,各ワープ内のス レッドの発行制御を行うものである.スレッド発行制御手法としては一部の演算器に集約してスレッド実 行を行うスレッドコンパクション,1warp を2つの warp に分割する warp 分割を検討する.シミュレー ションによる初期評価の結果,スレッドコンパクション,および warp 分割を適用した場合のリークエネ ルギー削減率はそれぞれ 46%,71%になった.また両者を組み合わせた場合のリーク電力削減率は74%に なることがわかった.

1. はじめに

近年,GPUの演算処理能力が注目されている.モバイル デバイスやデスクトップ計算機での画像処理だけでなく, その高い演算処理能力を活用して,HPC アプリケーショ ンでの利用が広まっている.

GPU は大量の SIMD 型演算器を搭載することで,高性 能を達成する.SIMD 型の演算器は1命令で複数のデータ に対して演算処理を行えるため,制御回路を小さくでき る利点がある.そのため,演算あたりの電力効率は良いプ ロセッサである.一方で,非常に多数の演算器を搭載し ているため,合計の消費電力が大きいことが GPU の問題 点としてあげられる.文献 [1] によると NVIDIA GeForce GTX 280 の消費電力は平均で 172[W] であり,その内訳 はアイドル時電力が 83[W],演算にかかる動的消費電力が 37[W],メモリアクセスにかかる動的消費電力が 52[W] で ある.従って,アイドル時の消費電力が比較的大きいこと がわかる.

ー般的にアイドル時の消費電力には,処理要求を受け付けるための回路動作やクロックゲーティングができない ことで消費される動的消費電力,およびリーク電流による リーク消費電力が含まれる.近年は,半導体プロセスの微 細化に伴うリーク電流の増大が問題となっており, リーク 消費電力を削減することは GPU にも効果的であると考え られる.

LSI においてリーク消費電力を削減する手法としてパ ワーゲーティング (PG) がある.PG は LSI の回路ブロッ クに対して電源遮断用のスイッチを設け,動作の必要がな いときに電源供給を遮断することでリーク電力を削減する 手法である.従来の GPU には Streaming Multiprocessor (SM)単位での粗粒度の PG が実装されている [4].しかし, SM 単位の PG では電源 ON/OFF の際の時間的・電力的 なオーバーヘッドが大きく,また,SM 全体を使用しない 場合にしか PG を適用できない.

本稿では,SM に対して処理が割り当てられている場合 でも,1) 同じワープ内で分岐命令の分岐方向が異なる場 合には,SIMD 演算器内で演算処理を行う必要のない演算 器が存在する,2) メモリアクセス待ちによるストール発 生時には長期間演算器がアイドルになることがある点に着 目し,SM 単位ではなく,SM 内の SIMD 演算器の各演算 器単位で細粒度にPGを適用することを考える.これによ り,従来の SM 単位よりも PG を行う機会を多くすること ができ,より効果的にリーク消費電力を削減できると期待 できる.

ただし,電源の ON/OFF によるスリープモードとアク ティブモードの移行には電力的なオーバーヘッドが生じる ため,頻繁なモード遷移を行うと,かえって消費電力が増

¹ 電気通信大学 大学院情報システム学研究科 Graduate School of Information Systems, The University of Electro-Communications

 $^{^{}a)}$ matsumoto@hpc.is.uec.ac.jp

情報処理学会研究報告 IPSJ SIG Technical Report

大する可能性がある.そのため,時間的に細粒度に PG を 行うに際には,電力オーバーヘッドよりもリーク電力の削 減量が大きくなるように,1回電源遮断した際の PG サイ クルを確保する必要がある.そこで本稿では,SM 内の各 演算器がなるべく長期間 PG できるようにするためのス レッド発行制御手法を検討し,その初期評価を行う.

2. GPU アーキテクチャと細粒度パワーゲー ティング

本章では, GPU において細粒度 PG を行う上での技術 背景として, GPU アーキテクチャと細粒度 PG の概要に ついて述べる.

2.1 GPU アーキテクチャ

一般的な GPU のハードウェアモデルを図1に示す.
GPU は計算処理を行う Streaming Multiprocessor (SM)
と, SM に対してデータを供給する Interconnection Network, L2 Cache, GDDR DRAM による記憶階層で構成されている. SM は SIMD ユニット,制御回路,レジスタ, Shared memory, L1 Cache から構成されている. GPU では複数のスレッドが1つの warp という単位でまとめられ,発行制御や演算の実行が行われる.複数の warp は SM 内の CTA (Cooperative Thread Array) に格納され,その中から実行可能なものがラウンドロビン方式で選択され,SIMD ユニットで演算が実行される.ここで,warp 内の全スレッドは同じ命令が実行される.また,warp 中のそれぞれのスレッドが SIMD ユニット内のどの演算器で実行されるかはあらかじめ固定されている

GPU で分岐命令が実行され, warp 内のスレッドが異な る分岐方向の命令を実行する必要がる場合には, warp 内 の全スレッドが両方の分岐先コードを実行し,本来処理す べき命令以外はマスクされる

2.2 細粒度 PG 手法

パワーゲーティング (PG) は動作させる必要のない回路 ブロックへの電源供給を遮断することで当該回路のリー ク電力を削減する手法である.例えば,GPUの SIMD ユ ニットの各演算器に対して PG を適用する場合には図2の ように,それぞれの演算器とグラウンド線との間にスリー プトランジスタを挿入する.各スリープトランジスタはス リープシグナル (sleep signal) によって制御され,各演算 器の電源供給の ON/OFF が行われる.

細粒度に PG を実行する場合はモード切り替え時オー バーヘッドに注意する必要がある.オーバーヘッドにはス リープトランジスタの駆動やスリープ信号の伝搬, VGND に溜まった電荷の放電などのエネルギー的・時間的オー バーヘッドが含まれる.

図3にPGのモード切替時のタイミングチャートを示す.



図 1 GPU アーキテクチャモデル



図 2 SIMD 演算器における PG





時刻 T0 で回路ブロックで実行するべき処理 (work) がなく なり,同時にフラグ (activity) がアイドル状態であること を示している.ここで,時刻 T1 でスリープ信号をアサー トすることにより当該回路の電源を遮断する.この際,先 に述べたようにスリープトランジスタの制御などにより,



図 4 SIMD ユニット実行 (発行制御前)

動的電力が消費される.なお,ここではまだ仮想グラウン ド線 (virtual GND) にリーク電流が流れるため,T2 まで リーク電流の削減はできない.次にT3 で実行が必要な処 理が現れたため,スリープから復帰する.この時,同じく スリープトランジスタの制御のために動的電力が消費され る.また,Virtual GND に溜まった電荷を放電する必要が あるため,実行再開はすぐにできず,T4 まで遅延が発生 してしまう可能性がある.

PGにより T2 から T4 までがリーク電力を削減できる期間となるが,実際の消費エネルギー削減は,PG のモード切り替えに必要な動的消費電力分を差し引いて考えなければならない.このオーバーヘッドが吊り合う時間(T2 からT2'まで)をBreak Even Time(BET)と呼ぶ.もし,BETより短い期間でPGをしてしまうと,かえって消費電力の増大を招く.そのため,PG を行う場合にはBET を考慮する必要がある.なお,このBET は半導体製造プロセスや対象回路の温度・構造に依存して変化する.

3. スレッド発行制御手法

本章では GPU の SIMD 演算器に細粒度 PG を適用した 場合に,リーク電力削減効果を増大させるための2つのス レッド発行制御手法について述べる

3.1 PG 向けスレッドコンパクション

これまでにも, GPU上で SIMD ユニットの演算器使用 率を向上させるためのスレッド発行制御に関する研究は多 く行われている [5][7][8][9].本提案手法は,文献 [5] で提案 されているスレッドコンパクションを応用し,一部の演算 器のアイドル時間が長くなるようにスレッド発行制御を行 い,PG サイクルを増やすことで効率的なリーク電力削減 を狙う.

前述のように, ある warp 内で各スレッドの分岐方向が 異なる分岐命令があると, その warp は分岐先の両パスを 実行するため, 各演算器で演算が実行されないサイクル



図 5 SIMD ユニット実行 (発行制御後)

が多く発生する.この場合,使用されない演算器の電源供 給を遮断し,リーク電力を削減できる可能性がある.しか し,各 warp で使用されない演算器の位置がばらばらであ ると,各演算器の ON/OFF の切り替えが頻発し,効率的 な PG はできない.例えば,図4の例では,演算器が使用 されないサイクルは多くあるが,10 サイクル程度の BET を仮定すると,PG によりスリープモードに移行させるこ とでリーク電力削減効果がある演算器は2つのみである.

そこで,本提案手法では,warp中のスレッドをコンパ クションすることにより,長期間スリープにできる演算器 数を増やすことを狙う.スレッドコンパクションは,warp 中に使用されないスレッド実行レーンがある場合に,warp 内でスレッドの発行を一部の演算器に集約して行う手法で あり,もともとは空いた演算器で異なる warpのスレッド を実行することで,演算器の利用効率を高めるために提案 された手法である[5].本論文の提案手法では,空いた演算 器に対して PGを適用することで,当該演算器のスリープ サイクルを長期化することができると考えられる.図4の warp 実行の例に対して,スレッドコンパクションを適用 した場合の実行の様子を図5に示す.コンパクションによ り,リーク電力削減効果がある演算器は4つに増加してい る.このように,本提案手法により PG によるリーク電力 削減効果の増大が期待される.

スレッドコンパクションを実現するための演算器とレジ スタの構成を図6に示す.図6(a)は通常のGPUの構成で あり,演算器の数だけ個別にレジスタが実装されている. 各スレッドはIDによりあらかじめ使用するレジスタ・演 算器が定められているので,使用する演算器を動的に変更 することはできない.そこで本提案手法図6(b)のように 演算器の数に合わせたポート数を持つ共有レジスタを設け ることで各スレッドが使用する演算器がを柔軟に変更する ことができるようになる.なお,同様の構成は文献[5]で も提案されている. IPSJ SIG Technical Report



図 7 SIMD ユニット実行 (warp 分割後)

3.2 warp 分割

GPU での warp の実行は図4のスレッドが発行されてい ない空きサイクル(点線で囲まれているスロット)のように 処理すべき warp が存在せず,演算器が使用されないサイ クルが発生する.処理すべき warp が存在しない原因とし てはプログラム上のスレッド数が十分でない場合,キャッ シュミスによるストール,他の SM 上のスレッドとの同期 が発生した場合などがあげられる.

そこで,本提案手法では,演算器の使用率が低い場合に, 1warp を2つの warp に分割することで,使用する演算器 を半分にし,長期間スリープにできる演算器数を増やすこ とを狙う.図4の warp の実行例に対して warp 分割を行っ た際の実行の様子を図7に示す.warp 分割により,リーク 電力削減効果がある演算器は9つに増加している.このよ うに,本提案手法によりPGによるリーク消費エネルギー 削減効果を増大させることができると期待される.

一方で,演算器の使用率が高い場合に warp 分割をして しまうと,後続の warp 実行が遅れ,性能が低下してしま う可能性がある.性能低下を防ぎつつリーク電力を削減す るためには,実行時の演算器の使用率を予測し,それに応 じて warp 分割を行う必要がある.演算器の使用率の予測 には,コンパイル時の静的解析によるものと,実行時に演

表1評	価に用いた	GPU	の仮定
-----	-------	-----	-----

Streaming Multiprocessor (SM) 数	15	
スレッド数/SM	1536	
SIMD 幅	16	
SIMD ユニット数/SM	2	
SIMD パイプライン数	32	
Clock speed (core)	$1.4[\mathrm{GHz}]$	
Clock speed (memory)	$3.7[\mathrm{GHz}]$	
レジスタサイズ/SM	128[KB]	
シェアードメモリサイズ/SM	16[KB]	
L1 キャッシュサイズ/SM	48[KB]	
L2 キャッシュサイズ	768[KB]	
メモリーチャネル数	6	
メモリーコントローラ	FR-FCFS	
DRAM リクエストキューサイズ	32	

表 2 ベンチマーク

名前	内容	命令数
BFS	幅優先探索	17M
MUM	タンパク質・RNA・DNA 間の類似解析	77M
LPS	三次元ラプラス変換	82M
LIB	モンテカルロシミュレーション	907 M
NN	神経回路シミュレーション	68M

算器の使用履歴をハードウェアで取得しつつ予測する方法 などが考えられる.なお,後述の初期評価では,1warpの スレッド数をあらかじめ半分に制限することで,warp分 割による効果を見積ることとする.演算器の使用率に応じ たwarp分割の制御については今後の課題である.

4. 評価

スレッド発行制御の有無で演算器のアイドルサイクル の変化および PG による消費エネルギー削減効果を評価 するために,提案手法をサイクルレベルシュミレータの GPGPU-Sim (version 3.1.1)[8] に実装して評価を行った. 評価に用いた GPU の仮定は NVIDIA GeForce GTX 480 に従い (表1)のようにした.また,ベンチマークについて は CUDA で記述された 6 つのベンチマークプログラムを 用いて評価を行った(表 2). BET については, 50 サイク ル,100 サイクル,200 サイクルの3 パターンを評価する. なお,各演算器がアイドルになった際にPGをするかどう かの制御は理想的にできるものとし, BET 以上のアイド ル時のみPGをする仮定して評価を行った.リーク消費エ ネルギー削減効果の見積りには,シミュレータから得られ た演算器の使用とアイドルサイクルのログを利用し, BET 以上のアイドルサイクルの合計から PG 回数 × BET を差 し引いて,正味のリーク電力削減サイクルを求め,合計実 行サイクル数に対するリーク電力削減サイクルを評価指標 とする.

提案手法のスレッドコンパクションについては,ログを解 析することでコンパクションした際の使用演算器を導出し



図 8 nomal (BFS)



スレッドコンパクション (BFS) 図 9

評価を行った.また, warp 分割については, GPGPU-Sim のパラメーター設定で 1warp 内のスレッド数を 32 スレッ ドから 16 スレッドに変更してシミュレーションを行うこ とで,静的に warp 分割を行ったと仮定して評価を行った.

評価結果 5.

5.1 演算器ごとのリーク消費エネルギー変化

まず,提案手法を用いることで,SIMD内の各演算器の リーク消費エネルギーがどのように変化するかを調査する.

図8,図9,図10,図11に,提案手法を用いない場合 (normal),スレッドコンパクションを用いた場合(comp), warp 分割を用いた場合 (div), スレッドコンパクションと warp 分割の両者を用いた場合 (mix)の,合計実行サイクル 中でリーク消費エネルギーが削減できるサイクルの割合を 示す.なお,これらはベンチマーク中のBFSを実行し,0 番 SM コアの1 つの SIMD ユニット内の演算器の結果であ る. 横軸は当該 SIMD ユニット内に 16 個ある演算器の ID を示している.また,各演算器について BET が 50 サイク ル,100 サイクル,200 サイクルの場合の結果を示している.

図8のスレッド発行制御前の結果を見ると,各演算器の リーク消費エネルギー削減は均等になる傾向があることが わかる.これは,各演算器でのスレッド実行が均等に行わ れるためである.一方,図9のスレッドコンパクションを 適用した場合の結果を見ると,演算器 ID が1に近い側の 演算器では normal に比べてリーク消費エネルギーが増加



Vol.2013-ARC-204 No.2

2013/3/26



し,反対に演算器 ID が 16 に近い側の演算器ではリーク消 費エネルギー削減効果が増大している.これは,スレッド コンパクションにより,スレッドの各演算器への割り当て が ID-1 の側に偏ったためである.この結果,およそ ID-1 から ID-6 の演算器では短い期間のアイドルが増加する-方, ID-7から ID-16では長い期間のアイドルが多くなる. そのため,後者のリークエネルギーの削減割合が大きく なっている.

図 10 の warp 分割を適用した場合には, ID-1 から ID-8 までの演算器では短い期間のアイドルが増加するが, ID-9 から ID-16 の演算器は全く使用しないため,全サイクルで リーク消費エネルギーを削減できている,一方で,ID-1か ら ID-8 までのリーク消費エネルギーは normal に比べて増 加する傾向にある.

図 11 のスレッドコンパクションと warp 分割の両者を適 用した場合の結果を見ると,warp分割の効果により,ID-9 から ID-16 の演算器は全サイクルでリーク消費エネルギー が削減できている.また,スレッドコンパクションの効果 により,図10に比べてID-1からID-3までの演算器のリー ク消費エネルギー削減効果は減少するが, ID-4 から ID-8 までの演算器のリーク消費エネルギー削減効果は増加して いることがわかる.

なお,全ての結果において,当然ながら BET が大きく なるほど PG によるリーク消費エネルギー削減効果は減少 してしまうことも,図より見てとれる.

情報処理学会研究報告

IPSJ SIG Technical Report



図 12 各種法ごとのリークエネルギー削減率 (BET=50)







5.2 各手法のリーク電力削減効果の比較

図 12,図 13,図 14にBETを50サイクル,100サイク ル,200 サイクルとした場合のベンチマークごとのリーク 消費エネルギー削減サイクルを示す.

全体としてスレッドコンパクション, warp 分割を行う ことで手法でリークエネルギー削減の効果が得られること がわかる.また,スレッドコンパクションと warp 分割の 両者を適用するとその効果が増加することがわかる.

BFSとMUMは分岐命令を含んでいる.また,演算器の 使用率が低いベンチマークであるため, BET が 100 サイ クルの場合に normal に比べて comp の効果は BFS で 19 ポイント増加, MUM で 6 ポイント増加した. div の効果 は BFS で 23 ポイント増加, MUM で 17 ポイント増加と なった.また, compとdivの相乗効果により normal に比



べて mix では BFS で 27 ポイント増加, MUM で 18 ポイ ントの増加となった.

LPS と LIB は warp 中のスレッド数が多く演算器の使用 率が高いベンチマークのため, comp のみによるリークエ ネルギー削減の効果は1ポイント未満であった.一方で, warp 分割では使用する演算器数を半分にするためリーク 電力削減の効果は BET が 100 サイクルの場合に LPS で 62 ポイントの増加, LIB で 53 ポイントの増加となった.

NNは warp 中のスレッド数が少なく, normal の場合に おいても特定の演算器のみでスレッドの実行が行われるた めリークエネルギー削減量の向上はBET が 100 サイクル の場合に mix でも 2 ポイント増程度であり,他のベンチ マークに比べ提案手法の効果が現れない結果となった.

warp 分割手法を演算器使用率が高い場合に用いると性能 低下を引き起こす可能性がある.そこで図 15 に, normal に対する warp 分割を用いた場合の各ベンチマークの相対 性能を示す.なお,スレッドコンパクションは性能に影響 しないため、ここでは評価結果を示していない. warp 分割 を適用し,1ワープのスレッド数を32スレッドから16ス レッドに削減した場合でも、演算器の使用率が低いBFS、 MUM, NN ではほとんど実行サイクルが変わっていない. そのため, このようなベンチマークでは warp 分割手法に よりほとんど性能低下なく大幅なリーク消費エネルギー削 減効果を得ることができる.一方で,演算器の使用率が高 い LPS では 2.1 倍, LIB では 1.8 倍の実行サイクル数にな る結果となった.このことからも,演算器使用率が高い場 合には warp 分割を用いるべきではなく,適応的に利用す ることが重要である.この, warp 分割手法の適応的な制 御については今後の課題である.

6. 関連研究

Wang ら [3] は GPU に搭載されているキャッシュを PG することで,走行時のリーク電力削減を行なっている.特 に,L1,L2 キャッシュのそれぞれに個別の省電力モード を用意することで PG の効果を高める手法を提案してい る.理論値では53%のリーク電力を削減できると報告され ている.

Rhuら [7] は,演算器使用効率の向上を目的としたスレッ ドコンパクション手法について,スレッドスケジューリン グ上のデメリットを考慮しつつ,コンパクションを行うか どうかを適応的に判断することにより,性能を向上させる 手法を提案している.

関ら [2] は, MIPS R3000 互換のプロセッサにおいて, コアの粒度ではなく,演算器の粒度でPG を行う手法を提 案している.演算器が長期間アイドルになる場合にのみス リープを行えるように,OS やコンパイラがハードウェアに よる PG を制御する手法などが述べられている.武藤ら [6] は,同じく MIPS R3000 プロセッサにおいて,演算器の アイドル履歴に基づいたスリープ制御手法を検討し,評価 を行なっている.この中で提案している Index 方式では, スリープ時間の履歴情報に基づいて Long Sleep か Short Sleep かを判定しスリープ制御を行うことにより,Non PG と比較して ALU で 35%のリーク電力を削減している.

本稿では,コンパクションなどを含めたスレッドの発行 制御を工夫することで,GPU内のSIMD演算器の各演算 器単位で細粒度にPGを行いリーク電力を削減することを 目的としている.この点で,上記の研究とは大きく異なる ものである.

7. おわりに

本研究では,GPU における SIMD 演算器を対象に,細 粒度パワーゲーティングによるリーク削減を目的としたス レッド発行制御手法を提案し,その初期評価を行った.提 案手法は各演算器がなるべく長期間 PG できるようにする ために,スレッドコンパクションと warp 分割を行うもの である.

シミュレーションによる初期評価の結果,BET が 100 サイクルの場合に通常のスレッドスケジューリングの場 合のリーク電力削減率は 42%,スレッドコンパクション, および warp 分割を適用した場合のリーク電力削減率はそ れぞれ 46%,71%になった.また両者を組み合わせた場合 のリーク電力削減率は74%になることがわかった.また, その際の性能は,演算器使用率が低い場合にはほとんど変 化しない結果となった.以上より,提案手法によりスレッ ド発行制御を適切に用いることで,性能低下を抑えつつ, リーク電力を効率的に削減できることがわかった.

今回は,実際に動的にスレッド発行制御を行ったもので はなく,ログの解析や静的に発行幅を設定することで,提 案手法のリーク電力削減効果を見積もった.今後,実際に シミュレータ上にスレッド発行制御機構を実装すること や,アプリケーションの特性に応じて制御を使い分けるこ とで性能低下を抑えつつ効率的にPG が行えるような手法 の開発に取り組んでいく予定である.

謝辞 本研究の一部は,科学技術振興機構・戦略的創造

研究推進事業 (CREST) の研究プロジェクト「ポストペタ スケールシステムのための電力マネージメントフレーム ワークの開発」,ならびに JSPS 科研費 24680004 の助成に より行われたものである.

参考文献

- Sunpyo Hong, Hyesoon Kim. An integrated GPU power and performance model. Proceedings of the 37th annual international symposium on Computer architecture (ISCA '10), pp.280-289, 2010.
- [2] 関直臣ら、"MIPS R3000 プロセッサにおける細粒度動的スリープ制御の実装と評価"、電子情報通信学会論文誌 J93-D 巻 6 号, pp.920-930, 2010 年.
- [3] Yue Wang, Soumyaroop Roy, Nagarajan Ranganathan. Run-time power-gating in caches of GPUs for leakage energy savings. Design, Automation & Test in Europe Conference & Exhibition (DATE12),pp.300-303,2012.
- [4] Po-Han Wang, Chia-Lin Yang, Yen-Ming Chen, Yu-Jung Cheng. Power gating strategies on GPUs. ACM Transactions on Architecture and Code Optimization (TACO) Volume 8 Issue 3 Article 13, 2011.
- [5] Wilson W. L. Fung, Ivan Sham, George Yuan, Tor M. Aamodt. Dynamic Warp Formation and Scheduling for Efficient GPU Control Flow. Proceedings of the 40th Annual IEEE/ACM International Symposium on Microarchitecture (MICRO 40), pp.407-420.
- [6] 武藤 徹也, 宇佐美 公良. "細粒度パワーゲーティングに おける履歴に基づいたスリープ制御方式の検討と評価"
 電子情報通信学会技術研究報告. VLSI 設計技術 110 巻, pp.31-36, 2011 年
- [7] Minsoo Rhu, Mattan Erez. CAPRI: Prediction of Compaction-Adequacy for Handling Control-Divergence in GPGPU Architectures. Proceedings of the 39th annual international symposium on Computer architecture (ISCA '12), pp.61-71, 2012.
- [8] Yuan, G.L. Fung, W.W.L. Wong, H. Aamodt, T.M. Analyzing CUDA workloads using a detailed GPU simulator. Proceedings of the Performance Analysis of Systems and Software 2009 (ISPASS), pp.163-174, 2009
- [9] Wilson W. L. Fung, H. Aamodt, T.M. Thread block compaction for efficient SIMT control flow. Proceedings of the 2011 IEEE 17th International Symposium on High Performance Computer Architecture (HPCA '11), pp.25-36, 2011