

ウェイ適応型キャッシュの高エネルギー効率化のための デッドブロック早期追だしポリシー

東方 雄亮[†] 佐藤 雅之[†] 江川 隆輔^{†,††}
滝沢 寛之^{†,†††} 小林 広明^{††,†††}

1. はじめに

近年、キャッシュメモリ（以下、キャッシュ）の低消費エネルギー化が求められている。低消費エネルギーを実現するためのキャッシュとして、ウェイ適応型キャッシュ¹⁾が提案されている。ウェイ適応型キャッシュは、スレッドの性能維持に必要なキャッシュウェイへの電源供給を停止することによって、キャッシュの低消費エネルギー化を実現する。しかし、キャッシュのデータ管理に LRU 置換ポリシーを用いているため、再利用されないブロックであるデッドブロックがキャッシュを占有する場合がある。ウェイ適応型キャッシュは、これらの性能向上に貢献しないデッドブロックも含めて、性能維持に必要なキャッシュ容量を見積もる。そのため、デッドブロックを保持するためのキャッシュ領域において電力が浪費され、ウェイ適応型キャッシュのエネルギー効率が著しく低下する。

本研究では、ウェイ適応型キャッシュのさらなるエネルギー効率向上を実現するために、デッドブロックを考慮したデータ管理ポリシーを提案する。提案ポリシーでは、デッドブロックの早期追出しによって、キャッシュを占有するこれらのブロックを削減し、電源供給を停止できるウェイ数を増加させる。その結果、ウェイ適応型キャッシュにおいてさらなるエネルギー効率向上が期待できる。

2. デッドブロックのキャッシュ占有

ウェイ適応型キャッシュでは、データ管理ポリシーとして LRU 置換ポリシーが用いられている。LRU 置換ポリシーでは、下位の階層へ書き戻されるブロックを選択するための優先度の順序列をキャッシュセット毎に考える。8 ウェイセットアソシアティブキャッシュの優

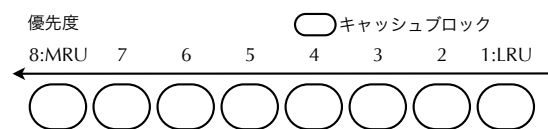


図 1 優先度の順序列

先度の順序列を図 1 に示す。キャッシュミスの場合、優先度が最も高い位置 (MRU) にアクセスされたブロックが挿入され、最も低い位置 (LRU) のブロックが下位のメモリ階層に書き戻される。キャッシュヒットの場合、アクセスされたブロックは MRU に昇格される。このような動作によって、最も長い間アクセスされていないブロックを優先的に追出し、時間的局所性に基づくデータアクセスを高速化できる。

しかし、LRU 置換ポリシーでは、将来再利用されないデッドブロックの場合でも、これらのブロックは MRU に置かれる。この結果、これらのブロックが追い出されるまでには長い期間が必要になるため、キャッシュ中におけるデッドブロックの割合が増大する。その結果、デッドブロックを保持するためのキャッシュ領域が浪費される。そのため、ウェイ適応型キャッシュにおいてエネルギー効率を高めるためには、再利用されるブロックが追い出されることを抑制しつつ、デッドブロックを早期に追い出す置換ポリシーが必要である。

デッドブロックの中でも、挿入されてから一度も再利用されないデッドオンフィルブロックが特に多い。デッドオンフィルブロックを早期に追い出すポリシーには、動的挿入ポリシー (以下、DIP)²⁾がある。DIP は、MRU と LRU に挿入位置を一定確率で切り替える二方向性挿入ポリシーと LRU 置換ポリシーを動的に選択することによってキャッシュミス減らす置換ポリシーである。しかし、いずれの置換ポリシーが選択されても、デッドオンフィルブロックが MRU へ挿入される場合がある。また、DIP における昇格は LRU 置換ポリシーと同様であるため、デッドブロックが MRU に昇格され、キャッシュに長時間保持される。以上より、DIP を適用したウェイ適応型キャッシュでは、デッドブロックのキャッシュ占有率をさらに削減する余地があると考えられる。

3. デッドブロック早期追出しポリシー

デッドオンフィルブロックを早期に追い出すためには、挿入位置を LRU へ近づける必要がある。しかし、挿入位置を LRU へ過剰に近づけるとブロックが再利用される前に追い出される恐れが高まる。したがって、再利用されるブロックを保持しつつ、デッドオンフィルブロックを早期に追い出せるように優先度の順序列での適切な挿入位置を選択する必要がある。参照の局所性に基づき、ブロックが挿入された位置から LRU

[†] 東北大学大学院情報科学研究科
^{††} 東北大学サイバーサイエンスセンター
^{†††} 科学技術振興機構戦略的創造研究推進事業

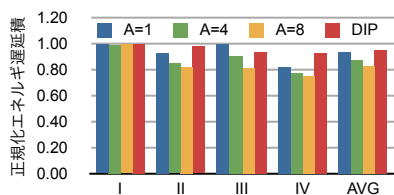


図2 $R = \infty$ におけるエネルギー遅延積

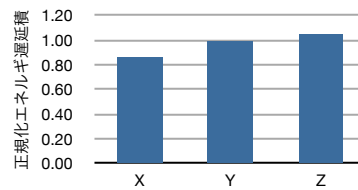


図3 $A = 8, R = 1$ におけるエネルギー遅延積

に近づくにつれて、ブロックに対する再参照回数は少なくなる傾向にある。そのため、優先度の順序列のいかなる位置にブロックが挿入されても、初回再参照回数はLRUが最小となる。したがって、優先度の順序列におけるLRUへの初回再参照回数を観測することによって、再利用されるブロックが最低優先度の位置まで降りてきていることを検知できる。そして、再利用されるブロックを保持しつつデッドオンフィルブロックを早期に追い出すために、最低優先度の位置以上に挿入位置をLRUに近づけないようにする。LRUにおける初回再参照回数を C 、初回再参照回数の過不足を A とする時、LRUへの初回再参照が多発している場合($C \geq A$)には、挿入位置をMRU方向へ1つ移動する。一方で、LRUへの初回再参照が全くない場合($C < A$)には、挿入位置をLRU方向へ1つ移動する。以上のように適切な挿入位置の選択を行う。

また、昇格後に再利用されないデッドブロックをキャッシュから早期に追い出すためには、キャッシュヒットしたブロックを必要以上に昇格させないことが重要である。しかし、昇格を全くしないと再利用されるブロックも早期に追い出される恐れがある。したがって、昇格を停止する適切な再参照回数を選択する必要がある。本研究では、キャッシュヒットしたブロックを R 回の再参照までMRUに昇格させ、 $R+1$ 回目の再参照では昇格させないようにする。

4. 評価実験

本評価では、LRU置換ポリシおよびDIPを用いたウェイ適応型キャッシュと、提案手法を用いたウェイ適応型キャッシュのエネルギー遅延積を比較する。なお、ベンチマークにはSPEC CPU2006を用いる。

まず、挿入方式を変更した場合におけるエネルギー遅延積を図2に示す。なお、各手法のエネルギー遅延積はLRU置換ポリシを用いた場合のエネルギー遅延積によって正規化されている。また各ベンチマークは、再利用されるブロックの割合が多く、再参照の間隔が長い場合にIカテゴリへ、短い場合にIIカテゴリへ分類されている。さらに、再利用されるブロックの割合が少なく、再参照の間隔が長い場合にIIIカテゴリへ、短い場合にIVカテゴリへ分類されている。図2より $A=8$ の場合において、LRU置換ポリシに対して16.6%、DIPに対して12.4%エネルギー遅延積が削減されていることがわかる。したがって、デッドオンフィルブロックの削減によって、ウェイ適応型キャッシュのエネルギー効率向上が明らかになった。なお、Iカテゴリではエネルギー遅延積が全く削減されない。この

原因として、デッドオンフィルブロックのキャッシュ占有率が低いことと、再参照の間隔が長いために挿入位置をLRUに近づけられないことが挙げられる。

昇格方式を変更した場合におけるエネルギー遅延積を図3に示す。なお、大部分のブロックは1回から4回までの再参照でデッドブロックになる。そのため時間的局所性より、1回MRUへ昇格されると、追い出されるまでに多くのブロックがデッドブロックになると考えられる。したがって、2回目以降昇格をしない $R=1$ の場合において昇格方式を変更した提案手法の評価を行う。エネルギー遅延積は $A=8, R=\infty$ の場合における提案手法のエネルギー遅延積によって正規化されている。また、各ベンチマークは、エネルギー遅延積の削減効果に応じて3つのカテゴリに分類されている。Xカテゴリはエネルギー遅延積が削減されるカテゴリであり、Yカテゴリは変化しないカテゴリである。また、Zカテゴリはエネルギー遅延積が増大するカテゴリである。図3より、昇格方式を変更することによって、Xカテゴリでは14.0%エネルギー遅延積が削減されていることがわかる。しかし、Zカテゴリではエネルギー遅延積が5.2%増大する。昇格回数の制限が有効であるベンチマークの詳細な特徴分析と、動的な R の設定方法については今後の課題とする。

5. おわりに

ウェイ適応型キャッシュのエネルギー効率向上を目的として、デッドブロックをキャッシュから早期に追い出す置換ポリシを提案した。提案手法によってエネルギー遅延積を最大50.6%、平均18.0%削減できた。したがって、提案手法によるウェイ適応型キャッシュのエネルギー効率向上の可能性が示された。

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によるものである。

参考文献

- 1) Kobayashi, H., Kotera, I. and Takizawa, H.: Locality Analysis to Control Dynamically Way-Adaptable Caches, *ACM SIGARCH Computer Architecture News*, Vol. 33, No. 3, pp. 25–32 (2005).
- 2) Qureshi, M. K., Jaleel, A., Patt, Y. N., Steely, S. C. and Emer, J.: Adaptive Insertion Policies for High Performance Caching, *ACM SIGARCH Computer Architecture News*, Vol. 35, No. 2, p. 381 (2007).