

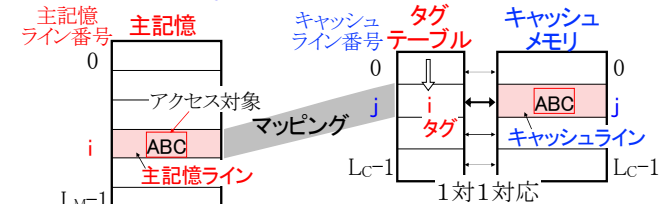
# 計算機方式論

## 第13章 キャッシュ -マッピング、書込アクセス-

1

## キャッシュと主記憶とのマッピング

- 主記憶ライン  $i$  とキャッシュライン  $j$  との対応付けは、**タグテーブル**なるマッピング表を使った連想写像方式(フルアソシアティブマッピング、セットアソシアティブマッピング)などで行う。  
主記憶ライン  $i$  がキャッシュライン  $j$  にコピーされているとき、**タグテーブル**の第  $j$  エントリには**タグ  $i$** (主記憶ライン番号)が入る。

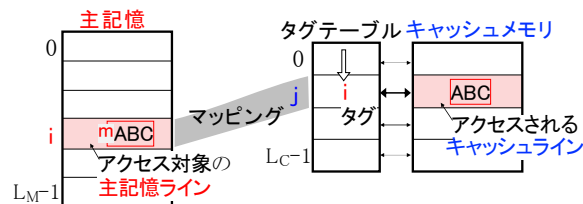


ラインサイズ:  $S_L$  主記憶容量:  $S_M$  キャッシュ容量:  $S_C$   
主記憶ライン数:  $L_M = S_M / S_L$  キャッシュライン数:  $L_C = S_C / S_L$

2

## キャッシュラインの決定(連想写像方式の場合)

- アクセスする主記憶番地  $m$  から、主記憶のライン番号  $i$  を求め ( $i = m / S_L$ , ライン内番地  $= m \% S_L$ )、 $i$  を**タグ**として、**タグテーブル**を連想検索(フルアソシアティブマッピング)する。
- ① **タグ  $i$** と一致するエントリ  $j$ があれば、**ヒット**で、キャッシュライン  $j$  をアクセス。  
② **タグ  $i$** と一致するエントリがなければ、**ミスヒット**で、**ミスペナルティ処理**を行う。



3

## マッピング方式

- ① **ダイレクトマッピング**(direct mapping)
- ② **フルアソシアティブマッピング**  
(full associative mapping)
- ③ **セットアソシアティブマッピング**  
(set associative mapping)

4

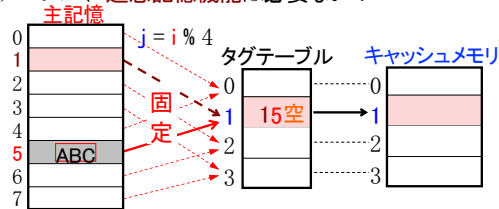
## ①ダイレクトマッピング

- 主記憶ライン  $i$  に対するキャッシュライン  $j$  は **予め一意に定めておく**。  
例えば、 $j = i \% L_c$  としたとき、

- ① タグテーブル  $j$  番目に主記憶ラインのタグ  $i$  が在るときは、**ヒット**で、**キャッシュライン  $j$**  をアクセス。
- ②  $j$  番目が**空き**のときは、**ミスヒット**で、**主記憶ライン  $i$**  をロードし、**タグ  $i$**  を書き込む。
- ③  $j$  番目に**違うタグ**が在るときも、**ミスヒット**で、そのキャッシュラインを**追い出した後**、**主記憶ライン  $i$**  をロードし、**タグ  $i$**  を書き込む。

×ライン置換の選定の自由度が低い。

◎タグテーブルに**連想記憶機能**は**必要ない**！



5

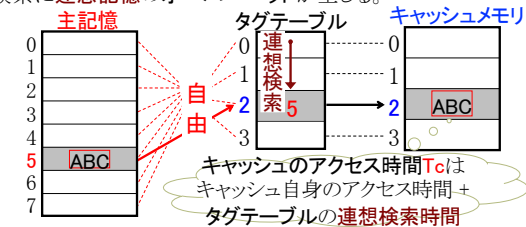
## ②フルアソシアティブマッピング

- 主記憶ラインをどのキャッシュラインにも**自由**にマッピング。  
タグテーブル中、主記憶ラインの**タグ**を**連想検索**で探す。  
タグが見つければ、**ヒット**でアクセス。なければ**ミスヒット**で、**空きキャッシュライン**があれば、**主記憶ライン**をロードし、タグテーブルに**タグ**を書き込む。  
**空き**がなければ、**ライン置換アルゴリズム**で**追い出すキャッシュライン**を決め、そこに**主記憶ライン**をロードし、タグテーブルに**タグ**を書き込む。

◎ライン置換アルゴリズムの選定の自由度が高い。

×タグテーブル全体に**連想記憶機能**が必要なため、**実装コストが高い**。

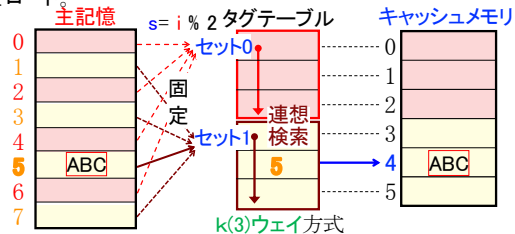
×タグ検索に**連想記憶**のオーバーヘッドが生じる。



6

## ③セットアソシアティブマッピング

- ダイレクトマッピングとフルアソシアティブマッピングを組み合わせた方式。
- キャッシュメモリを  $N$  個の**セット**に分割し、主記憶ライン  $i$  から**セット**へのマッピングはダイレクト方式で行う。例えば、**セット番号  $s = i \% N$** 。
- 各**セット**は、対応するタグテーブルに**連想記憶機能**をもたせてあり、**セット  $s$**  中を**タグ  $i$**  で**連想検索**で探す。
- タグ  $i$  が見つければ**ヒット**でアクセス、見つからなければ**ミスヒット**で、**空きキャッシュライン**があれば**主記憶ライン**をロード、**空き**がなければ、**ライン置換アルゴリズム**で**追い出すキャッシュライン**を決め、そこに**主記憶ライン**をロード。



7

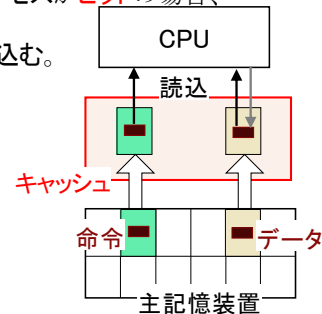
## kウェイセットアソシアティブマッピング

- 1**セット**を  $k$  個の**キャッシュライン**で構成するとき、**kウェイセットアソシアティブマッピング**( $k$ -way set associative mapping)という。  
 $k$ は2～16程度(前図は**3ウェイ**)
- キャッシュメモリの**セット数  $N = Lc/k$**
- $N$  個の**連想記憶タグテーブル**が必要だが、フルアソシアティブマッピングと比べ、  
テーブルのエントリの数は  $N$  分の1のため、  
**検索時間が短く**、ハードウェア実装コストが格段に下がる。
- $k=1$ ウェイ、 $N = Lc$ セットの場合、**ダイレクトマッピング**。
- $k=Lc$ ウェイ、 $N=1$ セットの場合、**フルアソシアティブマッピング**。

8

## キャッシュの読込アクセス

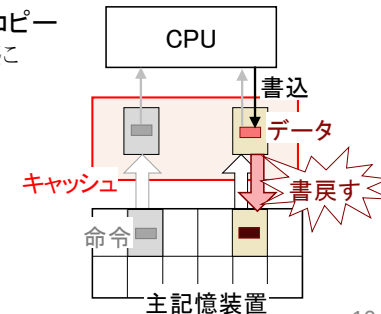
- **命令**キャッシュは、読込アクセスだけ。
- **データ**キャッシュは、読込と書込の2つのアクセス。
- 命令やデータの**読込**アクセスが**ヒット**の場合、キャッシュラインからその命令・データを読み込む。



9

## キャッシュの書込アクセス

- データの**書込**アクセスが**ヒット**の場合、**キャッシュ**の内容を書き換える(対象のデータの書き換え)
- キャッシュは主記憶のコピーなので、変更を主記憶に反映するため、**主記憶**に書き戻す！これを、**主記憶の更新**という！



10

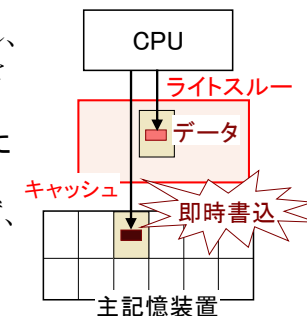
## 主記憶の更新

- 「**主記憶ライン**と**キャッシュライン**との内容の**同一性**」を**コヒーレンシ**(coherency)といい、これを保つため、キャッシュ制御機構は、**書込アクセスのヒット時には主記憶更新**(主記憶ラインへの書き戻し)を行わなければならない。
- **いつ**の時点で主記憶ラインを更新するかで、2つの方式がある。
  - ①**ライトスルー**(ストアスルー)  
キャッシュの書き込みと**同時**に主記憶更新。
  - ②**ライトバック**(ストアバック,コピーバック)  
キャッシュの書き込み**時**には主記憶更新しないで、当該キャッシュラインが**追い出し**の対象になったとき、主記憶更新。

11

### ①ライトスルー(write-through)

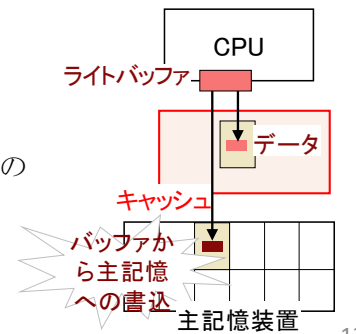
- キャッシュの書き込みと**同時に**主記憶にも書き込む(主記憶更新する)。
- コヒーレンシを常に保証し、主記憶更新のタイミングを図らなくてよい。
- 書込アクセスが、**実質的に主記憶アクセス**なので、キャッシュ効果があがらず、**書込アクセスの性能が改善されない。**



12

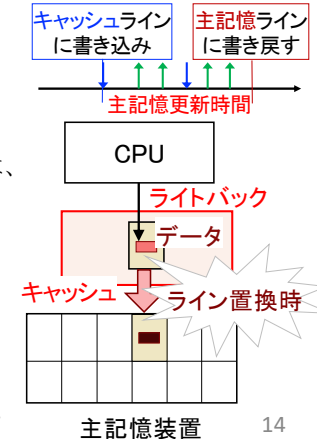
## ライトスルー (置いてきぼり制御)

- プロセッサにライトバッファ(キュー)を置き、書込データをバッファに一時記憶する。
- 主記憶へのアクセスを後回しにするため高速化が望める。
- バッファの読み出し、マルチプロセッサ方式のときに問題が生じる。



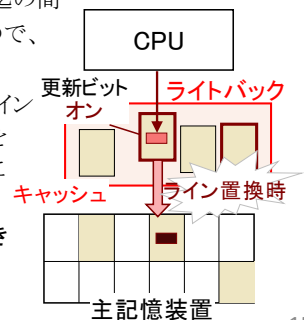
## ②ライトバック(write-back)

- キャッシュの書き込み時には、主記憶更新しない。
- 書き込み対象を含むラインを主記憶を更新する**時点**は、
- **アソシアティブマッピング**方式では、空きラインが必要になり、**ライン置換アルゴリズム**に当該キャッシュラインが選ばれ、追出し対象になったとき (**ライン置換時**)。
- **ダイレクトマッピング**方式では、当該キャッシュラインが選ばれたが、**キャッシュタグ**が異なるとき。



## ライトバック-更新ビット

- ライトバックは、**ヒット**時にはキャッシュ性能を十分発揮でき、主記憶へのアクセス回数を減らすことができる。
- キャッシュ更新から主記憶更新迄の間は**コヒーレンシが壊されている**ので、それを管理し保つ機構が必要。
- **タグテーブル**にそのキャッシュラインへの書き込みを示す**更新ビット**を用意し、書き込まれたら**“オン”**にする。**ライン置換時**に追い出し対象に選ばれたら、**“オン”**のときだけ主記憶ラインに書き戻し、**“オフ”**のときは書き戻さない。



## 書込アクセスがミスヒットの場合

- **ノーライトアロケート**(no write allocate)  
主記憶への書き込みだけで、主記憶ラインのキャッシュへの読み出しは行わない。
- **ライトアロケート**(write allocate)  
当該の主記憶ラインをキャッシュに読み出し、つぎに**キャッシュ**と**主記憶**の対象を書き換える。

