

計算機方式論

第14章 キャッシュ -置換アルゴリズム等-

1

ライン置換アルゴリズム

- キャッシュ制御機構は、**読出アクセスのミスヒット時**に、アクセス対象を含む主記憶ラインをキャッシュに読み込む。
- フル/セットアソシアティブマッピングでは、キャッシュに空きラインがない場合、**空きにするキャッシュライン**を選び、その更新ビットが“**オン**”のときそのキャッシュラインを主記憶に**追い出し**、更新ビットが“**オフ**”のときそのキャッシュラインを**廃棄**する。
- **書込アクセスのミスヒット時**にも、ライトアロケート方式では、主記憶ラインをキャッシュに読み込むため、空きラインがない場合、キャッシュから**追い出す**または**廃棄する**ラインを選ぶ。
- どのラインをキャッシュから**追い出す**または**廃棄する**かは**ライン置換アルゴリズム**で決め、ハードウェアで実現する。

2

色々なライン置換アルゴリズム

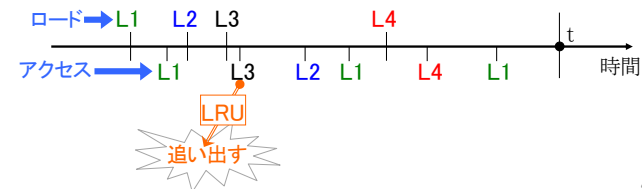
- ①LRU(Least Recently Used)
- ②FIFO(First In First Out)
- ③FINUFO(First In Not Used First Out)
- ④LFU(Least Frequently Used)
- ⑤RANDOM

- キャッシュラインの種々の**属性**(更新ビット、キャッシュへのコピー時刻、アクセス時刻、アクセス回数等)は、**タグテーブル**に記述する。

3

①LRU(Least Recently Used)

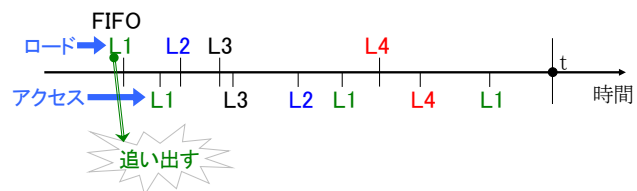
- 最後に参照されてからの経過時間の最も長い**キャッシュライン**を追い出す。
- 参照局所性を使っているので、ヒット率は良好だが、アクセス時刻を管理する機構が必要になり、判定時間がオーバーヘッドになる。



4

②FIFO(First In First Out)

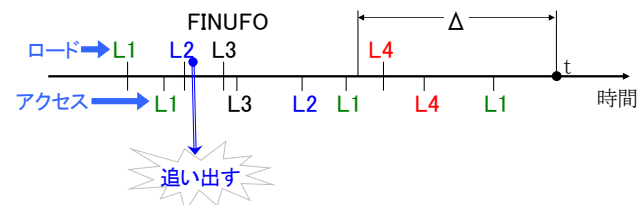
- 一番最初にキャッシュにコピーしたラインを追い出す。
- LRUに比べ、参照局所性は劣るが、機構は簡単になる。



5

③FINUFO(First In Not Used First Out)

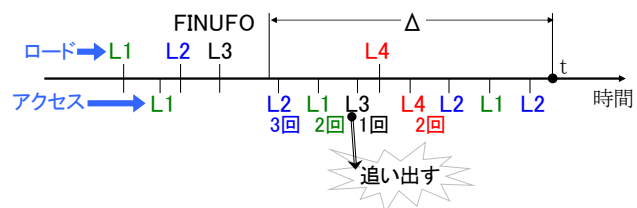
- 一定時間アクセスのないラインの中で、最初にキャッシュにコピーしたラインを追い出す。
- WS(Working Set)法のひとつ。



6

④LFU(Least Frequently Used)

- 一定時間内でアクセス回数が最小のラインを追い出す。
- WS(Working Set)法のひとつ。



7

⑤RANDOM

- 無作為にラインを決める。機構が簡単。

8

物理キャッシュと論理キャッシュ

- 仮想記憶を使用する計算機
番地変換機構 論理番地⇒物理番地
- 番地変換機構**を
CPU-?-キャッシュ-?-主記憶
のどこに?入れるかで、2つの方式がある。
 - ①物理キャッシュ
 - ②論理キャッシュ

9

物理キャッシュ

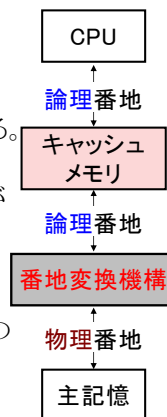
- 番地変換機構**が
CPUとキャッシュ間に入る。
- キャッシュは**物理番地**で指定。
- キャッシュにアクセスするとき
番地変換が必要となり、
キャッシュアクセス時間が長くなる。
- 実アドレスキャッシュともいう。



10

論理キャッシュ

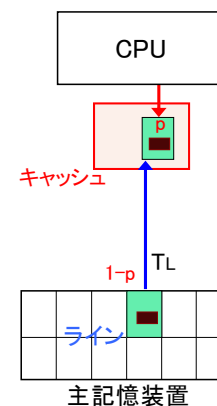
- 番地変換機構**が
キャッシュと主記憶間に入る。
- キャッシュは**論理番地**で指定。
- キャッシュのアクセスが**速やか**に行える。
- ラインのマッピング機構に番地変換機構が必要となり、コヒーレンスの保持が複雑な処理になる。
- 仮想記憶空間を切り替える度に、
キャッシュの内容を**無効**にするため
フラッシュを行うので、キャッシュ容量の
大きいときはオーバーヘッドになる。
- 仮想アドレスキャッシュともいう。



11

2階層キャッシュへ

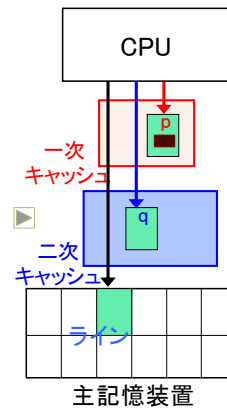
- キャッシュの実効アクセス時間 T_e
 $T_e = T_c + (1-p)T_L$
- キャッシュ容量を増やす!**
⇒キャッシュアクセス時間 T_c も増大
⇒ T_e が減るとは限らない!
 T_L は主記憶なので大きい
⇒高速の**二次キャッシュ**を設ける



12

2階層キャッシュ

- **キャッシュメモリ** (チップ内) と主記憶との間に **高速の二次キャッシュ** を設ける
- ⇒ **一次キャッシュ** でミスヒットしても、**二次キャッシュ** でヒットすれば、 T_e を縮めることができる!
- 二次キャッシュの容量は、**数Mバイト** (intel core i7ではL3に相当)



13

2階層キャッシュの実効アクセス時間

- 2階層キャッシュのアクセス手順

一次キャッシュ

ヒット

一次キャッシュラインをアクセス

ミスヒット

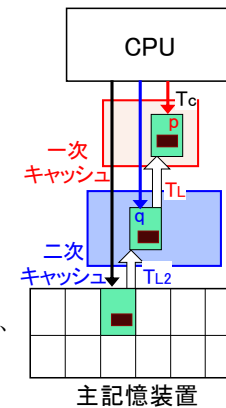
二次キャッシュ

ヒット

一次キャッシュにラインをコピー、
一次キャッシュラインをアクセス。

ミスヒット

主記憶ラインを、二次キャッシュに、
その二次キャッシュラインを
一次キャッシュに、コピー。
一次キャッシュラインをアクセス。



14

2階層キャッシュの実効アクセス時間

- 2階層キャッシュの実効アクセス時間 T_e

$$T_e = pT_c + (1-p)\{q(T_c + TL_1) + (1-q)(T_c + TL_1 + TL_2)\}$$

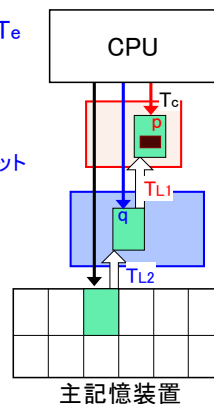
$$= T_c + (1-p)TL_1 + (1-p)(1-q)TL_2$$

p : 一次キャッシュのヒット率

q : 二次キャッシュのヒット率

TL_1 : 二次から一次へのコピー時間

TL_2 : 主記憶から二次へのコピー時間



15

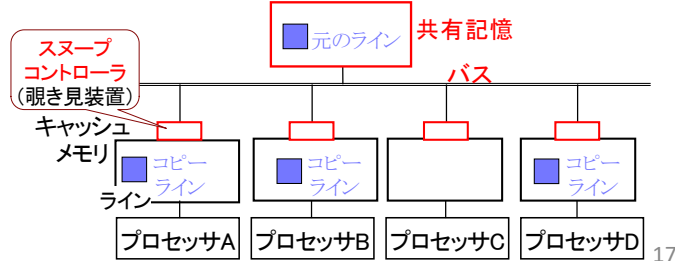
共有記憶型並列計算機

- **共有記憶型並列計算機**では、複数個のプロセッサによる**共有記憶**へのアクセスが集中!
.... プロセッサ個々に共有記憶を**分散**させたり、
プロセッサ個々に**キャッシュメモリ**を設けて回避
- **共有バス方式**
共有バス上に複数のプロセッサを配置。
スヌープキャッシュ法等がある。
- **ディレクトリ方式**
プロセッサが一般のネットワーク上に配置。
共有記憶は、各プロセッサ近くに**分散配置 (分散共有)**され、
各プロセッサに付加された**ディレクトリ**が、各ライン情報をもつ。

16

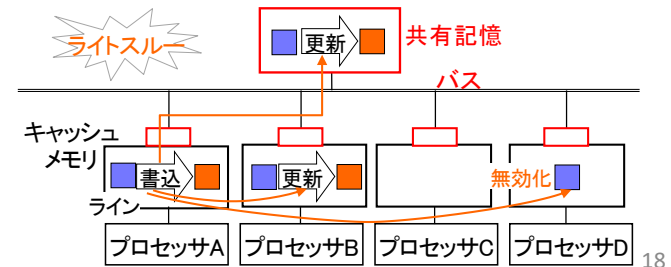
スヌープキャッシュ

- バス接続の**共有記憶**に適した管理方式
- 共有記憶**ラインを各プロセッサのキャッシュにコピーしアクセス
- 共有記憶へのアクセスやその番地等は**バス**を通して、すべてのプロセッサが同時に知ることができる(**スヌープコントローラ**)。
- 各プロセッサがバスを監視していれば、自分のもつキャッシュラインの**更新**やどのプロセッサが**同じコピー**をもつかを管理できる。



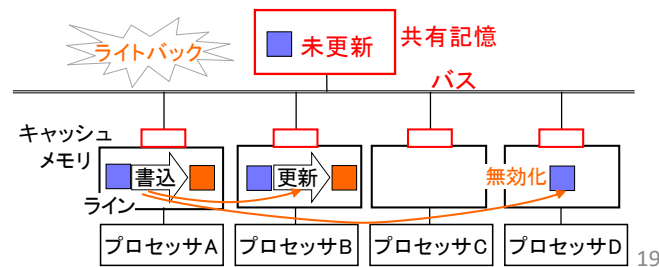
書き込み-ライトスルー

- ①プロセッサAのラインに**書き込み**
- ②共有記憶のラインを**更新**
- ③他のプロセッサのキャッシュ中のラインは、**更新**するか、**無効化**する。



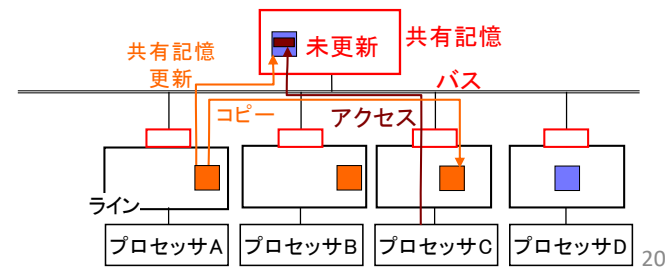
書き込み-ライトバック①

- ①プロセッサAのラインに**書き込み**
- ②他のプロセッサのキャッシュ中のラインは、**更新**するか、**無効化**する。



書き込み-ライトバック②共有記憶更新

- ①プロセッサCがキャッシュにないラインを**共有記憶**に**アクセス**。その共有記憶ラインが**未更新**のとき、
- ②最新のラインをもつプロセッサAのキャッシュから**コピー**
- ③プロセッサAのキャッシュラインで共有記憶ラインを**更新**



Intel Core i9-12900KS Q1'22 specifications

- Core i9, 16core, 3.4GHz, 128GB

L1 キャッシュ

命令キャッシュ 16x32KB, 8way set associative

データキャッシュ 16x48KB, 8way set associative

L2 キャッシュ 10x1.25MB, 8way set associative

L3 キャッシュ 30MB, 16way set associative

ラインの大きさ 64B