

計算機方式論

第13章 キャッシュ

-性能の指標-

1

キャッシュ性能の指標

- ①ヒット率(hit ratio)
- ②アクセス時間
- ③ミスペナルティ時間
- ④主記憶更新時間
- ⑤キャッシュ容量
- ⑥ラインサイズ

2

①キャッシュ性能の指標-ヒット率 p (hit ratio)

- CPUが要求するアクセス対象(番地)のキャッシュ中に存在する確率。
- 高いほど良い。
- 対象への実効アクセス時間(平均アクセス時間) T_e
ヒット率 p
キャッシュへのアクセス時間 T_c
主記憶にアクセスしたとき要する時間 T_{miss} としたとき、
 $T_e = pT_c + (1-p)T_{miss}$

3

ヒット率 p と実効アクセス時間 T_e の例

- $T_e = pT_c + (1-p)T_{miss}$
 $T_c = 5\text{ns}$ 、 $T_{miss} = 31T_c = 155\text{ns}$ とすると、
 $T_e = (31-30p)T_c$

$p = 100\%$	$T_e = T_c = 5\text{ns}$ キャッシュへのアクセス時間(最良)
$p = 90\%$	$T_e = 4 T_c = 20\text{ns}$
$p = 50\%$	$T_e = 16 T_c = 80\text{ns}$
$p = 0\%$	$T_e = T_{miss} = 31T_c = 155\text{ns}$ (最悪)

4

② キャッシュ性能の指標-アクセス時間 T_c

- T_c = キャッシュメモリ **自身** のアクセス時間
+ キャッシュラインへの **マッピング時間**
- キャッシュメモリ **自身** のアクセス時間は、構成する記憶素子で決まる…数nsの **SRAM** が使われる。
- **マッピング時間** とは、アクセス対象の主記憶番地から、キャッシュラインを決定する時間。連想検索方式のマッピング(アソシアティブ方式)のときは、キャッシュライン数が増えると、マッピング時間が増える。
- ヒット率 p が高い程、 T_c は実効アクセス時間 T_e に影響する。

$$T_e = pT_c + (1-p)T_{miss}$$

5

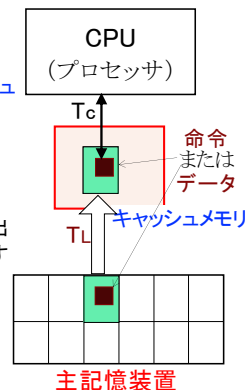
③ キャッシュ性能の指標-ミスペナルティ時間 T_{miss}

- ミスヒット時の処理時間で、主として、**ライン置換時間**。
- **ライン置換アルゴリズム**。
- 割り込み処理性能。
- 主記憶-キャッシュ間のデータ転送速度に依存。

6

ミスヒット時の実効アクセス時間 T_e

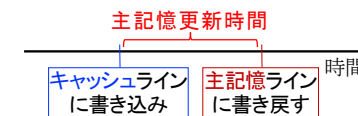
- 対象が **ミスヒット** 時、
 - **空のキャッシュライン** が在れば、**主記憶ライン** を時間 T_L でコピーし、**キャッシュライン** 中の対象を時間 T_c でアクセス。
 $T_{miss} = T_L + T_c$
 - **空のキャッシュライン** がなければ、アソシアティブマッピング方式では、**ライン置換アルゴリズム** が、時間 T_R で、**追出すキャッシュライン** を選び主記憶に書き戻すことで **空きキャッシュライン** をつくって、そこに **主記憶ライン** を時間 T_L でコピーし、対象を時間 T_c でアクセス。
 $T_{miss} = T_R + T_L + T_c$
- $T_{miss} = T_R + T_L + T_c$ とすると、
 $T_e = pT_c + (1-p)T_{miss} = T_c + (1-p)(T_L + T_R)$



7

④ キャッシュ性能の指標-主記憶更新時間

- **更新されたキャッシュライン** を主記憶ラインに **書き戻す** までの時間。
- 短い程有効
- **書き込みアクセス** が多いときほど、主記憶を更新することも多くなる。
- 主記憶の素子、ラインサイズに依存。



8

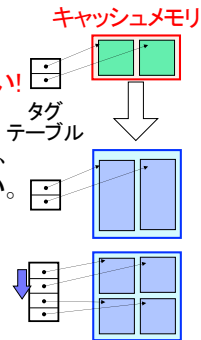
⑤ キャッシュ性能の指標-キャッシュ容量

- キャッシュ容量を増やしたとき、ヒット率 p は増えるが、実効アクセス時間 T_e が減るとは限らない!!

- (1) ラインサイズを大きくすると、主記憶-キャッシュ間転送時間 T_L が増え、実効アクセス時間 T_e が減るとは限らない。

$$T_e = T_c + (1-p)(T_L + T_R)$$

- (2) ラインサイズを同じか小さくすると、ライン数が増え、連想検索を行うマッピング方式では、タグ検索を行うため、 T_c が増え、実効アクセス時間 T_e が減るとは限らない。

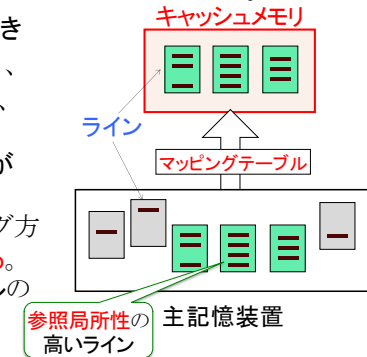


9

⑥ キャッシュ性能の指標-ラインサイズ

- ラインが主記憶-キャッシュ間の転送単位のため、そのラインの大きさでつぎのトレードオフがある。

- ラインサイズが大きいとき (キャッシュ容量は同じ)、
 - (1) 参照局所性が高ければ、ヒット率が高い。
 - (2) キャッシュ内のライン数が少なくなり、連想検索を行うマッピング方式では、検索時間が減る。また、マッピングテーブルの構成コストも安くなる。

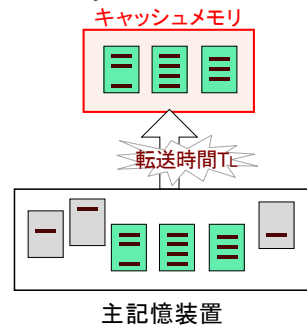


10

⑥ キャッシュ性能の指標-ラインサイズ

- ラインが主記憶-キャッシュ間の転送単位のため、その大きさでつぎのトレードオフがある。

- ラインサイズが大きいとき (キャッシュ容量は同じ)、
 - (3) 主記憶-キャッシュ間のデータ転送時間 T_L がかかる。

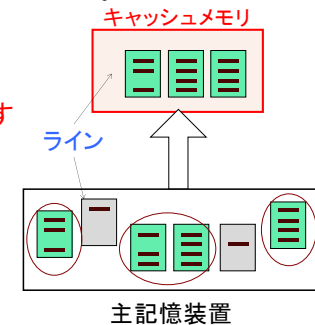


11

⑥ キャッシュ性能の指標-ラインサイズ

- ラインが主記憶-キャッシュ間の転送単位のため、その大きさでつぎのトレードオフがある。

- ラインサイズが大きいとき (キャッシュ容量は同じ)、
 - (4) キャッシュ内ライン数が少なすぎると、参照局所性を活用できなくなり、ヒット率が低下する。



12

⑥キャッシュ性能の指標-ラインサイズ(まとめ)

- ラインが主記憶-キャッシュ間の転送単位のため、ラインの大きさでつぎのトレードオフがある。
- ラインサイズが**大きい**とき(キャッシュ容量一定)は、
(1)参照局所性が高ければ、**ヒット率 p** は高し $T_e = T_c + (1-p)(T_L + T_R)$
(2)キャッシュ内のライン数が少なくなり、マッピングテーブルの**構成コスト**が安くなる。連想検索を行うマッピング方式では、**検索時間**(キャッシュアクセス時間 T_c)が減る。一方、
(3)主記憶-キャッシュ間の**データ転送時間 T_L** がかかる。
(4)キャッシュ内ライン数が少なすぎると、参照局所性を活用できなくなり、**ヒット率 p** が低下する。
- ラインサイズが**小さい**ときは、この長短が逆になる。
- マイクロコンピュータの**ラインサイズ**は、**32~128B**である。
- **キャッシュ容量**は、**32K~256KB**。

13

Intel Core i7-4770 Q2'13 specifications

- Core i7, 4core, 3.4GHz, 32GB

L1キャッシュ

命令キャッシュ**4x32KB**, 8way set associative

データキャッシュ**4x32KB**, 8way set associative

L2キャッシュ **4x256KB**, 8way set associative

L3キャッシュ **8MB**, 16way set associative

ラインの大きさ **64B**

14

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**

15