

Hiérarchie mémoire : structure et applications

problème : temps de fonctionnement
 répertoire $\leq 1 \mu\text{s}$ | RAM $\approx 500 \text{ cycles}$
 comment rendre la mémoire plus réactive ?

idée : dupliquer partiellement la mémoire dans un cache plus petit, proche et réactif. Ainsi le contenu qui y est stocké peut être accédé et modifié beaucoup plus vite.
 (localité temporelle : on manipule régulièrement les mêmes données)

I structure du cache

1) structure en bloc.

- Localité spatiale : quand on accède à une adresse on accède aux adresses contiguës.

↳ le cache charge et évacue le contenu de la mémoire par blocs

2) structure en onion

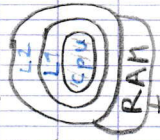
- Implanter plusieurs caches de taille croissante, et de réactivité décroissante traversés les uns après les autres lors de l'accès et de l'écriture dans la mémoire.

3) transparence.

- Le processeur accède et écrit dans la mémoire et ne sait pas si ses requêtes sont interceptées par un cache.

- Nécessité : le fonctionnement du cache est caractérisé par le matériel.

4) Métablanc



5) Performance

AMAT = Temps de recherche dans le cache + Temps d'échec \times (temps moyen en cas d'échec)

But : AMAT \ll Temps moyen d'accès sans cache.

II choix d'implémentation :

1) associativité



- fully associative : n'importe quel bloc peut être mis n'importe où dans le cache. Très dur à implémenter physiquement.

- directly mapped : un bloc ne peut être placé qu'à une seule position dans le cache dépendant de son adresse. Les adresses ne sont pas bien réparties et donc le même bloc du cache est utilisé en boucle.



- Set associative : pour chaque bloc de la mémoire il y a un ensemble de positions dans le cache où il peut être placé.



2) types d'adresses

On peut associer les blocs à leurs adresses physiques ou bien virtuelles. Dans le 1er cas, tester la présence du bloc dans le cache nécessite un appel à la MMU. Dans le second, le cache doit être purgé à chaque commutation.

3) Politique d'écriture et d'évacuation

a) écriture

- Write-through : chaque écriture est propagée en mémoire.
- Write-back : les écritures sont localisées au cache. Elles ne sont propagées que lors de l'évacuation de la ligne.

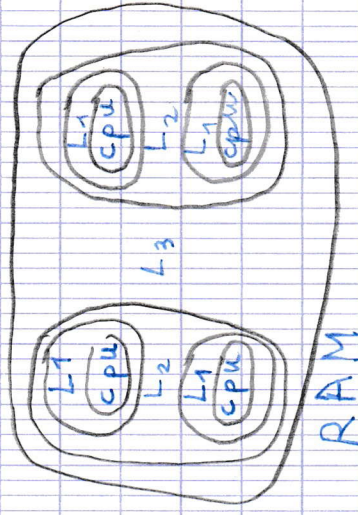
b) allocation:

- Write allcator: Lors d'un eche d'écriture, la ligne est chargée.
- Write non-allcator: Lors d'un eche d'écriture, la ligne n'est pas chargée.

4) multi² coeur.

a) localité du cache.

- un cache par coeur: Plus cher, des caches plus petits.
- un cache par puce: Plus loin, moins reactif.



b) Coherence.

- Les données dans un cache peuvent différer des données dans un autre cache ou dans la mémoire (charger une donnée qu'un autre coeur modifie dans la mémoire / écrire dans son cache mais pas dans la mémoire).

- Des²: Protocole MESI.

III Problematiques liées

1) Optimisations

a) séparer le cache données / cache instruction

- Les données et les instructions fonctionnent très différemment lors de l'exécution d'un programme. Ainsi on repare le cache en 2, un cache de données et un cache d'instruction utilisent des stratégies de caches très différentes.

b) préchargement

- On essaye de deviner le prochain bloc utilisé et on le charge à l'avance. Favoriser la localité spatiale au prix de la localité temporelle: préchargement. Favoriser la régularité d'une exécution: préchargement à pas.

c) Victim cache.

- Coller au cache L1 un cache associatif, le victim cache qui ramène les blocs évictés du cache L1. Perdre ainsi rapidement le cache L1. Ne ralentit pas les accès du cache L1. Diminue le problème lié au mauvais repartitionnement des adresses.

d) optimiser le code pour le cache.

- lors de la compilation, on optimise le code obtenu pour qu'il tienne au mieux avantage du cache. Ex: permutation de boucle.

2) Sécurité Dev 2: attaque Spectre.

3) Autres formes de cache:

a) TLB

- ici on veut gagner du temps en évitant un calcul pour la MMU.

b) SSD cache

- utiliser un SSD (plus rapide) comme cache pour une base de données dur.

c) SQL buffer cache

- stocker le résultat de requêtes SQL clés pour éviter leur calcul et simplifier les transactions.