

LRU: Partition on demand basis

Motivation: Higher demand for cache resources does not always correlate with a higher performance from additional cache resources.

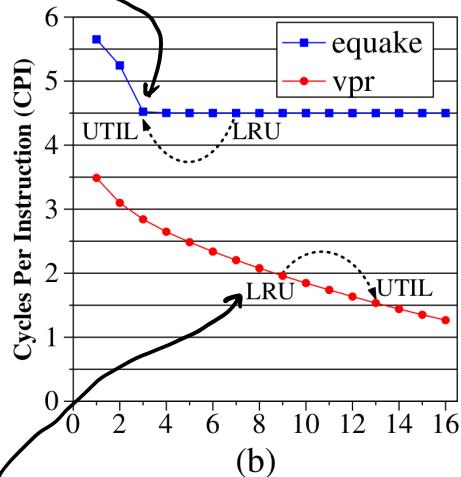
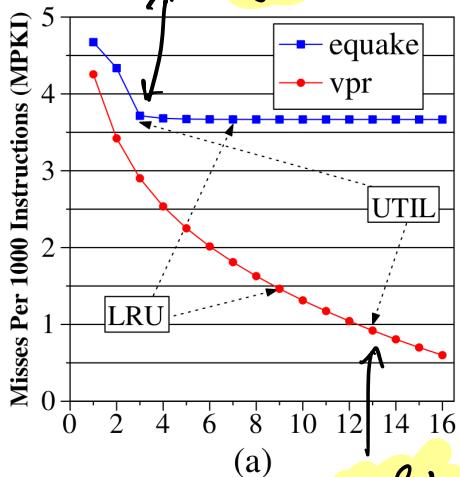
(Eg) A scanning will have high demand but there will be no reuse, hence it is not useful to allot it in cache.

↳ VCP (Utility-based Cache Partitioning)

↳ Low-overhead runtime

↳ Based on misses from each application given same demand the need.

Equake doesn't benefit beyond certain point



What we want to do?

→ Partition Cache based on utility (UTIL)

↳ Utility-based Cache Partitioning (VCP)

↳ UMON (Utility Monitoring)

↳ Lookahead

Makes sense to give more cache to app that will reuse data.

Types of utility:

→ Low utility:

- * Won't benefit from increasing the cache.
- * But will help if increased exponentially
- * Scanning / Thrashing

→ High utility

- * Benefits from increasing the cache

→ Saturating utility

- * small workload that fits in small cache. → won't benefit from increasing cache.

Combination of workloads:

Workload 1	Workload 2	
Low util	Low util	No issues
Saturating util	Saturating util	No issues
Saturating util	Low util	Cache may not accommodate working set of saturating util.
High util	Saturating util	Issue. Highly sensitive to cache size.
High util	Low util	

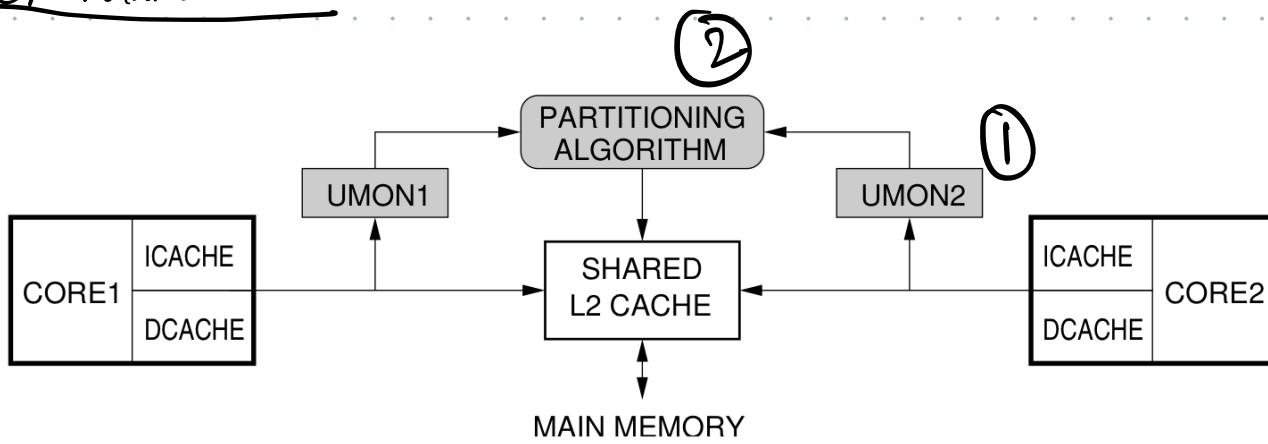
→ Reduction in MPKI \propto Reduction in CPI

Defining utility:

$$\rightarrow U_a^L = \underbrace{\text{miss}_a - \text{miss}_L}_{\text{as de}}$$

miss occurred
when a ways
are allocated

VCP Framework:



Totally there are 3 components:

1) **UMON**: (Utility Monitor)

→ collects util information

2) **Partitioning Algorithm**:

→ Partitions the ways based on info from UMON

3) **Replacement Engine**

→ Augmented to support the partitions made by partitioning algorithm!

① **UMON**:

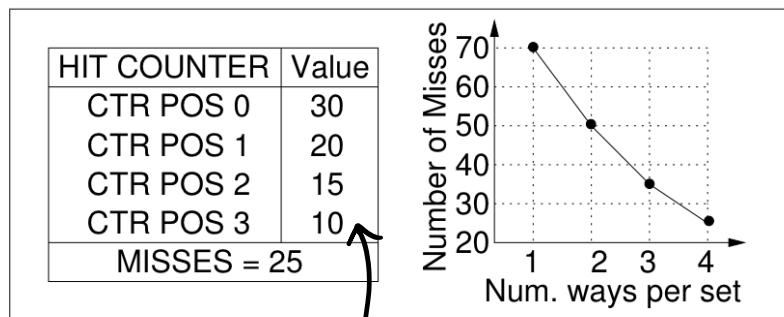
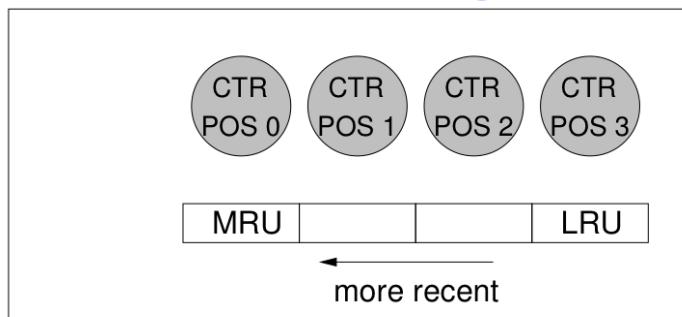
→ We only need tag directories

→ Naive but expensive way: check misses for 1–16 ways.

↳ Instead exploit streak property of LRU.

→ Maintain a counter for each way. Whenever there is a hit in a way, increment the corresponding counter.

These counter values indicate the number of misses if this way didn't exist.

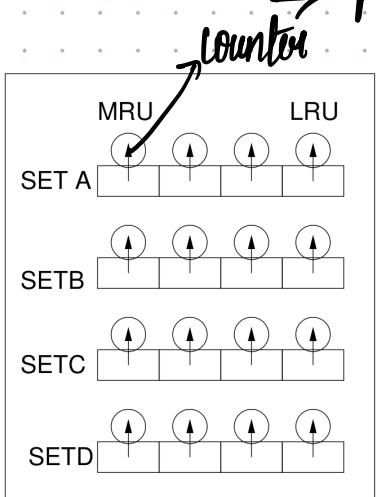


To get misses for way=2, add 10+15. Take into account the cumsum.

Therefore, in order to obtain misses for all ways, we just have to do it only for 16-ways, with only tag dirs.

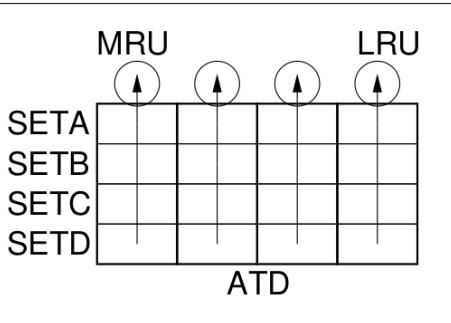
How does the circuit of UMON work?

- Auxiliary Tag Directory (ATD): Has same associativity as the main tag array
- Hit counters



* UMON - Local

- Has counter for each block
- Total "sets x ways" counter
- A lot of overhead.



* UMON-global

- Has one counter per way.
- Total "ways" counter
- comparatively less number of counters

Fissile with UMON-global:

- has the extra overhead of maintaining ATD.
- 1 UMON per core.
- Eg.)

Tag = 4 bytes.

We need + (4 x cores) overhead!

Solution: UMON-DSS

= UMON-global + Dynamic Set Sampling (DSS).

→ DSS: The Cache behaviour can be approximated with only a few sets.

→ This way we don't have to maintain a full ATD, instead just maintain the ATD for the subsets.

② Partitioning algorithm:

$$U_{\text{tot}} = UA_i^u + VB_i^{(16-u)} \quad \forall i \in (1, 16-1)$$

- UA_i^u will give the utility info (#misses) based on the UMON-DSS when the partition is done at i^{th} way.
- The i with the maximum U_{tot} will be chosen as the partition.

? After certain interval, the hit counter in all the UMON are halved? What does this mean?

(3) Replacement engine

- 1 extra bit is append to TAG indicating the core which the block belongs to based on the partitioning algorithm
- On conflict miss, count the blocks which belong to the missed block's core.
 - * If its less than allotted, evict a CRU block from the other core blocks
 - * If its greater, then evict a CRU block from the current core block.

scalability:

partitioning techniques are same

→ The above mentioned partitioning examples work well when there are only 2 applications but gets difficult to scale.

↳ Greedy Algorithm

↳ Look Ahead Algorithm.

Here's how the utility (benefit) for both applications changes as they get more cache blocks:

Cache Blocks	App A (Miss Reduction)	App B (Miss Reduction)
1	1	2
2	3	3
3	6	5
4	7	6
5	8	8
6	9	9

This table shows that App A and App B experience different improvements in miss reduction as they receive more cache blocks.

Greedy Algorithm Example:

1. **Initial Setup:** Both App A and App B start with 0 cache blocks.
2. **First Allocation:** The algorithm looks at the benefit of giving 1 block to either App A or App B:
 - App A gets a miss reduction of 1.
 - App B gets a miss reduction of 2.
 - The algorithm allocates the first block to App B since it has a higher immediate benefit.
3. **Second Allocation:** Now, App B has 1 block. The algorithm looks again:
 - App A gets a miss reduction of 1 with 1 block.
 - App B gets an additional miss reduction of 1 (total reduction is now 3).
 - The algorithm still favors App B because both have the same benefit, but App B is chosen arbitrarily in the case of a tie.
4. **Third Allocation:** The process repeats, and App B keeps getting blocks as long as it shows equal or higher benefit than App A.

With this greedy approach, the algorithm makes short-term decisions based on immediate gains without looking ahead at how giving blocks to App A might eventually yield a higher total reduction.

Lookahead Algorithm Example:

The lookahead algorithm, on the other hand, takes into account **future utility** for the next few blocks, not just the immediate benefit.

1. **First Allocation:** The lookahead algorithm checks what happens if it gives more blocks to either App A or App B:
 - If App A gets the next block, its total reduction will increase to 3 if it receives 2 blocks.
 - If App B gets the next block, its reduction will increase to 5 if it receives 3 blocks.
 - The lookahead algorithm evaluates the cumulative benefit (over several blocks) and sees that while App B still has better immediate benefits, App A will catch up after receiving 3 blocks. So, it might give one or two blocks to App A to balance the longer-term gains.
2. **Second Allocation:** After taking future gains into account, it might allocate more blocks to App A earlier than the greedy algorithm would, allowing App A to gain more total cache miss reduction over time.

Mind map:

High demand
≠
Performance improve
with bigger cache
(Scan)

