

DENSITY

- WCAB: Write congestion aware bypass
- VHC: Virtual hybrid cache

SRAM:

- Low density
- High leakage

solution: NVM → STT RAM (Spin-Transfer-Torque)
→ PRAM (Resistive RAM)

Issues:

- Long write latency
- This latency puts pressure on the queue \Rightarrow congestion

Issues with current solutions:

- Reduces LLC density
- Increased error rate \rightarrow unreliable
- Lowering write latency reduces the density advantages

Writes to LLC in inclusive Vs exclusive:

- Inclusive:
- 1) LLC miss
 - 2) Dirty victims from L2

- Exclusive:
- 1) Hit in LLC \rightarrow As they are made I and pushed to updir
 - 2) Clean and dirty victims from L2.



writes in exclusive \gg writes in inclusive. But exclusion increase capacity.

WCAB: Write Congestion Aware Bypass:

→ Aim of existing bypass \Rightarrow Improve bit rate

→ As LLC capacity increases \Rightarrow The fraction of writes that won't be used, increase drastically.

Solution: More aggressive bypass



But this would decrease bit rate.

Tradeoff

Goal

Dynamic params

Request Queue:

$$(\text{Occupancy of request queue} \text{ is high}) + (\text{fraction of writes is above threshold}) = \text{congestion}$$

Counters

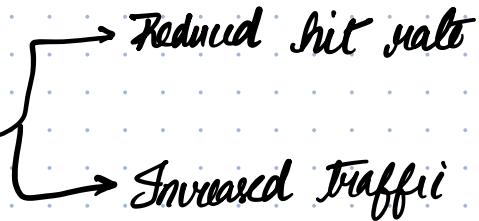
More congestion \rightarrow More WCAB bypass

Less congestion \rightarrow Less WCAB bypass

Running average \rightarrow Small spikes are not treated as congestion

Liveness:

→ Simple bypass



- Allocate observer sets in L2 & LLC \Rightarrow To track reuse behavior
- Hashed L2 access PC (Program Counter): Instruction address that last accessed the cache line in L2.

\hookrightarrow 5 buckets (0 - 100 %) increments of 20.

$\brace{ \text{liveness} }_{\text{buckets}}$

Liveness = How much of the evicted sets are recalled again

(Eg.) 20% = 20 sets were recalled out of the 100 evicted sets.

} only done for observer sets



for 21% - 50% \Rightarrow $\begin{cases} \text{Eviction} \rightarrow \text{Decrement by } 2 \\ \text{Recall} \rightarrow \text{Increment by } 10 \end{cases}$

WCAB: Write Congestion Aware Bypass

1) Detect congestion:

→ Monitor LLC request queue.

→ If the no. of write requests is above a fixed threshold, it indicates severe congestion

→ WCAB activates severe bypassing

2) Liveness score:

→ Don't bypass block with high liveness score

→ Observer list:

- Track how often these blocks are accessed after being written

→ Liveness score:

- Score is divided into bucket \Rightarrow 4 buckets

3) Bypass decision:

→ When congested is detected, writes with low liveness score are bypassed

4) This bypass is adapted overtime based on congestion to modulate the aggressiveness.

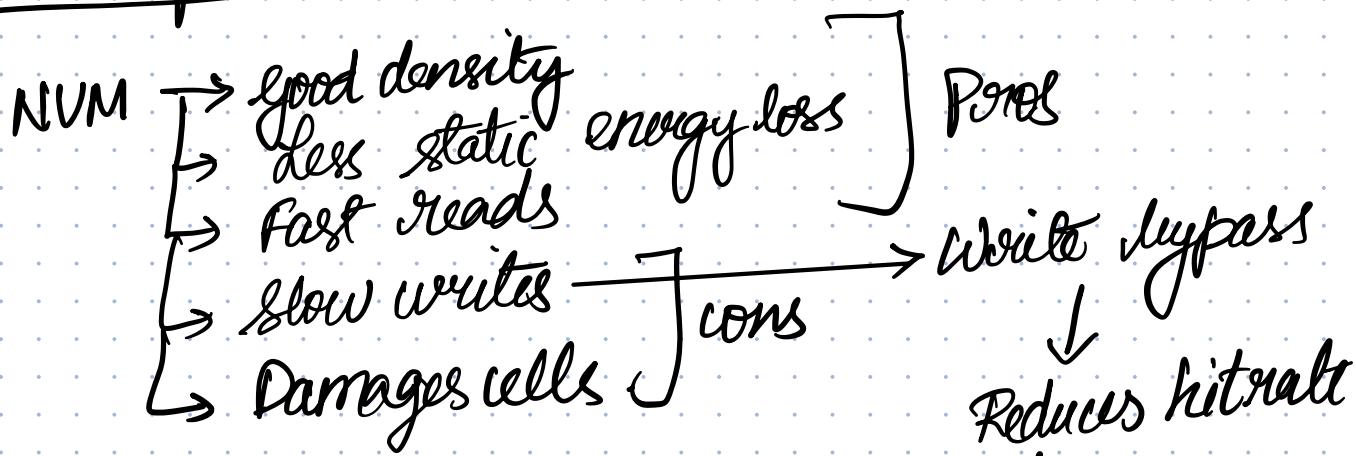
Virtual Hybrid Cache (VHC):

1) Frequently dirty fills: The writeback is not done to LLC everytime. Instead it is stored longer in L2 and these changes are merged & then written back to LLC.

2) Frequently clean fills: A copy of frequently clean

Blocks are stored both in L2E LLC to prevent unnecessary writebacks -

Mindmap:



DENSITY → Intelligent bypass

WCAB → Congestion detection → LLC request queue is full
+
·. Writes > threshold

WCAB → Liveness score → Observation sets

↓
Count access after first write
↓
4 buckets

↳ Bypass → If congestion \rightarrow bypasses
in queue
with low
liveliness
score

VTC \rightarrow Frequent dirty fill \rightarrow wait longer \rightarrow Merge
in L2 changes

↳ Frequent clean fill \rightarrow Have duplicates
in L2 & LLC.