# Paper Review Notes

| | Static RRIP | Dynamic RRIP | uses 2 bit per Cache block |
|---|---|---|---|
| Single Core | 4 % | 10 % | |
| Multi Core | 7 % | 9 % | |

---

LRU Poor perf

↳ less Cache Size

↳ Non temporal data replaces the active working set of application

## Access Patterns

$$[a_1 - - - a_k \; a_k \; a_{k-1} - - a_1]$$

Recent friendly ( for any k)

$$(a_1 - - - - a_k)^N \quad (\text{Thrashing Access pattern})$$
k > cache size

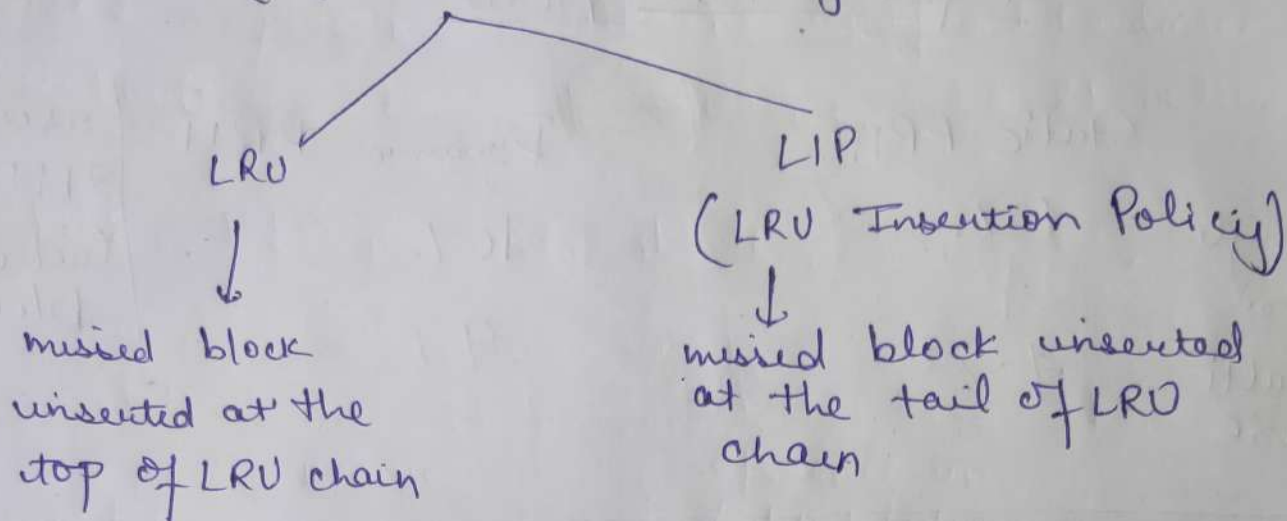$$(a_1 - - - a_k) \quad (\text{Streaming Access pattern})$$
$$k = \infty$$

$$[(a_1 - - - a_k, a_k - - - a_1)^A \; P_\varepsilon (a_1 - - - a_k \; a_{k+1} - - - a_m]^N$$

$$[(a_1 - - - a_k)^A \; P_\varepsilon (b_1 b_2, b_m)]^N$$

mixed Access pattern

# DIP (Dynamic Insertion Policy)

```
                    DIP (Dynamic Insertion Policy)
                   /                        \
              LRU                            LIP
               |                      (LRU Insertion Policy)
               ↓                              ↓
          missed block                 missed block inserted
          inserted at the              at the tail of LRU
          top of LRU chain             chain
```

o) DIP uses LIP when working set is larger than cache, LRU otherwise.

o) Both LRU & DIP performs poor in case of mixed access pattern.

o) Scan: Burst of references of data whose re-ref interval is in the distant future.

o) Real world application frequently suffer from scan. To improve cache performance we need scan resistant cache replacement policy

We have 2 cases

o) 
```
[ Data  ]  →  [ Scan ]  ———→  [ Data ref ]
[  ref  ]                     [   D2    ]
   D1
```

o) 
```
[ Data ]  →  [ Scan ]  ———→  [ Data  ]
[  ref  ]                     [  ref   ]
   D1                            D1
```

In case 1 : Decision during Scan did wt matter

But in case 2 it matter.

If we try to apply LIP during Scan & LRU rest. This will give good performance

## Related Work

**LFU** → degrades performance in recency

**Self tuning adaptive Policy** : Significantly inc. the hardware overhead & complexity

**Hybrid Cache Replacement** :- Uses set duelling to dynamically choose b/w multiple replacement policies. It provide scan resistant but require hardware & verification overhead of 2 caches

**Dead block prediction** : predicts which of the block will be dead from RRIP chain. Victum selection policy select dead block closer to tail. It improve cache performance but require additional hardware overhead for dead block predictor.

**Pseudo LIFO** :

**RRIP :** high performing practical scan resistant Cache replacement with no significant hardware overheads or changes to existing cache structure.

**NRU :** 1 bit per cache block.



bit = 0                           bit = 1
means reference in          reference in
near future                     distant future

∞

If all bit are 0, change all bits to 1.

**Problems :**

↳ Always predict re-ref. in near future for incoming block

↳ If predict distant then do poor on near-immediate re-ref. Interval

**Static RRIP**

↳ granularity of re-ref pred is inc.

M bit ⟶ $2^M$ bit granularity

0 ⟶ near immediate

$2^M - 1$ ⟶ distant future

$\hookrightarrow$ Always insert at long re-reference interval

$\hookrightarrow$ Author used $2^M - 2$ as l r ref I.

[prevent scan from polluting the Cache]

[& pred giving initial $2^M - 2$ instead of
$2^M - 1$ gives Cache to learn some turns]

## Victim Selection Policy

$\hookrightarrow$ Select with $RRPV = 2^M - 1$

$\hookrightarrow$ Ties broken by starting victim search by from fixed location

$\hookrightarrow$ If no block found increment RRPV of all block by 1.

## Hit promotion Policy

## Hit Priority

$\hookrightarrow$ on hit change $RRPV = 0$

(Prioritize cache block that receive hit over not receiving hits)

$\hookrightarrow$ degrade when CB only hit onces

## Freq Priority

$\hookrightarrow$ on hit decrease RRPV by 1

$\hookrightarrow$ goal is to prioritize replacement of infrequently re-ref each.

SRRIP Scan resistant length

$$S_{err} \quad S_{err} \propto (2^M - 1)(A - \omega)$$

Can be increased by increasing M.

This can also cause inefficient Cache Utilization
When a cache block receive its last hit and
RRPV became Zero

## Dynamic RRIP

↳ SRRIP is vulnerable to thrashing

↳ To Avoid thrashing we can use

Bimodal RRIP

Insert Majority                    few with
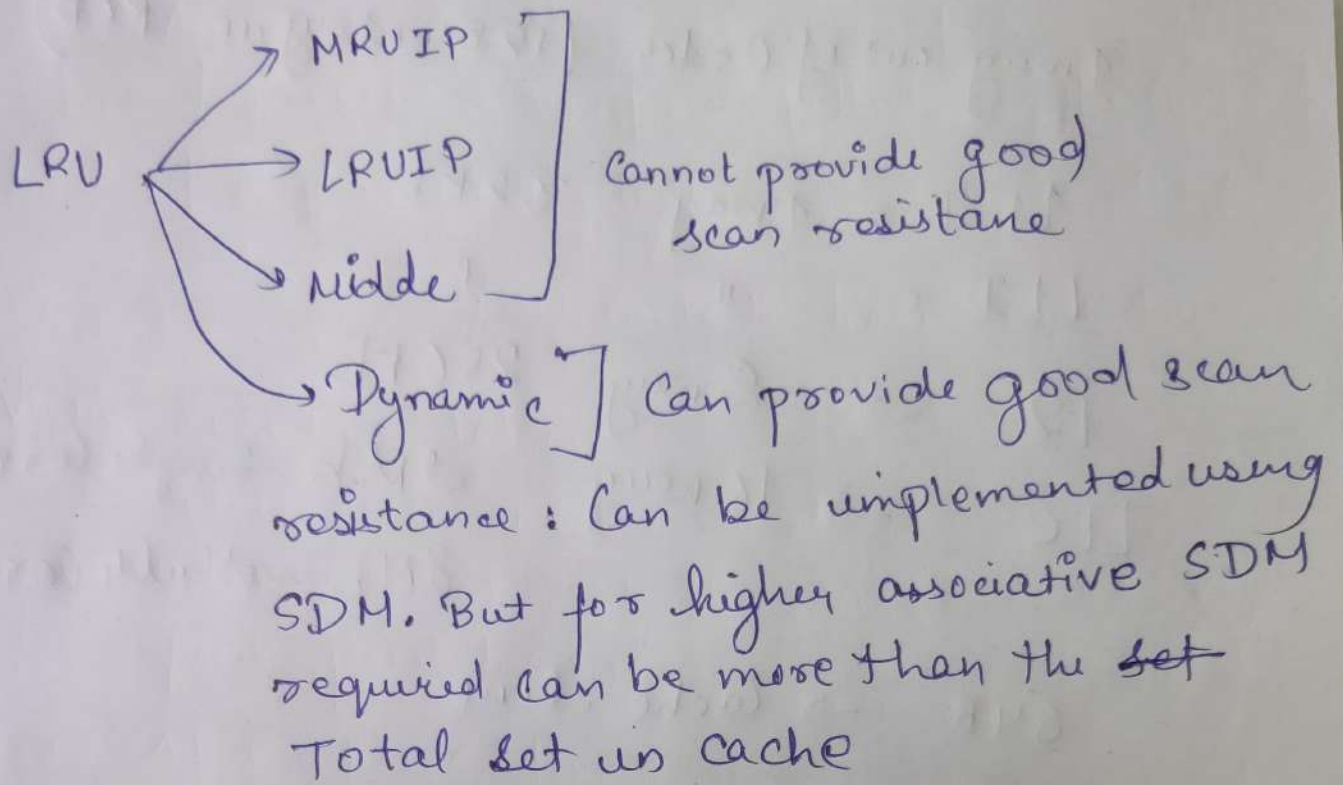CB with distant                   long re-ref interval
reref interval

( for non thrashing it can lead to
Performance degrade )

We choose dynamically between SRRIP & BRRIP
using set dueling

# SRRIP Vs LRU

LRU
→ MRUIP
→ LRUIP      } Cannot provide good
→ Middle        scan resistance

→ Dynamic ] Can provide good scan
resistance : Can be implemented using
SDM. But for higher associative SDM
required can be more than the ~~set~~
Total set in cache

## Extending RRIP to Shared Caches.

↳ ~~Same~~ We can see access as mixed pattern

↳ Extend DRRIP to Shared caches is
same as extending DIP to Shared
Caches.

↳ TADRRIP → uses 2 SDM per Application
to determines its optimal policy.

# Experimental

4 way out of Order (128 entry in ROB)

L1 I = 4 way   32kb
L1 D =   "       "

L2 → 8 way   256 kb
LLC → 16 way   2 Mb   Single Core
                8 Mb   Multi Core

64 B → Cache line

only demand ref change LRU state

Load latencies →

L1 → 1
L2 = 10
L3 = 24
Main Mem = 250

Support 32 outstanding musses to memory

Benchmarks used —

↳ vulnerable to replacement decision

## Results and Analysis

→ Bi taking m=2,3 performs best.

→ Best Insertion position for every M is $2^M - 2$

→ SRRIP-HP performs better than SRRIP-FP which seems little counter intuitive

→ SRRIP performs good on cache size (LLC) range from 512Kb to 8MB with optimal m = 2,3

→ Can Easily be extended to shared caches due to its good performance on mixed pattern.

→ Perform good even changing cache level to 2 (L₁ & LLC)

→ performs 4% & 10% better than LRU on single Core

→ perform 7% & 9% better on multi-Core than LRU.