

# Managing Distributed, Shared L2 Caches through OS-Level Page Allocation

January 21, 2020

In the last episode...

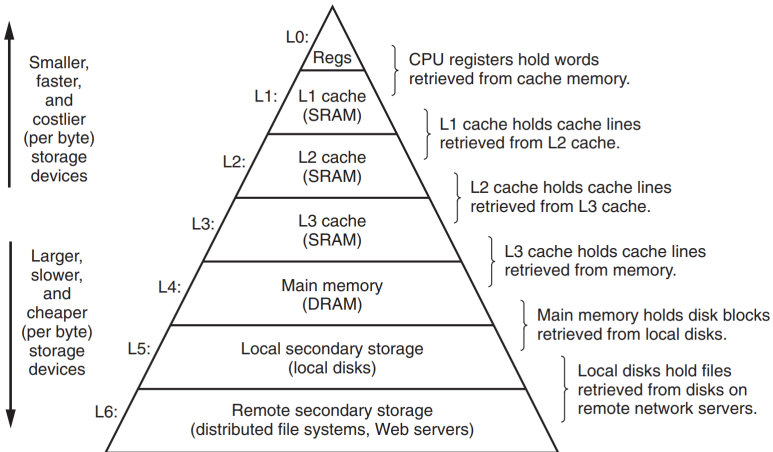
**Problem:** CPUs are cheap and fast.

Memories are cheap, fast, capacious (choose two out of three).

Fast CPU with slow memory doesn't make sense.

**Idea:** Put fast memory between CPU and main memory.  
Fast – so small.

CPU looks for data in the cache; if it finds it, it gets it from there and continues to work; if it doesn't find it, it looks for it in main memory.



memory	latency	size
L1	4 cycles	32KiB
L2	10 cycles	256KiB
L3	40-75 cycles	8MiB
DRAM	600 cycles	8GiB
HDD	$\infty$	$\infty$

INTEL® 01 PHILIPPINES  
1133/256/133/1.475



7130B200-0273  
PENTIUM III SL56Q

It was working



It was working; but now...



How to live?

How to live?

Private cache?

How to live?

Private cache? Shared cache?

Private cache

## Private cache

Each cache slice is associated with a specific processor core and replicates data freely as the processor accesses them.

## Private cache

Each cache slice is associated with a specific processor core and replicates data freely as the processor accesses them.

Pros: low hit latency



## Private cache

Each cache slice is associated with a specific processor core and replicates data freely as the processor accesses them.

Pros: low hit latency

Cons: capacity misses

## Private cache

Each cache slice is associated with a specific processor core and replicates data freely as the processor accesses them.

Pros: low hit latency

Cons: capacity misses

## Shared cache

## Private cache

Each cache slice is associated with a specific processor core and replicates data freely as the processor accesses them.

Pros: low hit latency

Cons: capacity misses

## Shared cache

Single cache where each cache slice accepts only an exclusive subset of all memory blocks.

## Private cache

Each cache slice is associated with a specific processor core and replicates data freely as the processor accesses them.

Pros: low hit latency

Cons: capacity misses

## Shared cache

Single cache where each cache slice accepts only an exclusive subset of all memory blocks.

Pros: better overall utilization of on-chip caching capacity, enforcing cache coherency becomes simpler

## Private cache

Each cache slice is associated with a specific processor core and replicates data freely as the processor accesses them.

Pros: low hit latency

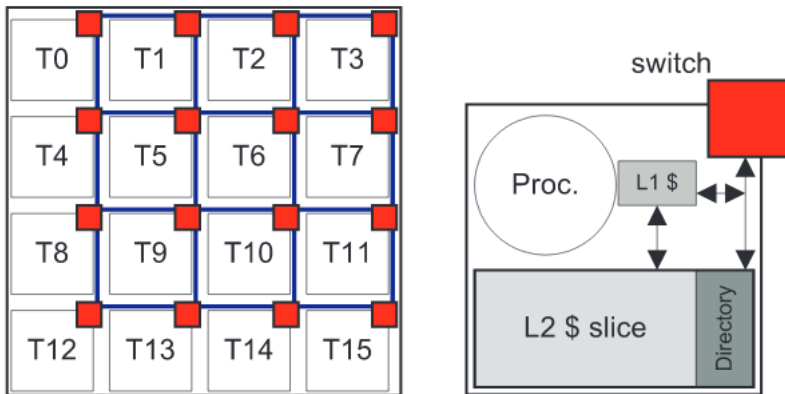
Cons: capacity misses

## Shared cache

Single cache where each cache slice accepts only an exclusive subset of all memory blocks.

Pros: better overall utilization of on-chip caching capacity, enforcing cache coherency becomes simpler

Cons: cache hit latency will be longer



**Figure 1. An example 16-core tiled processor chip and its tile (core).**

Idea: give OS control over cache.

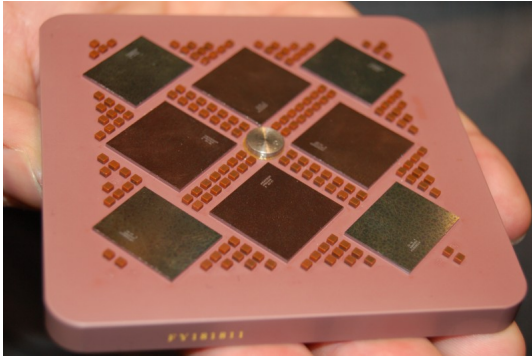
Idea: give OS control over cache.

Using simple shared cache hardware we can implement a private caching policy, a shared cache policy, or a hybrid of the two



Idea: give OS control over cache.

Using simple shared cache hardware we can implement a private caching policy, a shared cache policy, or a hybrid of the two **without any hardware support.**



IBM Power 5!

$$S = A \bmod N$$

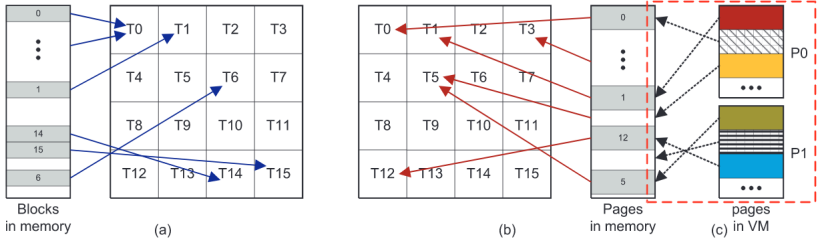
$$S = A \bmod N$$

Where  $S$  stands for the cache slice number,  $A$  for the memory block address, and  $N$  for the number of cache slices.

Problem: contiguous memory blocks are hosted by different cache slices

Problem: contiguous memory blocks are hosted by different cache slices

Solution: replacing  $A$  with physical page number (PPN)



**Figure 2. (a) Physical memory partitioning and mapping to cache slice at the cache line granularity. (b) Physical memory partitioning and mapping at the memory page granularity. (c) Virtual to physical page mapping (P0, P1: process 0 and 1, VM: virtual memory).**

$CG_i$  ( $0 < i < N - 1$ ) – congruence group

$CG_i = \{ \text{physical page (PPN} = j) \mid \text{pmap}(j) = i \}$

OS has control on pmap.



$CG_i$  ( $0 < i < N - 1$ ) – congruence group

$CG_i = \{ \text{physical page (PPN} = j) \mid \text{pmap}(j) = i \}$

OS has control on pmap.

Private caching: for a page requested by a program running on core  $i$ , allocate a free page from  $CG_i$

$CG_i$  ( $0 < i < N - 1$ ) – congruence group

$CG_i = \{ \text{physical page (PPN = } j) \mid \text{pmap}(j) = i \}$

OS has control on pmap.

Private caching: for a page requested by a program running on core  $i$ , allocate a free page from  $CG_i$

Shared caching: for a requested page, allocate a free page from all the congruence groups using random selection, round-robin etc...

$CG_i$  ( $0 < i < N - 1$ ) – congruence group

$CG_i = \{ \text{physical page (PPN} = j) \mid \text{pmap}(j) = i \}$

OS has control on pmap.

Private caching: for a page requested by a program running on core  $i$ , allocate a free page from  $CG_i$

Shared caching: for a requested page, allocate a free page from all the congruence groups using random selection, round-robin etc...

Hybrid caching: partition  $\{CG_i\}$  into  $K$  groups ( $K < N$ ); then define a mapping from a core to a group; for page requested by program running on core  $i$ , allocate a free page from the group that core  $i$  maps to

P	1	2	3
1	2	3	4
2	3	4	5
3	4	5	6

Case 1 (avg. distance = 3)

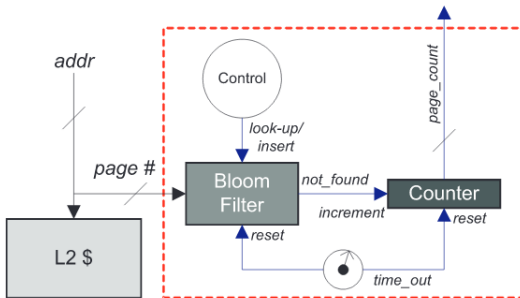
1	2	3	4
P	1	2	3
1	2	3	4
2	3	4	5

Case 2 (avg. distance = 2.5)

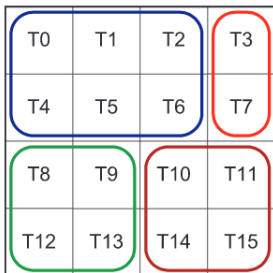
2	1	2	3
1	P	1	2
2	1	2	3
3	2	3	4

Case 3 (avg. distance = 2)

**Figure 3. Program (“P”) and data locations determine the minimum distance to bridge them.**



**Figure 4. A Bloom filter based monitor mechanism to count actively accessed pages.**



$G_0 = \{T_0, T_1, T_2, T_4, T_5, T_6\}$

$G_1 = \{T_3, T_7\}$

$G_2 = \{T_8, T_9, T_{12}, T_{13}\}$

$G_3 = \{T_{10}, T_{11}, T_{14}, T_{15}\}$

**Figure 5. A virtual multicore (VM) example.**

	PRIVATE	SHARED (w/ line interleaving)	OS-BASED
<i>Hardware (tile)</i>	Similar to a conventional uniprocessor core with two-level caches; coherence mechanism ( <i>e.g.</i> , directory) to cover L1 and L2 caches	Coherence enforcement is simpler and is mainly for L1 because L2 is shared by all cores	(Same as SHARED); simple hardware-based performance monitoring mechanism will help reduce monitoring overhead
<i>Software</i>	(NA)	(NA)	OS-level support, esp. in the page allocation algorithm
<i>Data Proximity</i>	Data items are attracted to local cache slices through active replication; limited caching space can result in performance degradation due to capacity misses	Fine-grained cache line interleaving results in non-optimal data distribution; there is no explicit control over data mapping	Judicious data mapping though page allocation can improve data proximity
<i>Network Traffic</i>	High coherence traffic ( <i>e.g.</i> , directory look-up and invalidation) due to data replication [23]; increased off-chip traffic due to high on-chip miss rate	High inter-tile data traffic due to remote L2 cache accesses, often $2\times$ to $10\times$ higher than PRIVATE [30]; lower off-chip traffic due to larger caching capacity	Low off-chip traffic like SHARED; improved program-data proximity through page allocation and process scheduling leads to lower inter-tile traffic than SHARED

**Table 1. Comparing private caching, shared caching, and OS-based cache management approaches.**

## Test setup

CPU: 16 tiles (4x4); 16kB L1 I/D caches, 512kB L2 cache slice

L1 caches are four-way set associative with a 32-byte line size

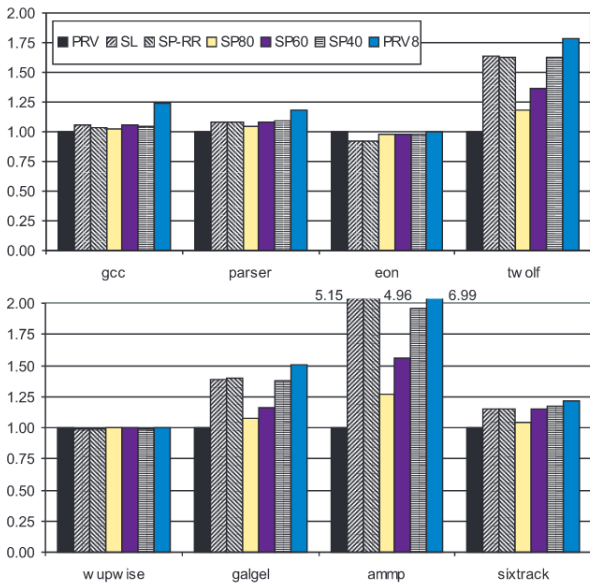
Each 8-cycle L2 cache slice is 8-way set associative with 128-byte lines, two-cycle latency per each hop

2-GB off-chip main memory with 300 cycles latency



NAME	DESCRIPTION	INPUT
<i>gcc</i>	gcc compiler	reference ( <i>integrate.i</i> )
<i>parser</i>	English parser	reference
<i>eon</i>	probabilistic ray tracer	reference ( <i>chair</i> )
<i>twolf</i>	place & route simulator	reference
<i>wupwise</i>	quantum chromodynamics solver	reference
<i>galgel</i>	computational fluid dynamics	reference
<i>ammp</i>	ODE solver for molecular dynamics	reference
<i>sixtrack</i>	particle tracking for accelerator design	reference
<i>fft</i>	fast Fourier transform	4M complex numbers
<i>lu</i>	dense matrix factorization	512×512 matrix
<i>radix</i>	parallel radix sort	3M integers
<i>ocean</i>	ocean simulator	258×258 grid

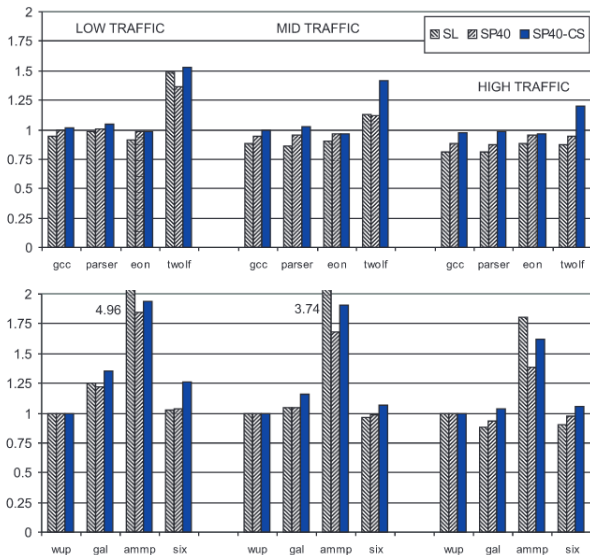
**Table 2. Benchmark programs.**



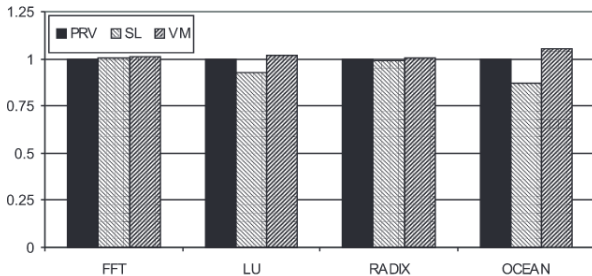
**Figure 6. Single program performance ( $time^{-1}$ ) of different policies, relative to *PRV*.**

		<i>PRV</i>	<i>SL</i>	<i>SP-RR</i>	<i>SP40</i>	<i>SP60</i>	<i>SP80</i>	<i>PRV8</i>
load miss rate (%)	<i>gcc</i>	2.9	0.1	0.5	1.8	2.1	2.8	0.1
	<i>parser</i>	6.6	0.5	0.6	2.6	3.7	5.8	0.4
	<i>eon</i>	0.0	0.0	0.0	0.0	0.0	0.0	0.0
	<i>twolf</i>	16.3	0.1	0.1	1.6	7.3	13.1	0.0
	<i>wupwise</i>	25.0	25.0	25.0	25.0	25.0	25.0	25.0
	<i>galgel</i>	6.3	0.1	0.1	0.9	3.4	5.0	0.1
	<i>ammp</i>	46.6	0.1	0.4	18.9	26.4	34.9	0.1
	<i>sixtrack</i>	13.5	0.5	0.5	1.4	3.2	10.4	0.5
on-chip network traffic	<i>gcc</i>	10.8	270.4	261.7	135.9	76.0	55.6	0.4
	<i>parser</i>	8.7	96.8	96.5	40.4	18.9	18.2	0.5
	<i>eon</i>	0.0	86.9	90.2	23.7	20.4	17.7	0.0
	<i>twolf</i>	35.0	138.2	150.1	67.8	48.4	37.8	0.1
	<i>wupwise</i>	35.1	39.4	39.9	20.3	15.6	10.1	0.1
	<i>galgel</i>	38.0	412.0	406.6	185.8	132.3	76.2	0.6
	<i>ammp</i>	441.7	810.9	803.4	424.6	361.9	306.9	0.5
	<i>sixtrack</i>	9.6	57.2	60.9	22.0	18.9	15.8	0.4

**Table 3. L2 cache load miss rate and on-chip network traffic (message-hops) per 1k instructions.**



**Figure 7. Performance sensitivity to network traffic. Performance relative to *PRV*, no traffic case.**



**Figure 8. Performance of parallel workloads, relative to *PRV*.**

(standing ovation)