

N

1 NUMA Caches

- 2 ALESSANDRO BARDINE, PIERFRANCESCO FOGLIA, COSIMO
3 ANTONIO PRETE, MARCO SOLINAS
4 Università di Pisa, Pisa, Italy

5 Definition

6 NUMA is the acronym for Non-Uniform Memory
7 Access. A NUMA cache is a cache memory in which
8 the access time is not uniform but depends on the posi-
9 tion of the involved block inside the cache. Among
10 NUMA caches, it possible to distinguish: (1) the NUCA
11 (Non-Uniform Cache Access) architectures, in which
12 the memory space is deeply sub-banked, and the access
13 latency depends on which sub-bank is accessed; and (2)
14 the shared and distributed cache of a tiled Chip Mul-
15 tiprocessor (CMP), in which the latency depends on
16 which cache slice has to be accessed.

17 Discussion

18 Introduction

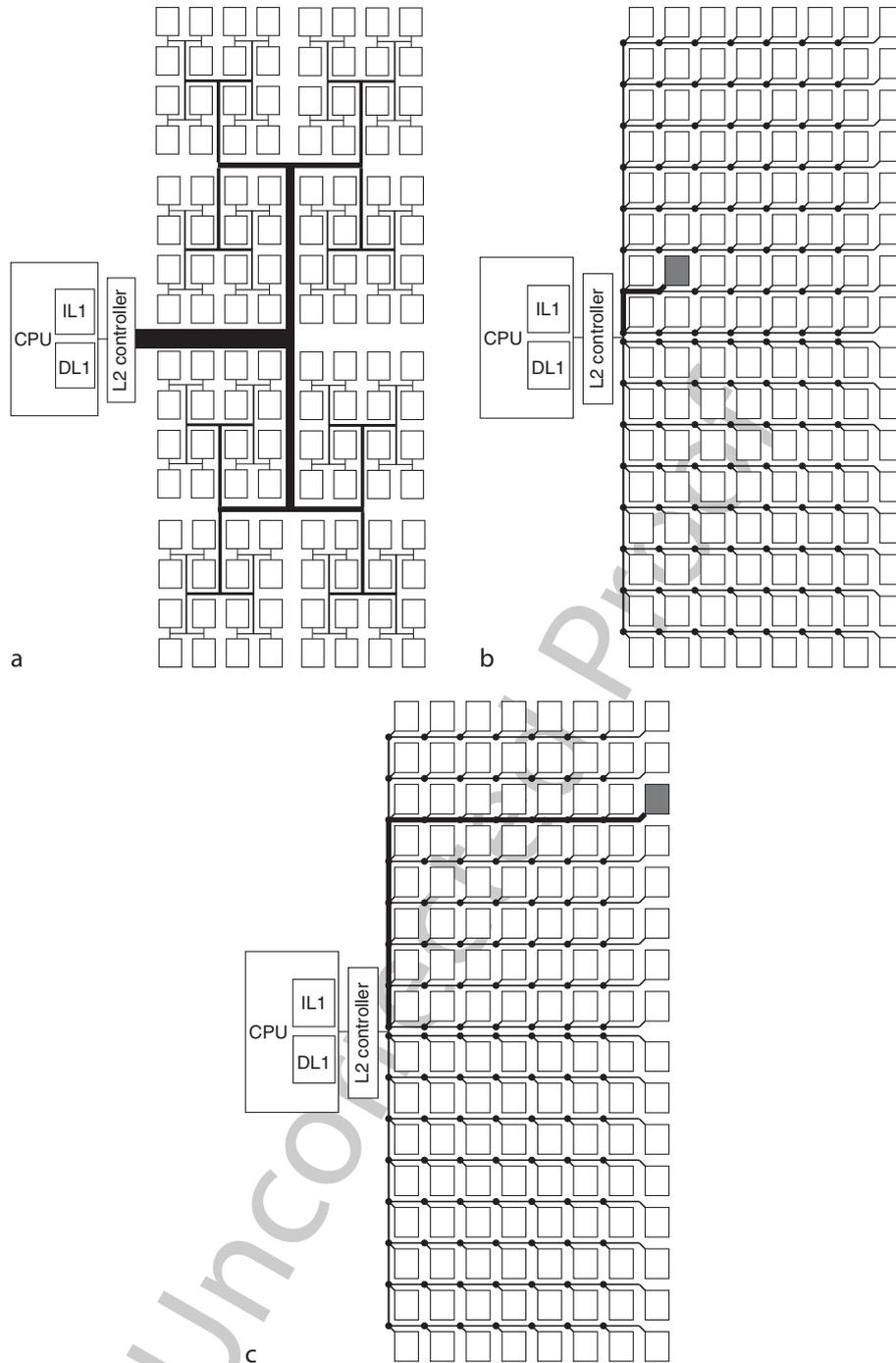
19 For the past decades, microprocessors’ overall perfor-
20 mance has been improved thanks to the continuous
21 reduction of transistor size obtained in silicon fabri-
22 cation technology. This scaling contributed in both (1)
23 allowing designers to put on a chip more and more tran-
24 sistors, and thus to implement on the same die more and
25 more complex microarchitectures, up to arrive to single
26 Chip Multiprocessor (CMP) systems, and (2) increasing
27 processors’ clock frequencies.

28 As the memory bandwidth requirements of cores
29 have increased as well, the increased number of on-chip
30 transistors has also been used to integrate on the same
31 die deeper and larger memory hierarchies. However,
32 reducing feature sizes has also introduced the *wire-*
33 *delay* problem: sizes of on-chip wires have been also
34 reduced, resulting in larger delay for signals propaga-
35 tion [1]. In particular, considering a huge traditional

on-chip cache, the wire delay and the increase of clock 36
frequency lead to an increase of the latency experienced 37
on each access. The bulk of access time involves signals’ 38
routing to and from cache banks and not the bank access 39
themselves. Exploiting a non-uniform access time for 40
on-chip caches is a viable solution to the wire-delay 41
problem. 42

By allowing non-uniform access time also to cache 43
memories, it is possible to design a cache based on inde- 44
pendently accessible banks, connected via a scalable 45
connection fabrics, typically a Network-on-Chip (NoC) 46
[2, 4, 5], in which the access latency is proportional to 47
the physical distance between the requesting unit and 48
the accessed bank. This organization is called NUCA 49
(Fig. 1b and c): banks close to the requesting unit can 50
be accessed faster than a traditional monolithic cache. If 51
the cache management policies succeed in keeping the 52
most performance impacting blocks in such banks, the 53
overall performance will be boosted. NUCA architec- 54
tures have been proposed for both single-core [2, 3, 6] 55
and multi-core [4, 5, 7–13] environments. 56

Needs for scalability in CMP designs has lead to the 57
design of CMP tiled systems that are composed by many 58
identical nodes, called tiles. Each tile is composed by 59
one cpu, its private cache levels, and the LLC. Nodes are 60
typically connected through a NoC. In these systems, 61
the LLCs can be used as a globally shared cache, the LLC 62
inside each node being a slice of the whole cache. Hence, 63
the access time depends on respective position of the 64
requesting node and of the memory slice involved in the 65
access, and this results in an non uniform access time. 66
Those systems results to be implicitly NUMA caches as a 67
consequence of their tiled architecture, differently from 68
NUCA architectures whose design explicitly focuses on 69
the non-uniformity in access time in order to tolerate 70
the wire-delay effects. 71



NUMA Caches. Fig. 1 Examples of systems adopting a traditional sub-banked UCA cache **(a)** and a NUCA cache **(b and c)**. In all the designs the memory space is partitioned in 128 banks (white squares in the picture). UCA adopts the H-Tree connection model, while NUCAs use a network on chip made up of routers (black circles in the picture) and links. In this way, each bank is independently accessible from the others and the access latency of a block that is in a close to the controller bank **(b)** is lower than access latency that is in a far bank **(c)**

72 NUMA Caches in the Single-Core

73 Environment: NUCA

74 In a traditionally organized cache, each bank is con-
75 nected to the controller via a fixed length path, usually
76 organized according to the H-Tree model (Fig. 1a). Thus,
77 in a wire delay-dominated context, the access latency
78 of each bank is dominated by the constant length of
79 the path followed by control and data signals, and it
80 is independent from the physical position of the bank.
81 This organization is referred as UCA (Uniform Cache
82 Access) as its access time is *uniform*.

83 The basic idea of NUCA, first proposed by Kim et al.
84 [2], is to design the cache so that banks are connected
85 among them and with the controller via a connection
86 fabric, typically a NoC. Figure 1b and c show a sys-
87 tem equipped with a Level 2 NUCA, partitioned in
88 128 independent banks, connected with a partial 2-D
89 mesh NoC.

90 The adoption of a NoC allows to access each bank
91 independently from the others. In particular, in a
92 NUCA, there is a different access latency for each bank
93 depending on its physical distance from the controller.
94 Being in a wire delay-dominated environment, the sig-
95 nal propagation delay dominates overall access latency.
96 Hence banks that are closer to the controller (Fig. 1b)
97 can be accessed faster than the others that are farther
98 (Fig. 1c). As a consequence, the access time is *non-*
99 *uniform* as it is proportional to the physical distance to
100 be traversed by signals. This architectural design is com-
101 plemented by data management policies that maintain
102 critical blocks in faster banks. This allows better per-
103 formance with respect to UCA architectures by obtain-
104 ing a lower average cache access latency, thus hiding
105 wire-delay effects.

106 Data Management Policies: S-NUCA 107 and D-NUCA

108 When designing a NUCA, one of the main choices is
109 related to the *mapping rule* that dictates which bank
110 can hold each block. The simplest mapping rule is the
111 *static mapping*: a block can exclusively reside in a single
112 predetermined bank, basing on its memory address. A
113 NUCA cache adopting the static mapping is called Static
114 NUCA (S-NUCA). An alternative rule is the *dynamic*
115 *mapping*: a block can be hosted in a set of banks, called
116 *bankset*, and can be moved inside the bankset. A NUCA

cache adopting the dynamic mapping is called Dynamic
NUCA (D-NUCA) [2]. 117 118

S-NUCA 119

In an S-NUCA made up of N banks, the bank to be
accessed is select basing on the value of $\log(N)$ bits of
the address. Proposed policies for static mapping are
sequential mapping, which uses the most significant bits
of the index field of the physical address, and the *inter-*
leaved mapping, which uses the low-order bits of the
same field in order to distribute among different banks
blocks that are memory contiguous, thus reducing bank
contentions. 120 121 122 123 124 125 126 127 128

When searching for a block, the controller injects
into the NoC, a request that will be delivered to the per-
taining bank. If the block is present, i.e., the cache access
results in a *hit*, the bank sends it to the controller. If
the block is not present, i.e., there is a *miss*, the block
is retrieved from the main memory, stored in the cache
bank, and then delivered to satisfy the request. 129 130 131 132 133 134 135

Despite the simple design and management policies,
an S-NUCA usually exhibits lower-average access laten-
cies with respect to traditional UCA caches. However,
as the data mapping doesn't take into account the data
frequency usage, it may happen that more frequently
used data are mapped in banks that are far from the con-
troller, thus resulting in higher access latencies and poor
performances. 136 137 138 139 140 141 142 143

D-NUCA 144

In a D-NUCA, banks are grouped into banksets, and
each block can be stored in any of the banks that belong
to the pertaining bankset: for example, in the systems of
Fig. 1b and c, a bankset could be made up of all banks
belonging to the same row. The bankset is chosen bas-
ing on the address of the block itself, using $\log(\text{number}$
of bankset) bits of the index field, adopting either the
sequential or the interleaved mapping. Blocks can move
from one bank to another of the same bankset thanks
to the *block migration* mechanism. Migrating most fre-
quently used data from farthest banks to closest banks
allows to reduce the average access latency, and thus to
achieve better performance with respect to S-NUCA. 145 146 147 148 149 150 151 152 153 154 155 156 157

The three main topics in designing a D-NUCA are: 158

- *Bank Mapping*: how to map bankset to physical
banks? 159 160

- 161 • *Data Search*: how the possible locations are searched
 162 to find a block?
 163 • *Movement and replacement*: how and when the data
 164 should migrate? In which bank of the bankset a new
 165 block should be placed?

166 Bank Mapping

167 The bank mapping policy of a D-NUCA dictates how
 168 to group the available physical banks into banksets.
 169 Many different policies can be adopted, each offering
 170 a different trade-off in terms of implementation com-
 171 plexity and fairness among the various banksets of the
 172 latency distribution. Proposed mapping policies are [2]:
 173 (1) *simple mapping*, (2) *fair mapping*, and (3) *shared*
 174 *mapping*.

175 The simple mapping (Fig. 2a) maps a bankset to all
 176 the banks belonging to the same row. Such policy is
 177 simple to be implemented, but there are very different
 178 access latencies between banksets that are mapped to
 179 the central rows of the cache and those that are mapped
 180 in the other rows; to access them, a longer network path
 181 along the vertical direction must be traversed. In the fair
 182 mapping (Fig. 2b), this is avoided at the cost of addi-
 183 tional complexity: the banks are allocated to banksets
 184 so that the average access times to each bankset are
 185 equalized. In the shared mapping (Fig. 2c), the banksets
 186 share the banks closest to the controller so that a fast
 187 bank-access is provided to all the sets.

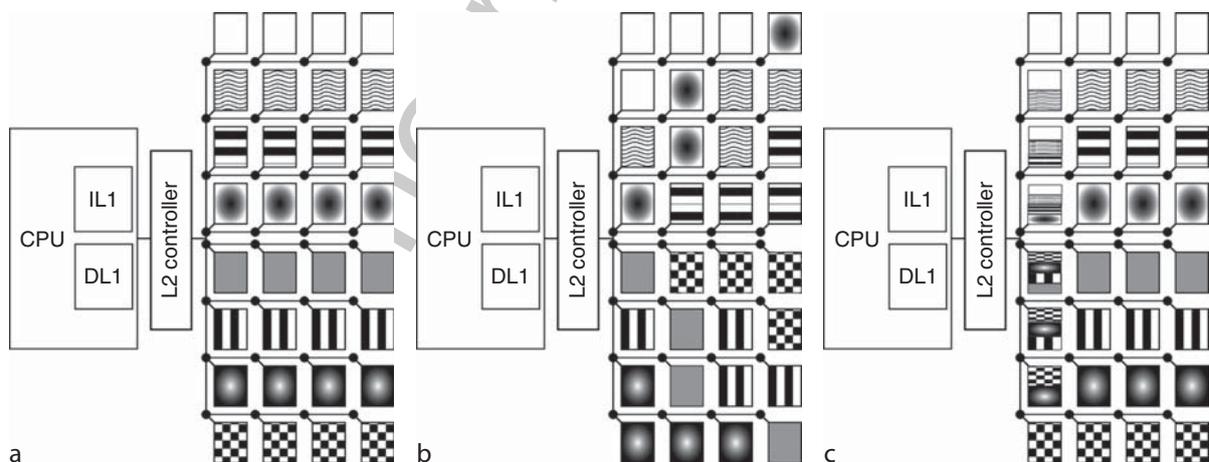
Data Search

188

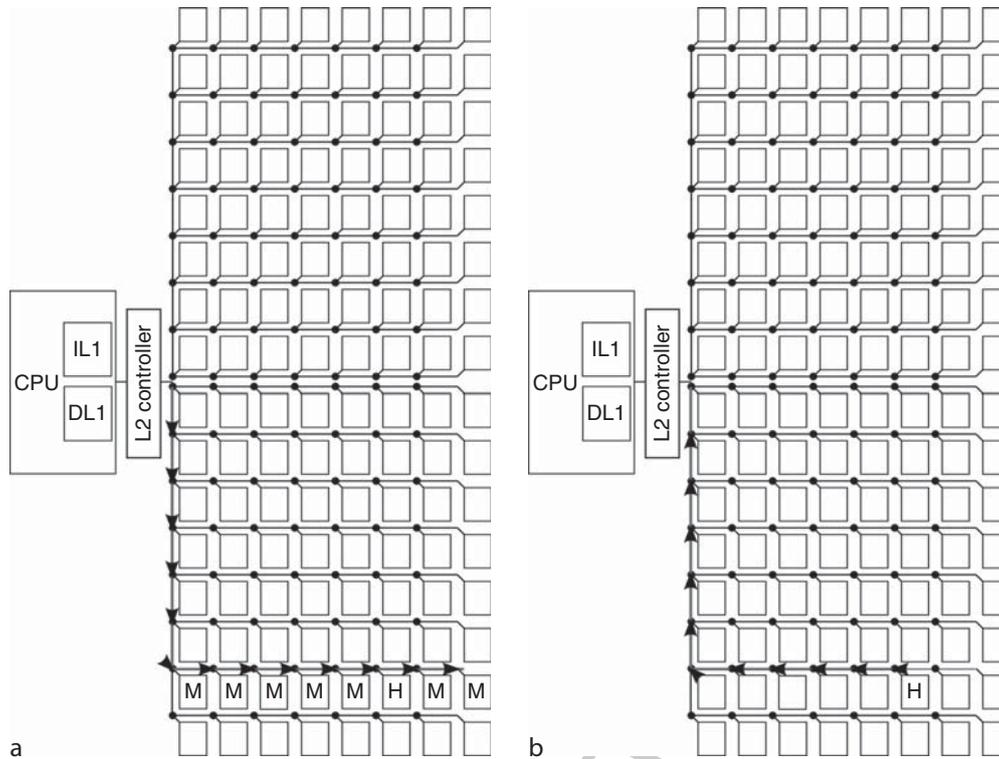
189 When the search is performed, all the banks of the
 190 bankset that can contain a referred block must be
 191 searched. The controller injects in the network a request
 192 that will be delivered to all the pertaining banks. Each
 193 bank performs a search inside itself. If all the banks miss
 194 the searched block, then there is a *cache miss*, and the
 195 block must be retrieved by the next level in the hierarchy
 196 memory. Alternatively, one bank contains the searched
 197 block, there is a *cache hit*, and the block is sent to the
 198 controller.

199 Examples of proposed search policies for D-NUCA
 200 caches are *multicast search* and *incremental search* [2].
 201 In the multicast search, the request is delivered in par-
 202 allel to all the possible banks, and each of them starts the
 203 search as soon as possible. Figure 3 shows an example of
 204 memory operation in a D-NUCA adopting the simple
 205 mapping scheme and the multicast search. As oppo-
 206 site, in the incremental search, the banks of a set are
 207 sequentially searched from the closest to the farthest.
 208 The delivery of the request to each bank is performed
 209 only if the previous bank has missed the block. In such
 210 a way, the number of bank accesses and the associ-
 211 ated energy consumption are reduced at the cost of an
 212 increase of the operation total latency.

213 Other solutions have been proposed [5] mixing the
 214 two techniques: the banks of one bankset are divided
 215 in two or more groups; inside each group the multicast
 216 search is adopted, while the various groups are searched
 217 incrementally.



NUMA Caches. Fig. 2 Three examples of bank mapping: simple mapping (a), fair mapping (b) and shared mapping (c)



NUMA Caches. Fig. 3 The multicast search operation in a D-NUCA adopting the simple mapping: the request first traverses the vertical links to the pertaining bankset, then traverses the horizontal links and routers where is replicated to reach all the banks of the bankset (a). The bank containing the block sends a reply to the controller (b) that follows the inverse path of the request message

218 Data Movement and Replacement

219 Basing on the locality principle, chances there are that,
 220 when a block is referred, it will be referred again in
 221 a short time. So when a block is accessed, it makes
 222 sense to move it in another bank of the bankset that is
 223 closer to the controller. This movement is called *data*
 224 *promotion* [2].

225 An ideally desirable policy would be to use LRU
 226 ordering to order the blocks in the banksets, with the
 227 closest bank holding the MRU block, second closest
 228 holding second most-recently used, etc. Maintaining
 229 this ordering would require a heavy block redistribution
 230 among banks on each access, with consequent impact
 231 on network traffic and energy consumption.

232 For this reason, more practical low-impact policies
 233 are employed usually based on *generational promotions*
 234 [2] of blocks basing on their access pattern. For exam-
 235 ple, in the D-NUCA cache of Fig. 1c on every hit, the

accessed block is moved in the left direction by one 236
 bank. In such a way, the next access to the same block 237
 will incur a lower latency. This migration policy is called 238
per-hit promotion. 239

More generally, designing a generational data promo- 240
 tion, in the case of a D-NUCA, requires the defini- 241
 tion of three elements: the *promotion trigger* (number 242
 of hits after which a promotion must be triggered), the 243
promotion distance (number of banks that the promot- 244
 ing block must be advanced), and the *initial placement* 245
 of a new block after a cache miss. 246

Another important design decision involves what to 247
 do with the block that, in the new bank, is eventually 248
 occupying the line in which a promoting block must 249
 be inserted. Basic choices for this design issue are the 250
demotion of the block in the bank previously hosting the 251
 promoted block or its eviction from the cache. Demotion 252
 requires more bank accesses and network traffic 253

254 than eviction; however, the latter affects performance as
 255 it evicts blocks without taking into account their actual
 256 usefulness.

257 Different choices result in different power, network
 258 traffic, and performance trade-off; however, the major-
 259 ity of the proposed D-NUCA designs [4–6, 8–10] adopt
 260 the per-hit promotion (i.e., one hit as promotion trigger
 261 and one bank as promotion distance) and insert new
 262 block in the farther from the controller bank. Increasing
 263 the value of promotion trigger has been shown to affect
 264 performances too much as it limits the promotion of
 265 newly loaded data that, basing on the temporal locality
 266 principle, are likely to be requested at the next accesses.
 267 At the same time, increasing the promotion distance
 268 involves the rapid demotion of data that, because of the
 269 spatial locality principle, may be needed again in the
 270 next future. Using different insertion point for newly
 271 loaded data means evicting from the cache blocks that
 272 are still useful as they occupy banks that are close to the
 273 controller.

274 NUMA Caches in the CMP Environment

275 The increase in the number of available on-chip transis-
 276 tors and the demand for increasing performance have
 277 driven the design of CMP systems that include more
 278 elaborating cores and a multilevel cache memory hier-
 279 archy in a single die. The adoption of a large on-chip
 280 multilevel cache allows to guarantee fast memory access
 281 to each core.

282 CMP caches are usually organized in two or more
 283 levels in which the Last-Level-Cache (LLC) can be
 284 shared among all on-chip cores, while all the other lev-
 285 els are private for each core. As the LLC is usually quite
 286 large, it is likely to be affected by the wire-delay prob-
 287 lem. For this reason, a CMP system can adopt a shared
 288 NUMA cache as its LLC [4, 5, 8–10].

289 System Topologies

290 Taking into account the respective position of cache
 291 banks and cpus, CMP systems can be classified in two
 292 main topology families: *dancehall* and *tiled*.

293 In the dancehall, the shared LLC cache is concen-
 294 trated in an area of the chip, while cpus stay at one or
 295 more sides of the cache. Taking into account a NUCA
 296 LLC CMP [4, 5, 7–10, 12], cpus and cache banks are

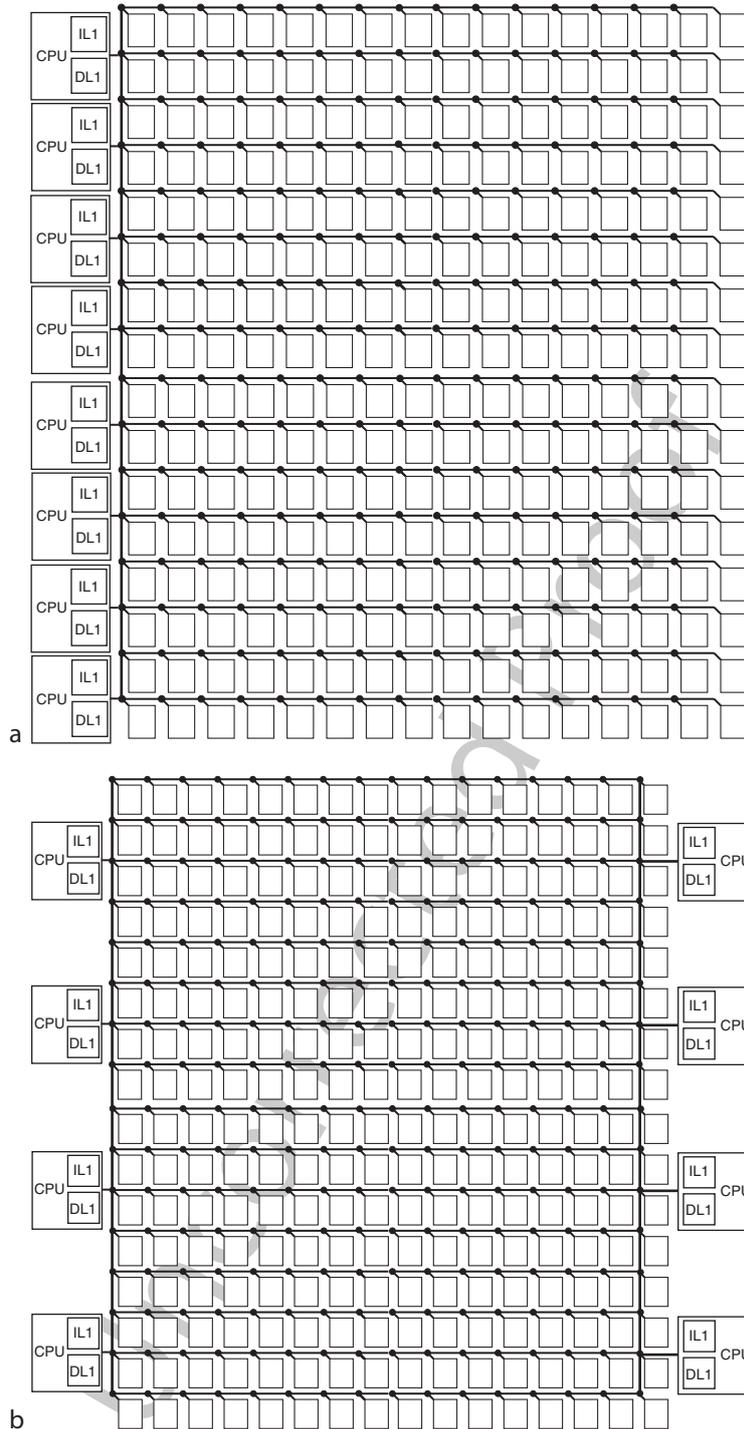
all connected via a NoC (or any other connection fab- 297
 ric), resulting in a NUCA design similarly to the sin- 298
 gle core case. Figure 4a shows an example of dancehall 299
 CMP in which cpus are plugged to the same side of the 300
 NUCA. Variations of this simple configuration have half 301
 of the cpus plugged at one side of the shared NUCA, 302
 and the others plugged to the opposite side (Fig. 4b), or 303
 all the cores distributed along all the sides of the banks 304
 matrix [5]. 305

The tiled topology, as shown in Fig. 5, is composed 306
 by many identical nodes, called tiles. Each tile is com- 307
 posed by one cpu, its private cache levels, and a LLC that 308
 can be used as a private cache or as a local slice of a glob- 309
 ally shared LLC [11–13]. In the last case, the access time 310
 to the shared LLC depends on the position of both the 311
 requesting node and the looked-up LLC slice. In partic- 312
 ular, for a core, an access to a remote slice will result 313
 in greater latency, with respect to an access to the local 314
 slice, due to network paths that vary with the physical 315
 distance from the local to the remote tile. 316

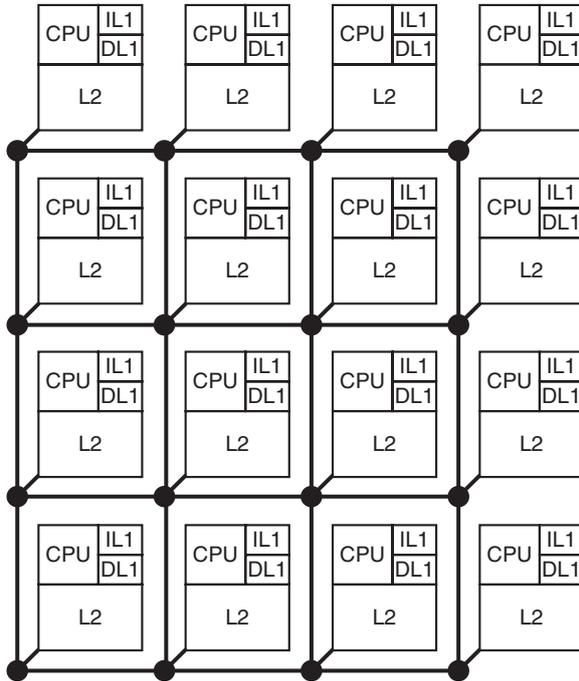
When designing a CMP system, choosing between 317
 such two families of topology represents a trade-off 318
 between performance and scalability. The three main 319
 components of the NUMA access time are: (1) the num- 320
 ber of NoC hops to be traversed to reach the hosting 321
 bank, (2) the latency of each hop, and (3) the access time 322
 to the banks that stores the block. Hence, if the bank to 323
 be accessed is small and close to the requesting node, 324
 the access latency will be limited. If the interested bank 325
 is huge and far from the requesting node, the cost of the 326
 cache access will increase proportionally to the number 327
 of hops to be traversed. 328

Adopting the dancehall topology results in a large 329
 number of small banks, each characterized by small 330
 response latencies. Moreover, as banks are small, NoC 331
 links are short because they have to surround cache 332
 banks. Thus they are characterized by a limited link 333
 traversal delay. Instead, with the tiled organization, 334
 there is a small number of greater and slower banks. 335
 Links are longer since they surround the entire tile, 336
 and consequently, their traversal delay is higher. Conse- 337
 quently, a request to any nonlocal slice traverses a small 338
 number of higher cost network hops. 339

Adopting a tiled topology will increase scalability at 340
 the cost of high cache access due to huge cache slice 341
 and high network hop traversal cost. Instead, adopt- 342
 ing a dancehall topology minimizes the latency of each 343



NUMA Caches. Fig. 4 NUCA based CMP systems adopting the *dancehall* topology, with 8 cpus with private L1 caches, and a shared L2 NUCA cache made of 16×16 banks. All the cpus, together with their private caches, can be plugged at the same side of the NUCA **(a)**, or at two opposite sides **(b)**



NUMA Caches. Fig. 5 A NUCA based CMP system adopting the *tiled* topology, with 16 cpus with private L1 cache, and a shared NUCA cache composed by 16 slices, one for each tile

344 hop, as well as the access time to each single bank,
 345 but presents a smaller degree of scalability than the
 346 tiled case.

347 Coherency in CMP Systems Adopting NUMA
348 Caches

349 In a CMP system, private cache levels must be kept
 350 coherent. Given that NUMA LLC adopt a commu-
 351 nication fabric which does not guarantee sequencing
 352 of coherence traffic (e.g., the NoC), a directory-based
 353 coherence protocol, similarly to the ones proposed for
 354 classical DSM systems, must be adopted. The *directory*
 355 is a node of the system that tracks which of the pri-
 356 vate caches hold a copy of each block. In CMP systems
 357 with NUMA LLC, the directory can be *centralized* or
 358 *distributed*.

359 A centralized directory is a monolithic node of the
 360 system separated from the remaining nodes. This solu-
 361 tion presents some scalability issues because all the
 362 nodes must access it to guarantee the correctness of
 363 memory operations. When the number of such nodes

increases, then directory contentions increase as well,
 364 resulting in higher response time and thus affecting
 365 performance. 366

The distributed directory is typically implemented
 367 inside each NUMA cache bank. Directory information
 368 is stored near the TAG field of each cached block, using
 369 some specific status bits. As a consequence, directory
 370 accesses are distributed among all banks and, when a
 371 *miss* in the last private level cache occurs, an unique
 372 access is used to retrieve both the block and the direc-
 373 tory information. 374

Mapping Policies 375

Also for the NUMA cache in CMP systems, the two pos-
 376 sible mapping policies are static and dynamic. Due to
 377 the presence of many traffic sources, in both cases the
 378 NoC traffic increases proportionally to the increase in
 379 the number of cores per chip. At the same time, the
 380 contention of shared blocks introduces new phenom-
 381 ena that are strictly tied to both the mapping policy and
 382 the topology. 383

Static Mapping 384

Static mapping has been adopted both in the *tiled* [11]
 385 and in the *dancehall* [8, 9, 12] designs. The static map-
 386 ping associates each memory address to a single cache
 387 bank, which is able to satisfy both read/write and coher-
 388 ence requests coming from any of the private next-
 389 level cache. 390

The path that request and coherence messages have
 391 to traverse to reach the interested bank depends on the
 392 respective position of requestor and bank. If the bank
 393 is far from the requestor, the message has to traverse
 394 many network hops before reaching the bank. If the
 395 bank is close to the requestor, the network path will be
 396 short. As the bank to be accessed is chosen basing on
 397 the block address, the banks that a processor accesses
 398 depend on how the *working set* of the running pro-
 399 cess/thread is mapped on memory, i.e., which are the
 400 physical addresses that are accessed. So, if the running
 401 process/thread's *working set* is mapped to remote banks,
 402 then the cache access latency will be high; otherwise
 403 it will be small. Consequently, both topology aspects
 404 and memory mapping rules contribute in affecting or
 405 enhancing overall performance. 406

407 Dynamic Mapping

408 Block migration can be used to improve performance
409 also in CMP systems with a NUMA LLC. *Dancehall* [4,
410 5, 7, 10, 12] and *tiled* [11–13] designs have been proposed
411 in which cache blocks are able to move among banks
412 aiming to reducing response latency.

413 As for the *dancehall* case, the design results in a
414 D-NUCA similar to the single core case. Also in this
415 context, the shared D-NUCA is usually considered as
416 partitioned in many banksets, and cache blocks are able
417 to migrate among banks belonging to the same bankset
418 [4, 5, 8–10]. In the systems shown in Fig. 4a and b,
419 banksets are represented by the rows of the NUCA bank
420 matrix. If D-NUCA succeeds in bringing the most fre-
421 quently accessed blocks near to the referring cpu(s),
422 overall performance will be boosted. However, the per-
423 formance gain that could be obtained with migration
424 comes with some design issues, as in a CMP, and there
425 are many traffic sources that complicate the managing
426 of the block migration process.

427 Blocks that are shared by two or more threads can
428 be accessed by different nodes if the sharers are running
429 on different cpus. This is not necessarily a disadvantage
430 but depends on the overall system topology. By consid-
431 ering the system of Fig. 4a, as all the cpus are placed at
432 the same D-NUCA side, all the blocks migrate toward
433 the same direction because the faster way is the same
434 for all cpus. Instead, in the system of Fig. 4b, the fastest
435 way for half of the cpus is the slowest for the others, and
436 vice versa. If the sharers are running on cpus plugged
437 at opposite D-NUCA sides, shared blocks will alterna-
438 tively migrate in both directions, and none of the cpus
439 will succeed in bringing such blocks in the respective
440 faster way. This phenomenon is known as *conflict hit*,
441 and if not properly managed, reduces the performance
442 gain deriving from block migration.

443 Another design issue of D-NUCA CMP systems
444 regards the *initial placement* of a new block after a cache
445 miss (i.e., choosing which bank of each bankset will
446 store incoming blocks). Considering the *dancehall* con-
447 figuration of Fig. 4a, choosing the fastest way as the
448 block entry point would interfere with most referred
449 blocks, and similar considerations can be done as in
450 the single core case. Instead, for the system shown in
451 Fig. 4b, the farthest way for half of the cpus is the clos-
452 est for the other half, so blocks loaded by half of the cpus

would interfere with most frequently used blocks of the 453
other half. Consequently, a possible trade-off could be 454
choosing the *initial placement* in the central banks. 455

Alternative schemes, for both the *dancehall* and 456
tiled, have been proposed [11, 12], in which blocks are 457
initially placed near the first requestor and stays there 458
until a different cpu requests the block. In this case, 459
the block migrates toward a bank (or a slice) that is 460
determined according to its memory address basing 461
on a static mapping policy. On subsequent requests, 462
the block is not further moved. Dually, migration 463
schemes have been proposed so that blocks are ini- 464
tially placed basing on their memory address and subse- 465
quently migrated basing on the position of the referring 466
cpu [13]. 467

Dynamic block movement designs must face the 468
false miss problem. This is a particular race condition 469
that can occur when a subsequent access to a migrating 470
block arrives when the block is still *on-the-fly*. A block 471
migration involves two banks: the *source* and the *desti-* 472
nation. If the new request is received by the *source* after 473
the block has been sent and by the *destination* before the 474
migrating block has arrived, both the accesses will result 475
in a miss; actually, the block is present in cache and must 476
not be retrieved by the next level memory. That has two 477
disadvantages: (1) affect performance because off-chip 478
accesses are always more expensive than cache accesses; 479
(2) the correctness of memory operations is affected as 480
there would be two (or more) unmanaged copies of the 481
same block. Techniques for resolving such problem have 482
been proposed, for example, a directory holding all the 483
cached block addresses can be accessed in order to rec- 484
ognize false misses [5], or a specific false miss avoidance 485
algorithm can be designed [10]. 486

Research Directions 487

Promising lines of studies for NUMA caches are: 488

1. Power-consumption saving techniques. While being 489
able to reduce wire-delay effects, NUCA caches still 490
suffer of another typical problem of nanometers 491
technologies: the static power consumptions due to 492
the increase of leakage currents. Besides, due to the 493
additional NoC traffic and bank access, D-NUCA 494
exhibit an increase also in the dynamic power con- 495
sumption [14]. Even if NUCA caches are made up of 496

497 traditional cache banks, the direct adoption of well-
 498 known techniques like, for example, drowsy mem-
 499 ories and decay cache lines, may not be effective
 500 particularly in D-NUCA because of the data move-
 501 ment that modify bank access pattern. An alterna-
 502 tive proposed technique [15, 16] exploits different
 503 usage of banks due to block migration to power-gate
 504 farther and less used banks, and has been intro-
 505 duced for both the single-core case and the CMP
 506 environment in which all the cpus are connected
 507 to the same cache side. However, when cpus are
 508 distributed around more NUCA sides due to the
 509 problems that arise in D-NUCA based CMP sys-
 510 tems, the concept of farthest banks is ambiguous and
 511 has to be redefined; thus the raw application of such
 512 technique is no longer possible.

513 2. Performance improvement via remapping policies.
 514 Memory remapping at compile time is tradition-
 515 ally used to improve cache memory performance,
 516 in particular for reducing conflict misses. In the
 517 context of NUCA caches, similar techniques can
 518 be used also for placing data block in banks (or
 519 banksets) that reduce their average access latency
 520 [17]. Typical samples of used metrics in mapping
 521 policies are the data usage frequency and their being
 522 shared or private. Remapping can be implemented
 523 also via hardware, for example, by designing specific
 524 protocols that dynamically change the hosting bank
 525 (in the case of S-NUCA, similarly to [12]) or bankset
 526 (in the case of D-NUCA).

527 3. Efficient block migration schemes are an open issue
 528 for CMP caches. Generational promotion schemes
 529 as in the case of single-core D-NUCA are effective in
 530 *dancehall* systems in which cores are all plugged at
 531 the same cache side. In all the other cases, such sim-
 532 ple schemes are no longer effective due to the *conflict*
 533 *hit* problem that can arise on shared data. Alter-
 534 native solutions are likely to take into account the
 535 position of cores that share a specific block in order
 536 to migrate it in a position that optimizes overall per-
 537 formance (similarly to [13]). Moreover, blocks repli-
 538 cation could be evaluated as a viable solution to the
 539 *conflict hit* problem so that each copy can migrate
 540 toward a specific requestor. Such scheme must take
 541 care of the need of each application. In fact, while
 542 the convenience of replicating Shared-Read-Only
 543 block is obvious, there may be classes of application,

in which replicating also Shared-Read-Write blocks
 can contribute in improving performance.

Related Entries

►Cache Coherency	547
►Interconnection Networks/Topology/Routing	548
►Memory Models	549
►Memory Wall	550
►Network On Chip	551
►Shared Memory Multiprocessor	552
	553

Bibliographic Notes and Further Reading

The first NUCA cache architecture was proposed in
 [2], and demonstrates that a dynamic NUCA struc-
 ture achieves an IPC 1.5 times higher than a traditional
 Uniform Cache Architecture (UCA) when maintaining
 the same size and manufacturing technology. Lever-
 aging on alternatives data mapping policies, NuRapid
 [4] and Triangular D-NUCA cache [6] have been
 proposed as alternative designs for NUCA architec-
 tures in order to optimize the trade-off between area
 occupation and performance. An energy model for
 NUCA architecture is proposed in [14] together with an
 energy/performance trade-off evaluation for NUCAs
 and its comparisons with UCA. The Way Adaptable
 D-NUCA architecture is proposed in [15, 16] shows that
 it is possible to dynamically adapt the cache size to the
 application needs, thus consistently reducing NUCA
 static power consumption while marginally affecting
 the performances. The impact of NoC elements param-
 eters on NUCA performances have been analyzed in
 [18, 19]. The same work proposes an alternative design
 for NUCA caches that relaxes the constraints imposed
 on network routers latency by clustering multiple banks
 around each router.

In the context of on-chip multiprocessors, the
 behavior of D-NUCA as shared L2 caches has been
 analyzed in [4, 5, 7]; these works evaluate the perfor-
 mance of different data mapping and migration policies.
 Topology studies for S-NUCA and D-NUCA, related to
 the dancehall configurations, can be found in [8–10].
 These works also propose a directory coherence pro-
 tocol specifically suited for CMP NUCA environment
 and a solution to some design issue of block migration

588 policies in D-NUCA-based CMP systems. The analy-
 589 sis in [11] is based on the tiled architecture, and stud-
 590 ies also an intelligent block placement, while in [12] is
 591 proposed a scheme called SP-NUCA in which private
 592 blocks are dynamically recognized and then stored in
 593 S-NUCA banks that are very close to the owner. The sys-
 594 tem proposed in [13] is based on a tiled architecture and
 595 implements a migration mechanism that places blocks
 596 in slices depending on the sharers' position.

597 Bibliography

- 598 1. Matzke D (1997) Will physical scalability sabotage performance
 599 gains? *IEEE Comput* 30(9):37–39.
- 600 2. Kim C, Burger D, Keckler SW (2003) Non uniform cache archi-
 601 tectures for wire-delay dominated on-chip caches. *IEEE Micro*
 602 23(6):99–107.
- 603 3. Chisti Z, Powell MD, Vijaykumar TN (2003) Distance associa-
 604 tivity for high-performance energy-efficient non-uniform cache
 605 architectures. *Proc. 36th int. symp. on microarchitecture*. San
 606 Diego, CA, pp 55–66.
- 607 4. Huh J, Kim C, Shafi H, Zhang L, Bourger D, Keckler SW (2005)
 608 A NUCA substrate for flexible CMP cache sharing. *Proc. of the*
 609 *19th int. conf. on supercomputing*. Cambridge, MA, pp 20–22.
- 610 5. Beckmann BM, Wood DA (2003) Managing wire delay in large
 611 chip-multiprocessors caches. *Proc. of 37th int. symp. on microar-*
 612 *chitecture*. San Diego, CA, pp 55–66.
- 613 6. Foglia P, Mangano D, Prete CA (2005) A cache design for high
 614 performance embedded systems. *J Embedded Comput* 1(4):587–
 615 598.
- 616 7. Chisti Z, Powell MD, Vijaykumar TN (2005) Optimizing replica-
 617 tion, communication, and capacity allocation in CMPs. *Proc. of*
 618 *the 32nd int. symp. on computer architecture*. Madison.
- 619 8. Foglia P, Panicucci F, Prete CA, Solinas M (2009) Investigating
 620 design trade-off in S-NUCA based CMP systems. *Proc. work-*
 621 *shop on UNIQUE CHIPS and SYSTEMS (UCAS-5)*. Boston,
 622 vol. 1, pp 53–60.
- 623 9. Foglia P, Panicucci F, Prete CA, Solinas M (2009) An evaluation
 624 of behaviors of S-NUCA CMPs running scientific workload. *Proc.*
 625 *of the 12th euromicro conference on digital system design (DSD)*.
 626 Patras, Greece, pp 26–33.
- 627 10. Foglia P, Panicucci F, Prete CA, Solinas M (2009) Analysis of per-
 628 formance dependencies in NUCA-based CMP systems. *Proc. of*
 629 *the 21st international symposium on computer architecture and*
 630 *high performance computing (SBAC-PAD)*, Sao Paulo.
- 631 11. Hardavellas N, Ferdman M, Falsafi B, Ailamaki A (2009) Reactive
 632 NUCA: near-optimal block placement and replication in dis-
 633 tributed caches. *Proc. of the 36th international symposium on*
 634 *computer architecture (ISCA-09)*. Austin, pp 184–195.
- 635 12. Merino J, Puente V, Prieto P, Gregorio JÁ (2008) SP-NUCA:
 636 a cost effective dynamic non-uniform cache architecture. *ACM*
 637 *SIGARCH Computer Architecture News* 36(2):64–71, New York.
- 638 13. Hammoud M, Cho S, Melhem R (2009) ACM: an efficient
 639 approach for managing shared caches in chip multiprocessors.
 Proc. 4th int. conf. on high performance embedded architectures 640
 and compilers (HiPEAC-09). Paphos, Cyprus, pp 355–372. 641
14. Bardine A, Foglia P, Gabrielli G, Prete CA (2007) Analysis of 642
 static and dynamic energy consumption in NUCA caches: initial 643
 results. *Proc. of the MEDEA 2007 workshop*. Brasov, Romania, 644
 pp 105–112. 645
15. Bardine A, Comparetti M, Foglia P, Gabrielli G, Prete CA, 646
 Stenstrom P (2008) Leveraging data promotion for low power 647
 D-NUCA caches. *Proc. of the DSD '08. 11th EUROMICRO con-*
ference on digital system design architectures, methods and tools. 648
 Parma, Italy, pp 307–316. 650
16. Bardine A, Comparetti M, Foglia P, Gabrielli G, Prete CA (2010) 651
 Way-adaptable D-Nuca caches. *International Journals of High*
Performance Systems Architecture 2(3/4):215–228. 652
17. Bartolini S, Foglia P, Prete CA, Solinas M (2010) Feedback driven 654
 restructuring of multi-threaded applications for NUCA cache 655
 performance in CMPs. *22nd international symposium on com-*
puter architecture and high performance computing. Petropolis, 656
 Brazil. 657
18. Bardine A, Comparetti M, Foglia P, Gabrielli G, Prete CA (2009) 659
 Impact of on-chip network parameters on NUCA cache perfor- 660
 mance. *IET Computers & Digital Techniques* 3(5):501–512. 661
19. Bardine A, Comparetti M, Foglia P, Gabrielli G, Prete CA (2008) 662
 Performance sensitivity of NUCA caches to on-chip network 663
 parameters. *20th international symposium on computer architec-*
ture and high performance computing (IEEE SBAC-PAD 2008). 664
 October 29–November 1, Campo Grande, Brazil. 665
 666