

C

1 Cache-Only Memory Architecture

2 JOSEP TORRELLAS

3 University of Illinois at Urbana-Champaign 4231 Siebel
4 Center, M/C-258, Urbana, IL, USA

5 Synonyms

6 COMA

7 Definition

8 A Cache-Only Memory Architecture (COMA) is a
9 type of cache-coherent nonuniform memory access
10 (CC-NUMA) architecture. Unlike in a conventional
11 CC-NUMA architecture, in a COMA, every shared-
12 memory module in the machine is a cache, where each
13 memory line has a tag with the line’s address and state.
14 As a processor references a line, it transparently brings
15 it to both its private cache(s) and its nearby portion of
16 the NUMA shared memory (*Local Memory*) – possibly
17 displacing a valid line from its local memory. Effectively,
18 each shared-memory module acts as a huge cache mem-
19 ory, giving the name COMA to the architecture. Since
20 the COMA hardware automatically replicates the data
21 and migrates it to the memory module of the node that
22 is currently accessing it, COMA increases the chances of
23 data being available locally. This reduces the possibility
24 of frequent long-latency memory accesses. Effectively,
25 COMA dynamically adapts the shared data layout to the
26 application’s reference patterns.

27 Discussion

28 Basic Concepts

29 In a conventional CC-NUMA architecture, each node
30 contains one or more processors with private caches and
31 a memory module that is part of the NUMA shared
32 memory. A page allocated in the memory module of

one node can be accessed by the processors of all other 33
nodes. The physical page number of the page speci- 34
fies the node where the page is allocated. Such node is 35
referred to as the *Home Node* of the page. The physi- 36
cal address of a memory line includes the physical page 37
number and the offset within that page. 38

In large machines, fetching a line from a remote 39
memory module can take several times longer than 40
fetching it from the local memory module. Conse- 41
quently, for an application to attain high performance, 42
the local memory module must satisfy a large fraction 43
of the cache misses. This requires a good placement of 44
the program pages across the different nodes. If the pro- 45
gram’s memory access patterns are too complicated for 46
the software to understand, individual data structures 47
may not end up being placed in the memory module of 48
the node that access them the most. In addition, when a 49
page contains data structures that are read and written 50
by different processors, it is hard to attain a good page 51
placement. 52

In a COMA, the hardware can transparently elimi- 53
nate a certain class of remote memory accesses. COMA 54
does this by turning memory modules into large caches 55
called *Attraction Memory* (AM). When a processor 56
requests a line from a remote memory, the line is 57
inserted in both the processor’s cache and the node’s 58
AM. A line can be evicted from an AM if another line 59
needs the space. Ideally, with this support, the proces- 60
sor dynamically attracts its working set into its local 61
memory module. The lines the processor is not access- 62
ing overflow and are sent to other memories. Because a 63
large AM is more capable of containing a node’s current 64
working set than a cache is, more of the cache misses are 65
satisfied locally within the node. 66

There are three issues that need to be addressed in 67
COMA, namely finding a line, replacing a line, and deal- 68
ing with the memory overhead. In the rest of this article, 69
these issues are described first, then different COMA 70
designs are outlined, and finally further readings are 71
suggested. 72



73 Finding a Memory Line

In a COMA, the address of a memory line is a global identifier, not an indicator of the line's physical location in memory. Just like a normal cache, the AM keeps a tag with the address and state of the memory line currently stored in each memory location. On a cache miss, the memory controller has to look up the tags in the local AM to determine whether or not the access can be serviced locally. If the line is not in the local AM, a remote request is issued to locate the block.

COMA machines have a mechanism to locate a line in the system so that the processor can find a valid copy of the line when a miss occurs in the local AM. Different mechanisms are used by different classes of COMA machines.

One approach is to organize the machine hierarchically, with the processors at the leaves of the tree. Each level in the hierarchy includes a directory-like structure, with information about the status of the lines present in the subtree extending from the leaves up to that level of the hierarchy. To find a line, the processing node issues a request that goes to successively higher levels of the tree, potentially going all the way to the root. The process stops at the level where the subtree contains the line. This design is called Hierarchical COMA [2, 7].

Another approach involves assigning a home node to each memory line, based on the line's physical address. The line's home has the directory entry for the line. Memory lines can freely migrate, but directory entries do not. Consequently, to locate a memory line, a processor interrogates the directory in the line's home node. The directory always knows the state and location of the line and can forward the request to the right node. This design is called Flat COMA [12].

107 Replacing a Memory Line

The AM acts as a cache, and lines can be displaced from it. When a line is displaced in a plain cache, it is either overwritten (if it is unmodified) or written back to its home memory module, which guarantees a place for the line.

A memory line in COMA does not have a fixed backup location where it can be written to if it gets displaced from an AM. Moreover, even an unmodified line can be the only copy of that memory line in the system, and it must not be lost on an AM displacement. Therefore, the system must keep track of the last

copy of a line. As a result, when a modified or otherwise unique line is displaced from an AM, it must be relocated into another AM.

To guarantee that at least one copy of an unmodified line remains in the system, one of the line's copies is denoted as the *Master* copy. All other shared copies can be overwritten if displaced, but the master copy must always be relocated to another AM. When a master copy or a modified line is relocated, the problem is deciding which node should take the line in its AM. If other nodes already have one or more other shared copies of the line, one of them becomes the master copy. Otherwise, another node must accept the line. This process is called *Line Injection*.

Different line injection algorithms are possible. One approach is for the displacing node to send requests to other nodes asking if they have space to host the line [7]. Another approach is to force one node to accept the line. This, however, may lead to another line displacement. A proposed solution is to relocate the new line to the node that supplied the line that caused the displacement in the first place [8].

Dealing with Memory Overhead

A CC-NUMA machine can allocate all memory to application or system pages. COMA, however, leaves a portion of the memory unallocated to facilitate automatic data replication and migration. This unallocated space supports the replication of lines across AMs. It also enhances line migration to the AMs of the referencing nodes because less line relocation traffic is needed.

Without unallocated space, every time a line is inserted in the AM, another line would have to be relocated. The ratio between the allocated data size and the total size of the AMs is called the *Memory Pressure*. If the memory pressure is 80%, then 20% of the AM space is available for data replication. Both the relocation traffic and the number of AM misses increase with the memory pressure [8]. For a given memory size, choosing an appropriate memory pressure is a trade-off between the effect on page faults, AM misses, and relocation traffic.

160 Different Cache-Only Memory Architecture

161 Designs

162 Hierarchical COMA

163 The first designs of COMA machines follow what has
164 been called Hierarchical COMA. These designs orga-
165 nize the machine hierarchically, connecting the proces-
166 sors to the leaves of the tree. These machines include
167 the KSR-1 [2] from Kendall Square Research, which has
168 a hierarchy of rings, and the Data Diffusion Machine
169 (DDM) [7] from the Swedish Institute of Computer
170 Science, which has a hierarchy of buses.

171 Each level in the tree hierarchy includes a directory-
172 like structure, with information about the status of the
173 lines extending from the leaves up to that level of the
174 hierarchy. To find a line, the processing node issues a
175 request that goes to successively higher levels of the tree,
176 potentially going all the way to the root. The process
177 stops at the level where the subtree contains the line.

178 In these designs, substantial latency occurs as the
179 memory requests go up the hierarchy and then down
180 to find the desired line. It has been argued that such
181 latency can offset the potential gains of COMA relative
182 to conventional CC-NUMA architectures [12].

183 Flat COMA

184 A design called Flat COMA makes it easy to locate a
185 memory line by assigning a home node to each memory
186 line [12] – based on the line’s physical address. The line’s
187 home has the directory entry for the line, like in a con-
188 ventional CC-NUMA architecture. The memory lines
189 can freely migrate, but the directory entries of the mem-
190 ory lines are fixed in their home nodes. At a miss on a
191 line in an AM, a request goes to the node that is keeping
192 the directory information about the line. The directory
193 redirects the request to another node if the home does
194 not have a copy of the line. In Flat COMA, unlike in a
195 conventional CC-NUMA architecture, the home node
196 may not have a copy of the line even though no proces-
197 sor has written to the line. The line has simply been
198 displaced from the AM in the home node.

199 Because Flat COMA does not rely on a hierarchy to
200 find a block, it can use any high-speed network.

201 Simple COMA

202 A design called Simple COMA (S-COMA) [10] transfers
203 some of the complexity in the AM line displacement and

relocation mechanisms to software. The general coher- 204
ence actions, however, are still maintained in hard- 205
ware for performance reasons. Specifically, in S-COMA, 206
the operating system sets aside space in the AM for 207
incoming memory blocks on a page- granularity basis. 208
The local Memory Management Unit (MMU) has map- 209
pings only for pages in the local node, not for remote 210
pages. When a node accesses for the first time a shared 211
page that is already in a remote node, the processor suf- 212
fers a page fault. The operating system then allocates a 213
page frame locally for the requested line. Thereafter, the 214
hardware continues with the request, including locating 215
a valid copy of the line and inserting it, in the cor- 216
rect state, in the newly allocated page in the local AM. 217
The rest of the page remains unused until future 218
requests to other lines of the page start filling it. Sub- 219
sequent accesses to the line get their mapping directly 220
from the MMU. There are no AM address tags to check 221
if the correct line is accessed. 222

Since the physical address used to identify a line in 223
the AM is set up independently by the MMU in each 224
node, two copies of the same line in different nodes are 225
likely to have different physical addresses. Shared data 226
needs a global identity so that different nodes can com- 227
municate. To this end, each node has a translation table 228
that converts local addresses to global identifiers and 229
vice versa. 230

Multiplexed Simple COMA

231

S-COMA sets aside memory space in page-sized 232
chunks, even if only one line of each page is present. 233
Consequently, S-COMA suffers from memory frag- 234
mentation. This can cause programs to have inflated 235
working sets that overflow the AM, inducing frequent 236
page replacements and resulting in high operating sys- 237
tem overhead and poor performance. 238

Multiplexed Simple COMA (MS-COMA) [1] elim- 239
inates this problem by allowing multiple virtual pages 240
in a given node to map to the same physical page at 241
the same time. This mapping is possible because all the 242
lines on a virtual page are not used at the same time. 243
A given physical page can now contain lines belonging 244
to different virtual pages if each line has a short vir- 245
tual page ID. If two lines belonging to different pages 246
have the same page offset, they displace each other 247
from the AM. The overall result is a compression of the 248
application’s working set. 249



250 Further Readings

251 There are several papers that discuss COMA and
 252 related topics. Dahlgren and Torrellas present a more
 253 in-depth survey of COMA machine issues [3]. There
 254 are several designs that combine COMA and con-
 255 ventional CC-NUMA architecture features, such as
 256 NUMA with Remote Caches (NUMA-RC) [9], Reactive
 257 NUMA [5], Excel-NUMA [14], the Sun Microsystems’
 258 WildFire multiprocessor design [6], the IBM Prism
 259 architecture [4], and the Illinois I-ACOMA architecture
 260 [13]. A model for comparing the performance of COMA
 261 and conventional CC-NUMA architectures is presented
 262 by Zhang and Torrellas [15]. Soundarajan et al. [11]
 263 describe the trade-offs related to data migration and
 264 replication in CC-NUMA machines.

265 Related Entries

266 ►Cache-Coherent Non-Uniform Memory Access
 267 (CC-NUMA) architecture

AU1

268 Bibliography

- 269 1. Basu S, Torrellas J (1998) Enhancing memory use in simple coma:
 270 multiplexed simple coma. In: International symposium on high-
 271 performance computer architecture, Las Vegas, February 1998
- 272 2. Burkhardt H et al (1992) Overview of the KSRI computer system.
 273 Technical Report 9202001, Kendall Square Research, Waltham,
 274 February 1992
- 275 3. Dahlgren F, Torrellas J (1999) Cache-only memory architectures.
 276 IEEE Computer Magazine 32(6):72–79, June 1999
- 277 4. Ekanadham K, Lim B-H, Pattnaik P, Snir M (1998) PRISM: an
 278 integrated architecture for scalable shared memory. In: Interna-
 279 tional symposium on high-performance computer architecture,
 280 Las Vegas, February 1998
- 281 5. Falsafi B, Wood D (1997) Reactive NUMA: a design for unify-
 282 ing S-COMA and CC-NUMA. In: International symposium on
 283 computer architecture, Denver, June 1997
- 284 6. Hagersten E, Koster M (1992) WildFire: a scalable path for SMPs.
 285 In: International symposium on high-performance computer
 286 architecture, Orlando, January 1999
- 287 7. Hagersten E, Landin A, Haridi S (1992) DDM – a cache-only
 288 memory architecture. IEEE Computer 25(9):44–54
- 289 8. Joe T, Hennessy J (1994) Evaluating the memory overhead
 290 required for COMA architectures. In: International symposium
 291 on computer architecture, Chicago, April 1994, pp 82–93
- 292 9. Moga A, Dubois M (1998) The effectiveness of SRAM network
 293 caches in clustered DSMs. In: International symposium on high-
 294 performance computer architecture, Las Vegas, February 1998
- 295 10. Saulsbury A, Wilkinson T, Carter J, Landin A (1995) An
 296 argument for simple COMA. In: International symposium on
 297 high-performance computer architecture, Raleigh, January 1995,
 298 pp 276–285

AU2

11. Soundararajan V, Heinrich M, Verghese B, Gharachorloo K, 299
 Gupta A, Hennessy J (1998) Flexible use of memory for repli- 300
 cation/migration in cache-coherent DSM multiprocessors. In: 301
 International symposium on computer architecture, Barcelona, 302
 June 1998 303
12. Stenstrom P, Joe T, Gupta A (1992) Comparative performance 304
 evaluation of cache-coherent NUMA and COMA architectures. 305
 In: International symposium on computer architecture, Gold 306
 Coast, Australia, May 1992, pp 80–91 307
13. Torrellas J, Padua D (1996) The illinois aggressive coma multi- 308
 processor project (I-ACOMA). In: Symposium on the frontiers 309
 of massively parallel computing, Annapolis, October 1996 310
14. Zhang Z, Cintra M, Torrellas J (1999) Excel-NUMA: toward 311
 programmability, simplicity, and high performance. IEEE 312
 Trans Comput 48(2):256–264. Special Issue on Cache Memory, 313
 February 1999 314
15. Zhang Z, Torrellas J (1997) Reducing remote conflict misses: 315
 NUMA with remote cache versus COMA. In: International 316
 symposium on high-performance computer architecture, San 317
 Antonio, February 1997, pp 272–281 318