

Evaluation of Low-Overhead Organizations for the Directory in Future Many-Core CMPs*

Alberto Ros¹ and Manuel E. Acacio²

¹ Dpto. de Informática de Sistemas y Computadores
Universidad Politécnica de Valencia, 46022 Valencia, Spain

² Dpto. de Ingeniería y Tecnología de Computadores
Universidad de Murcia, 30100 Murcia, Spain
aros@gap.upv.es, meacacio@ditec.um.es

Abstract. If current trends continue, today's small-scale general-purpose CMPs will soon be replaced by multi-core architectures integrating tens or even hundreds of cores on-chip. Most likely, some of these many-core CMPs will implement the hardware-managed, implicitly-addressed, coherent caches memory model. Cache coherence in these designs will be probably maintained through a directory-based cache coherence protocol implemented in hardware. The organization of the directory structure will be a key design point due to the requirements in area that it will pose. In this work, we study the effects on performance, network traffic and area that the use of compressed sharing codes for the directory will have in many-core CMPs. In particular, we select two compressed sharing codes previously proposed in the context of large-scale shared-memory multiprocessors that have very small area requirements. Simulation results of 32-core CMPs show that degradations of up to 32% in performance and 350% in network traffic are experienced. Additionally, since some proposals for efficient multicast support in on-chip networks have recently appeared, we also consider the case of using this support in combination with the compressed sharing codes. Unfortunately, we found that multicast support is not enough to remove all the performance degradation introduced by the compressed sharing codes and barely can reduce network traffic.

1 Introduction

In the last years we have witnessed the substitution of single-core processors by multi-core ones. Following the Moore's Law that establishes that the number of transistors doubles every 18 months, it is expected that current small-scale general-purpose chip-multiprocessors (CMPs) will soon be followed by multi-core architectures integrating tens or even hundreds of cores on-chip [1]. Architectures of this type are usually known as many-core CMPs.

* We would like to thank anonymous reviewers for their suggestions. This research was supported by the Spanish MEC and MICINN, as well as European Commission FEDER funds, under Grants CSD2006-00046 and TIN2009-14475-C04, and PROMETEO from Generalitat Valenciana (GVA) under Grant PROMETEO/2008/060.

Many-core CMPs will be probably designed as arrays of identical or close-to-identical building blocks (tiles) connected over a switched direct network [2,3]. Tiled architectures provide a scalable solution for supporting families of products with varying computational power, managing the design complexity, and effectively using the resources available in advanced VLSI technologies. As an example, Intel has recently announced the 48-core Single-chip Cloud Computer [4], an experimental research microprocessor that has been developed in the context of the Tera-scale Computing Research Program. The Single-chip Cloud Computer consists of 24 tiles with two IA cores per tile, which are interconnected by means of a 24-router mesh network providing 256 GB/s bisection bandwidth.

On the other hand, if current trends continue, future many-core CMP architectures will implement the hardware-managed, implicitly-addressed, coherent caches memory model [5]. With this memory model, all on-chip storage is used for private and shared caches that are kept coherent in hardware by using a cache coherence protocol. In this way, each tile contains at least one level of cache memory that is private to the local core (the L1 in this work), and the first level of shared cache (commonly, the L2 cache) is physically distributed between the tiles of the system.

The cache coherence protocol will be a key design issue in these architectures since it will add requirements of area and energy consumption to the final design, and therefore, could restrict severely its scalability. When the number of cores is large, as is the case of many-core CMPs, the best way today of keeping cache coherence is by implementing a directory-based protocol, which reduces energy consumption compared to broadcast-based protocols by keeping track of the caches that hold copies of each block in a directory structure. In tiled CMPs, the directory structure is distributed between the L2 cache banks, usually included into the L2 tags' portion [3]. In this way, each tile keeps the sharing information of the blocks mapped to the L2 cache bank that it contains. This sharing information comprises two main components: the *state bits* used to codify one of the three possible states the directory can assign to the line (*Uncached*, *Shared* and *Private*), and the *sharing code*, that holds the list of current sharers. Most of the bits of each directory entry are devoted to codifying the sharing code.

In a traditional directory organization, each directory entry keeps track of the sharers of the corresponding memory block through a simple bit-vector (one bit per private cache). In Figure 1, we plot the area (in mm^2) that one 1MB 4-way L2 module would take as the number of cores grows from 2 to 256 (area estimations are based on CACTI. Refer to Section 4 for more details). As it can be seen, while the number of cores keeps below 16 the bit-vector sharing code barely impacts area requirements. However, from 16 cores and onwards, the use of bit-vectors would entail too much area overhead and more area efficient sharing codes would be required.

One approach for reducing directory area requirements in the context of traditional shared-memory multiprocessors is the use of compressed sharing codes. Compressed sharing codes store the directory information in a compressed way

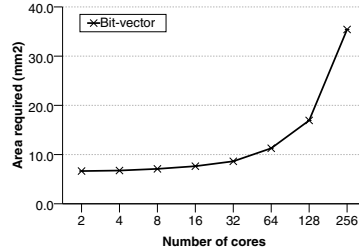


Fig. 1. Area (mm²) required for a 1MB cache module when the bit-vector sharing code is used

to use fewer number of bits, introducing a loss of precision compared to *exact* ones (e.g., bit-vector). This means that when this information is reconstructed, some of the cores codified in the sharing code are real sharers and must receive the coherence messages, whereas some other cores are not sharers actually and *unnecessary* coherence messages will be sent to them. Unnecessary coherence messages lead to increased miss latencies, since more messages are required to resolve caches misses. These messages also entail extra traffic in the interconnection network and useless cache accesses, which will increase energy consumption. Conversely, a bit-vector directory does not generate unnecessary coherence messages and thus shows the best results in terms of both performance and energy consumption.

In this work we study the effects on performance, network traffic and area required by the directory structure that the use of compressed sharing codes will have in many-core CMPs. In particular, we select two area-efficient compressed sharing codes previously proposed by us in the context of large-scale shared-memory multiprocessors, namely Binary Tree (BT) and Binary Tree with Symmetric Nodes (BT-SN) [6]. Simulation results of 32-core CMPs show that degradations of up to 32% in performance and 350% in network traffic are experienced. Additionally, since some proposals for efficient multicast support in on-chip networks have recently appeared [7], we also consider the case of using this kind of support in combination with the compressed sharing codes. Unfortunately, multicast support is not enough to completely remove the performance degradation that the compressed sharing codes introduce (performance degradations of 10% on average are still observed when BT is used) and barely can reduce network traffic.

The rest of the paper is organized as follows. First of all, we will give more details regarding the target CMP architecture in Section 2. Subsequently, in Section 3 we will present a couple of compressed sharing codes based on the concept of multilayer clustering. Next, in Section 4, we will describe the evaluation environment that we are assuming, and the results of the evaluation will be shown in Section 5. Finally, Section 6 closes the work.

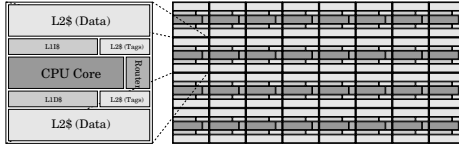


Fig. 2. Organization of the tile assumed in this work and a 4×8 tiled CMP

2 Base Architecture

A tiled CMP architecture consists of a number of replicated *tiles* connected over a switched direct network. Each tile contains a processing core with primary caches (both instruction and data caches), a slice of the L2 cache, and a connection to the on-chip network. Cache coherence is maintained at the L1 caches. In particular, a directory-based cache coherence protocol with directory information stored in the tags' part of the L2 cache modules is employed. The L2 cache is shared among the different processing cores, but it is physically distributed between them. Therefore, some accesses to the L2 cache will be sent to the local slice while the rest will be serviced by remote slices (L2 NUCA architecture [8]). Moreover, for simplicity the L1 and L2 caches are inclusive, that is to say, all the blocks included in any L1 cache keep an entry in the L2 cache. Figure 2 shows the organization of a tile (left) and a 16-tile CMP (right). From now on, we will use the terms tile and node interchangeably.

3 BT and BT-SN Compressed Sharing Codes

The two compressed sharing codes considered in this work (BT and BT-SN) were derived from the multi-layer clustering concept introduced in [6]. Multi-layer clustering assumes that nodes are recursively grouped into clusters of equal size until all nodes are grouped into a single cluster. Compression is achieved by specifying the smallest cluster containing all the sharers (instead of indicating *all* the sharers). Compression can be increased even more by indicating only the level of the cluster in the hierarchy. In this case, it is assumed that the cluster is the one containing the home node for the memory block. Although clusters can be formed by grouping any integer number of clusters in the immediately lower layer of the hierarchy, we analyze the case of using a value equal to two. That is to say, each cluster contains two clusters from the immediately lower level. By doing so, we simplify binary representation and obtain better granularity to specify the set of sharers. This recursive grouping into layer clusters leads to a logical binary tree with the nodes located at the leaves.

Since nodes are located at the leaves of a tree, the set of nodes (sharers) holding a copy of a particular memory block can be expressed as the minimal subtree that includes the home node and all the sharers. This minimal subtree is codified using the level of its root (which can be expressed using just $\lceil \log_2(\log_2 N + 1) \rceil$ bits). Intuitively, the set of sharers is obtained from the home node identifier by

Table 1. System parameters

32-core CMP			
GEMS Parameters		SICOSYS Parameters	
Processor frequency	4 GHz	Network frequency	2 GHz
Cache hierarchy	Inclusive	Topology	8x4 Mesh
Cache block size	64 bytes	Switching technique	Wormhole, Multicast
Split L1 I & D caches	128KB, 4 ways, 4 hit cycles	Routing technique	Deterministic X-Y
Shared unified L2 cache	1MB/tile, 4 ways, 7 hit cycles	Message size	4 flits data, 1 flit control
Memory access time	300 cycles	Routing time	2 cycles
		Link latency (one hop)	2 cycles
		Link bandwidth	1 flit/cycle

changing the value of some of its least significant bits to *don't care*. The number of modified bits is equal to the level of the above mentioned subtree. It constitutes a very compact sharing code (observe that, for a 128-node system, only 3 bits per directory entry are needed). This sharing code is known as *binary tree* or BT.

We also considered the concept of symmetric nodes of a particular home node. Assuming that 3 additional symmetric nodes are assigned to each home node, they are codified by different combinations of the two most-significant bits of the home node identifier (note that one of these combinations represents the home node itself). In other words, symmetric nodes only differ from the corresponding home node in the two most significant bits. Now, the process of choosing the minimal subtree that includes all the sharers is repeated for the symmetric nodes. Then, the minimum of these subtrees is chosen to represent the sharers. The intuitive idea is the same as before but, in this case, the two most significant bits of the home identifier are changed to the symmetric node used. Therefore, the size of the sharing code of a directory entry is the same as before plus the number of bits needed to codify the symmetric nodes (for 3 sym-nodes, 2 bits). This sharing code is known as *binary tree with symmetric nodes* or BT-SN.

4 Evaluation Environment

We perform the evaluation using the full-system simulator Virtutech Simics [9] extended with Multifacet GEMS 1.3 [10], that provides a detailed memory system timing model. Since the network modeled by GEMS 1.3 is not very precise, we have extended it with SICOSYS [11], a detailed interconnection network simulator. We simulate a 32-tile CMP architecture as the one described in Section 2. The values of the main parameters used for the evaluation are shown in Table 1. Cache latencies have been calculated using the CACTI 5.3 tool [12] for 45nm technology. We also have used CACTI to measure the area of a 1MB 4-way L2 cache bank that includes the different sharing codes assumed in this work. In this study, we assume that the length of the physical address is 44 bits, like in the SUN UltraSPARC-III architecture [13].

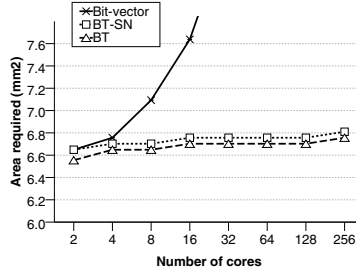


Fig. 3. Area (mm^2) required for a 1MB cache module when bit-vector, BT or BT-SN are used

The ten applications used in our simulations cover a variety of computation and communication patterns. *Barnes* (8192 bodies, 4 time steps), *FFT* (256K points), *Ocean* (258x258 ocean), *Radix* (1M keys, 1024 radix), *Ray-trace* (teapot), *Volrend* (head) and *Water-Sp* (512 molecules, 4 time steps) are scientific applications from the SPLASH-2 benchmark suite [14]. *Unstructured* (Mesh.2K, 5 time steps) is a computational fluid dynamics application. *MPGdec* (525_tens_040.m2v) and *MPGenc* (output of *MPGdec*), are multimedia applications from the APLBench suite [15]. We account for the variability in multithreaded workloads by doing multiple simulation runs for each benchmark in each configuration and injecting random perturbations in the memory systems timing for each run.

5 Evaluation Results

We start this section by comparing the area overhead introduced by the different organizations for the sharing code considered in this work. Next, we study the impact that the compressed sharing codes have on network traffic, considering both a network with and without multicast support. Finally, we compare the execution times for the three directory organizations.

5.1 Impact on Area Overhead

Figure 3 plots the total area (in mm^2) that would be required by a 1MB 4-way cache module when bit-vector, BT and BT-SN sharing codes are used. Due to the limited number of cores used in our simulations (32), we evaluate BT-SN assuming only one symmetric node. In this way, the size of BT-SN is equal to the size of BT plus 1 bit to codify whether the home node or the symmetric node is being used in the codification.

As shown in Figure 3 (and discussed in the introduction of this work), the area overhead that the bit-vector sharing code entails does not scale with the number of cores. Obviously, the size of the bit-vector (in bits) increases linearly with the

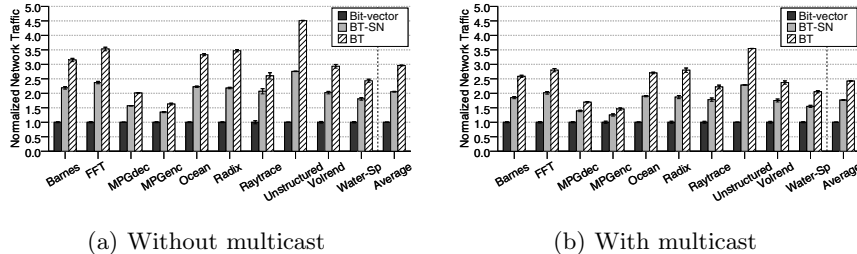


Fig. 4. Normalized network traffic for bit-vector, BT and BT-SN

number of cores. For this reason, the bit-vector could be a good option for a small number of cores. However, for 16 or more cores the increase in area that the bit-vector conveys makes it infeasible (the area overhead becomes almost 100% for the 64-core configuration). On the other hand, the size of BT and BT-SN barely increases with the number of cores. Moreover, the total number of bits needed by BT and BT-SN is very small in all cases ($\lceil \log_2(\log_2 N + 1) \rceil$ bits and $\lceil \log_2(\log_2 N + 1) \rceil + 1$ bits, respectively). In this way, the area overhead of BT and BT-SN is very low (less than 5% for the 256-core configuration) and keeps almost constant with the number of cores. This makes BT and BT-SN promising alternatives to bit-vector for future many-core CMPs, since besides introducing very small overheads in terms of area, they would allow to support families of CMPs with varying number of cores using the same tile structure (without requiring any modifications in the directory).

5.2 Impact on Network Traffic

Although compressed sharing codes can drastically reduce the size of the directory, their drawback is that they could increase the number of coherence messages as a consequence of the in-excess codification of the sharers that they perform. Increasing the number of coherence messages leads to more traffic being injected in the interconnection network of the CMP. Since previous works have identified the interconnection network as one of the most important elements of the CMP from the point of view of energy consumption (consuming almost 40% of the total energy budget in the Raw processor [16]), more traffic at the end means more energy.

Figure 4 shows the amount of network traffic that would be generated for bit-vector, BT and BT-SN for the 32-core CMP configuration assumed in this work. In particular, each bar plots the number of bytes transmitted through the interconnection network (the total number of bytes transmitted by all the switches) normalized with respect to the bit-vector case. We present results considering both a network with unicast support (a) and with multicast support (b).

As shown in Figure 4(a), the use of BT has severe impact on the amount of network traffic and degradations ranging from approximately 50% for *MPGenC* to 350% for *Unstructured* are found. The problem with BT is that when one of the

sharers is far from the home node in the logical tree, the root of the tree is selected as the minimum tree level covering both the home node and the sharer, which results in all cores being actually codified. We have found that this situation occurs frequently in most applications, which explains the significant amount of extra traffic for BT. In particular, the average number of coherence messages that are sent on a coherence event¹ increases from 2 in bit-vector to more than 20 in BT. On the contrary, when BT-SN is considered the tree level that covers all the sharers can be computed from either the home node or its symmetric node. This leads to noticeable reductions in the average number of coherence messages (12 in BT-SN), which leads to important savings in network traffic when compared with BT. Unfortunately, BT-SN does not mitigate completely the extra traffic introduced by BT and degradations of approximately 100% on average are still observed. Again, when two or more cores, distant in the logical tree, share a memory block, the root of the tree would be codified by BT-SN.

Obviously, the provision of multicast support at the interconnection network level can alleviate the levels of extra traffic. More specifically, in Figure 4(b) we show the results obtained when we take advantage of multicast support for sending coherence messages (invalidations and cache-to-cache transfer commands). Efficient implementations of such kind of multicast support in on-chip networks have recently been proposed [7]. Unfortunately, using multicast support for also the response messages is not a trivial issue. So, in this work we assume that responses to coherence commands are unicast messages. As it can be seen, the use of multicast support is a step forward in achieving the network traffic levels obtained by bit-vector, and it is especially useful when BT is considered (average traffic overhead is reduced from 200% without multicast support to 150%). Anyway, the fact that multicast support is available just for the coherence commands and not for their associated responses limits its benefits.

5.3 Impact on Execution Time

The degradations previously reported in terms of network traffic finally translate into increases in terms of execution time. In Figure 5 we show how the use of BT and BT-SN impacts applications' execution times, considering an interconnection network with and without multicast support, (a) and (b) respectively. Again, all results have been normalized with respect to the bit-vector case.

As observed in Figure 5(a), the use of BT without multicast support has important consequences on performance. In particular, the execution time grows from less than 10% for *Barnes* and *Water-Sp* to more than 30% for *Raytrace* (19% on average). In general, the greater number of messages that are needed with BT to resolve every coherence event leads to longer cache miss latencies, and therefore, execution times. Obviously, the extent of the degradation in execution time will depend on the particular characteristics of each application (L1

¹ By coherence event we refer to a situation where the home node must use the sharing code to send coherence messages (invalidations or cache-to-cache transfer commands).

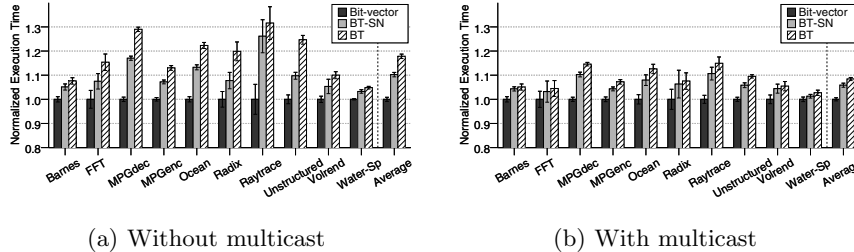


Fig. 5. Execution time for 32 cores

cache miss rate, average number of coherence messages per cache miss, kind of synchronization used, etc.). This is why there is no direct correlation between the amount of extra traffic reported in Figure 4(a) and the degradation in execution time shown in Figure 5(a). On the other hand, when BT-SN is used instead of BT, the average overhead in terms of execution time is reduced to a half (10%). In this case, significant reductions in execution time are observed for most applications. The exceptions are *Barnes* and *Water-Sp*, that hardly see their execution times reduced when BT-SN is used, even when significant savings in terms of network traffic were reported.

The effects of using multicast support with BT and BT-SN are analyzed in Figure 5(b). As before, multicast support has significant impact on execution time when BT is assumed. In this case, average degradation falls from 19% to less than 10%. Although all applications benefit from multicast support, *FFT*, *MPGdec*, *Radix*, *Raytrace* and *Unstructured* are the most affected (in all these cases performance degradation entailed by BT is reduced to more than a half). Finally, and as it was reported for network traffic, multicast support does not help much in reducing performance overhead when BT-SN is considered. In this case, what dominates cache miss latencies is the time taken to collect all responses to a coherence event, which is not optimized with the assumed multicast support.

6 Conclusions

The organization of the directory needed to maintain cache coherence will be a key design point in future many-core CMPs. In this work, we have analyzed the effects that the BT and BT-SN compressed sharing codes have on area, network traffic (as representative of the energy consumed in the interconnection network), and performance in the context of many-core chip-multiprocessors. In particular, we have found that although very area-efficient directories could be derived based on these two sharing codes (with area overheads of less than 5%), the degradations in terms of network traffic (200% for BT and 100% for BT-SN) as well as execution time (20% for BT and 10% for BT-SN) that they entail could preclude them from being employed in future many-core CMPs. Moreover,

we have studied the case of having an interconnection network with multicast support, and have found that although BT can significantly benefit from such kind of support (degradations in execution time and network traffic are reduced to 8% and 150% respectively), BT-SN barely finds any benefits from it. The reasons why multicast support is unable to remove completely the degradation that BT and BT-SN introduce are two. First, multicast support is only used for sending coherence commands but not for collecting the responses. And second, even if an efficient mechanism able to provide combined responses were used, more destinations for the coherence commands still implies more traffic and longer cache miss latencies.

References

1. Borkar, S.: Thousand core chips: A technology perspective. In: 44th Annual Design Automation Conference, pp. 746–749 (2007)
2. Taylor, M.B., Kim, J., Miller, J., et al.: The raw microprocessor: A computational fabric for software circuits and general purpose programs. *IEEE Micro* 22, 25–35 (2002)
3. Zhang, M., Asanović, K.: Victim replication: Maximizing capacity while hiding wire delay in tiled chip multiprocessors. In: 32nd Int'l Symp. on Computer Architecture (ISCA), pp. 336–345 (2005)
4. Intel Res.: Single-chip Cloud Computer (2010), <http://techresearch.intel.com/articles/Tera-Scale/1826.htm>
5. Leverich, J., Arakida, H., Solomatnikov, A., Firoozshahian, A., Horowitz, M., Kozyrakis, C.: Comparing memory systems for chip multiprocessors. In: 34th Int'l Symp. on Computer Architecture (ISCA), pp. 358–368 (2007)
6. Acacio, M.E., González, J., García, J.M., Duato, J.: A new scalable directory architecture for large-scale multiprocessors. In: 7th Int'l Symp. on High-Performance Computer Architecture (HPCA), pp. 97–106 (2001)
7. Rodrigo, S., Flich, J., Duato, J., Hummel, M.: Efficient unicast and multicast support for CMPs. In: 41st IEEE/ACM Int'l Symp. on Microarchitecture (MICRO), pp. 364–375 (2008)
8. Kim, C., Burger, D., Keckler, S.W.: An adaptive, non-uniform cache structure for wire-delay dominated on-chip caches. In: 10th Int. Conf. on Architectural Support for Programming Language and Operating Systems (ASPLOS), pp. 211–222 (2002)
9. Magnusson, P.S., Christensson, M., Eskilson, J., et al.: Simics: A full system simulation platform. *IEEE Computer* 35, 50–58 (2002)
10. Martin, M.M., Sorin, D.J., Beckmann, B.M., et al.: Multifacet's general execution-driven multiprocessor simulator (GEMS) toolset. *Computer Architecture News* 33, 92–99 (2005)
11. Puente, V., Gregorio, J.A., Beivide, R.: SICOSYS: An integrated framework for studying interconnection network in multiprocessor systems. In: 10th Euromicro Workshop on Parallel, Distributed and Network-based Processing, pp. 15–22 (2002)
12. Thoziyoor, S., Muralimanohar, N., Ahn, J.H., Jouppi, N.P.: CACTI 5.1. Technical Report HPL-2008-20, HP Labs (2008)
13. Horel, T., Lauterbach, G.: UltraSPARC-III: Designing third-generation 64-bit performance. *IEEE Micro* 19, 73–85 (1999)

14. Woo, S.C., Ohara, M., Torrie, E., Singh, J.P., Gupta, A.: The SPLASH-2 programs: Characterization and methodological considerations. In: 22nd Int'l Symp. on Computer Architecture (ISCA), pp. 24–36 (1995)
15. Li, M.L., Sasanka, R., Adve, S.V., Chen, Y.K., Debes, E.: The ALPBench benchmark suite for complex multimedia applications. In: Int'l Symp. on Workload Characterization, pp. 34–45 (2005)
16. Wang, H., Peh, L.S., Malik, S.: Power-driven design of router microarchitectures in on-chip networks. In: 36th IEEE/ACM Int'l Symp. on Microarchitecture (MICRO), pp. 105–111 (2003)