



Optimizing Software Applications for NUMA

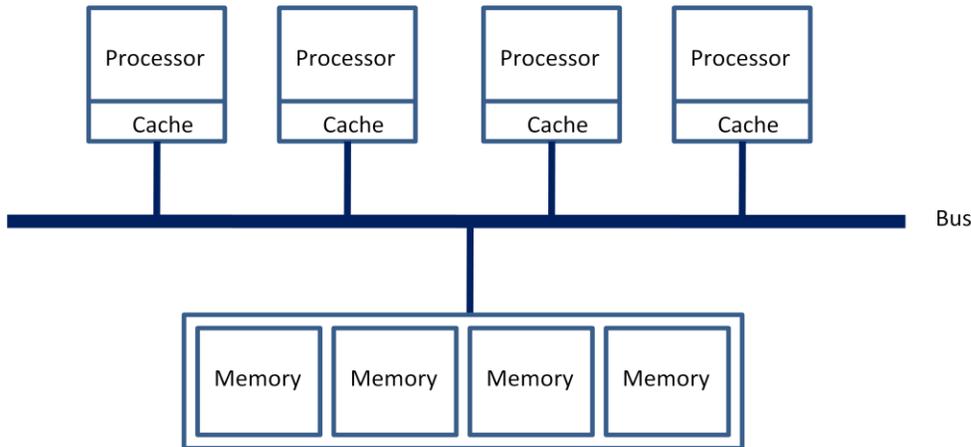
Introduction

In this brief technical paper, we provide an overview of the NUMA shared memory architecture and describe various techniques for optimizing application memory performance within a NUMA-based system. In particular, we discuss the role of processor affinity, memory allocation using implicit operating system policies, and the use of the system API's for assigning and migrating memory pages using explicit directives.

1. The Basics of NUMA

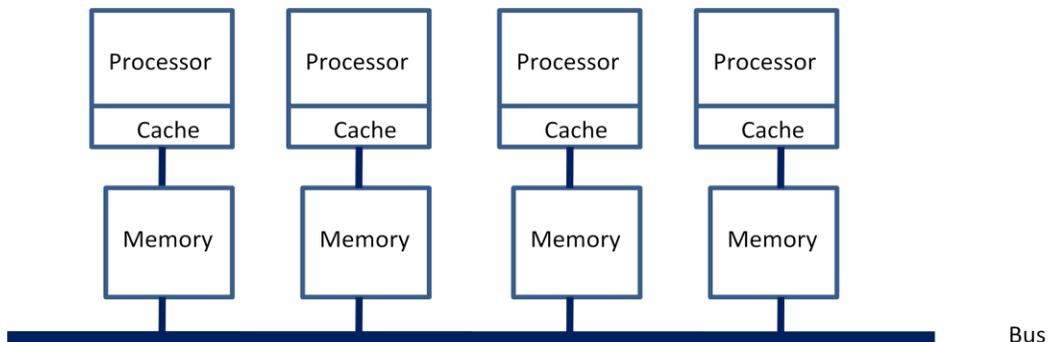
NUMA, or **Non-Uniform Memory Access**, is a shared memory architecture that describes the placement of main memory modules with respect to processors in a multiprocessor system. Perhaps the best way to understand NUMA is to compare it with its cousin **UMA**, or **Uniform Memory Access**.

In the UMA memory architecture, all processors access shared memory through a bus (or another type of interconnect) as seen in the following diagram:



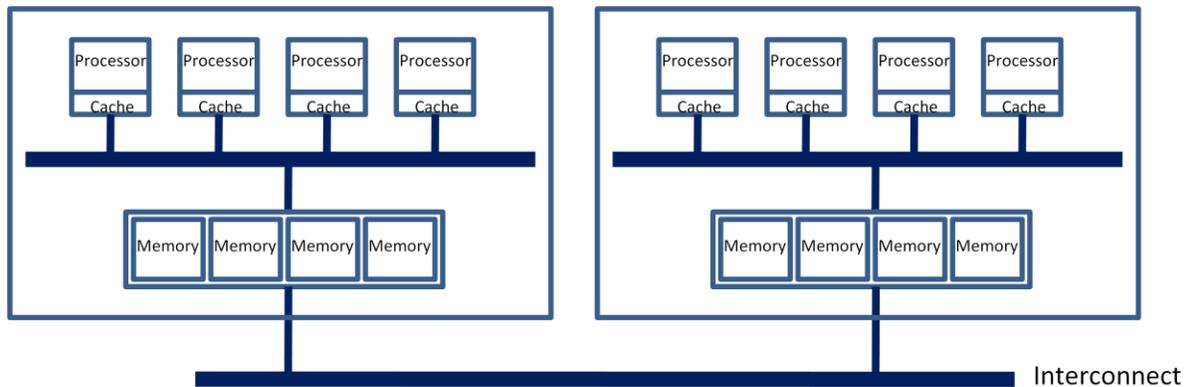
UMA gets its name from the fact that each processor must use the same shared bus to access memory, resulting in a memory access time that is uniform across all processors. Note that access time is also independent of data location within memory. That is, access time remains the same regardless of which shared memory module contains the data to be retrieved.

In the NUMA shared memory architecture, each processor has its own *local* memory module that it can access directly and with a distinctive performance advantage. At the same time, it can also access any memory module belonging to another processor using a shared bus (or some other type of interconnect) as seen in the diagram below:



What gives NUMA its name is that memory access time varies with the location of the data to be accessed. If data resides in local memory, access is fast. If data resides in remote memory, access is slower. The advantage of the NUMA architecture as a *hierarchical* shared memory scheme is its potential to improve *average case access time* through the introduction of fast, local memory.

Modern multiprocessor systems mix these basic architectures as seen in the following diagram:



In this complex hierarchical scheme, processors are grouped by their physical location on one or the other multi-core CPU package or “node”. Processors within a node share access to memory modules as per the UMA shared memory architecture. At the same time, they may also access memory from the remote node using a shared interconnect, but with slower performance as per the NUMA shared memory architecture.

Server platforms like Intel® Xeon using the Intel® Core i7 processors provide an example of this complex memory architecture, and for this reason our discussion will center on it henceforth. Note that such platforms employ a fast interconnect technology known as Intel® QuickPath Interconnect (QPI) to mitigate (but not eliminate) the problem of slower remote memory performance.

2. NUMA Advantages and Risks

The advantage of the NUMA shared memory architecture is its *potential* to reduce memory access time in the average case. By providing each node with its own local memory, memory accesses can take place in parallel and avoid throughput limitations and contention issues associated with a shared memory bus. In fact, memory constrained systems can theoretically improve their performance by up to the number of nodes on the system. For example, a memory-constrained dual processor system could conceivably double its performance if processors could access memory in a fully parallelized manner.

The downside of the NUMA architecture, however, is the cost associated when data is not local to the processor. In the NUMA model, the time required to retrieve data from an adjacent node within the NUMA model will be significantly higher than that required to access local memory. Furthermore, the time required to retrieve data from a non-adjacent node may be even higher, complicating memory performance and generating a hierarchy of access time possibilities. In general, as the *distance* from a processor increases, the cost of accessing memory increases.²

The key issue in determining whether the performance benefits of the NUMA architecture can be realized, then, is **data placement**. The more data can effectively be placed in memory local to the processor that needs it, the more overall access time will benefit from the architecture. Conversely, the more data fails to be local to the node that will access it, the more memory performance will suffer from the architecture. For this reason, the NUMA architecture can be said to provide the *potential* to reduce overall memory access times. To realize this potential, strategies are needed to ensure smart data placement. An application that effectively manages such placement is one that has been “optimized for NUMA”, is “NUMA-aware”, or is “NUMA-friendly”.

3. Strategies for NUMA Optimization

Two key notions in managing performance within the NUMA shared memory architecture are *processor affinity* and *data placement*.

3.1. Processor Affinity

Affinity refers to the persistence of association with a particular resource instance, despite the availability of another instance for the same purpose. Consider the case of processor affinity. Today's complex operating systems assign application threads to processor cores using a scheduler. A scheduler will take into account system state and various policy objectives (e.g., "balance load across cores" or "aggregate threads on a few cores and put remaining cores to sleep") and match application threads to physical cores accordingly. A given thread will execute on its assigned core for some period of time and then wait as other threads are given the chance to execute. If another core becomes available, the scheduler may choose to migrate the thread to insure timely execution and meet its policy objectives.

Thread migration from one core to another poses a problem for the NUMA shared memory architecture because of the way it disassociates a thread from its local memory allocations. That is, a thread may allocate memory on node 1 at startup as it runs on a core within the node 1 package. But when the thread is later migrated to a core on node 2, the data stored earlier becomes remote and memory access time significantly increases.

Enter processor affinity. Using a system API, or by modifying an OS data structure (e.g., affinity mask), a specific core or set of cores can be associated with an application thread. The scheduler will then observe this affinity in its scheduling decisions for the lifetime of the thread. For example, a thread may be configured to run only on cores 0 through 3, all of which belong to quad core CPU package 0. Henceforth, the scheduler will choose among these alternatives without migrating the thread to another package.

Exercising processor affinity insures that memory allocations remain local to the thread(s) that need them. Several downsides, however, should be noted. In general, processor affinity may significantly harm system performance by restricting scheduler options and creating resource contention when better resources management could have otherwise been used. For example, affinity restrictions may prevent the scheduler from assigning waiting threads to unutilized cores during a particular interval. Or, low priority threads may adversely impact high priority threads due to affinity restrictions that prevent adjustments through the use of additional cores. Processor affinity restrictions may even hurt the application itself when additional execution time on another node would have more than compensated for a slower memory access time.

Such downsides imply the need to think carefully about whether processor affinity solutions are right for a particular application and shared system context. Note, finally, that processor affinity APIs offered by some systems support priority "hints" and affinity "suggestions" to the scheduler in addition to explicit directives. Such suggestions may insure optimal performance in the common case yet avoid constraining scheduling options during periods of high resource contention.

3.2. Data Placement Using Implicit Memory Allocation Policies

In the simple case, many operating systems transparently provide support for NUMA-friendly data placement. When a single-threaded application allocates memory, the processor will simply assign memory pages to the physical memory associated with the requesting thread's node (CPU package), thus insuring that it is local to the thread and access performance is optimal.

Alternatively, some operating systems will wait for the first memory access before committing on memory page assignment.² To understand the advantage here, consider a multi-threaded application with a start-up sequence that includes memory allocations by a main control thread, followed by the creation of various worker threads, followed by a long period of application processing or service. While it may seem reasonable to place memory pages local to the requesting thread, in fact, they are more effectively placed local to the worker threads that will access the data. As such, the operating system will observe the first access request and commit page assignments based on the requester's node location.

These two policies together illustrate the importance of an application programmer being aware of the NUMA context of the program's deployment. If the page placement policy is based on first access, the

programmer can exploit this fact by including a carefully designed data access sequence at startup that will generate “hints” to the operating system on optimal memory placement. If the page placement policy is based on requester location, the programmer should insure that memory allocations are made by the thread that will subsequently access the data and not by an initialization or control thread designed to act as a provisioning agent.

Multiple threads accessing the same data are best co-located on the same node so that the memory allocations of one, placed local to the node, can benefit all. This may, for example, be used by prefetching schemes designed to improve application performance by generating data requests in advance of actual need. Such threads must generate data placement that is local to the actual consumer threads for the NUMA architecture to provide its characteristic performance speedup.

It should be noted that when an operating system has fully consumed the physical memory resources of one node, memory requests coming from threads on the same node will typically be fulfilled by sub-optimal allocations made on a remote node. The implication for memory-hungry applications is to correctly size the memory needs of a particular thread and to insure local placement with respect to the accessing thread.

For situations where a large number of threads will randomly share the same pool of data from all nodes, the recommendation is to stripe the data evenly across all nodes. Doing so spreads the memory access load and avoids bottleneck access patterns on a single node within the system.³

3.3. Data Placement Using Explicit Memory Allocation Directives

Another approach to data placement in NUMA-based systems is to make use of system APIs that explicitly configure the location of memory page allocations. An example of such APIs is the `libnuma` library for Linux.¹

Using the API, a programmer may be able to associate virtual memory address ranges with particular nodes, or simply to indicate the desired node within the memory allocation system call itself. With this capability, an application programmer can insure the placement of a particular data set regardless of which thread allocates it or which thread accesses it first. This may be useful, for example, in schemes where complex applications make use of a memory management thread acting on behalf of worker threads. Or, it may prove useful for applications that create many short-lived threads, each of which have predictable data requirements. Pre-fetching schemes are another area that could benefit considerably from such control.

The downside of this scheme, of course, is the management burden placed on the application in handling memory allocations and data placement. Misplaced data may cause performance that is significantly worse than default system behavior. Explicit memory management also presupposes fine-grained control over processor affinity throughout application use.

Another capability available to the application programmer through NUMA-based memory management APIs is memory page migration. In general, migration of memory pages from one node to another is an expensive operation and something to be avoided. Not only is there the cost of migrating the data, but all associated memory references must be discovered and modified to observe the new mapping. As the remapping is taking place, pages must be removed from operating system page lists and detached from normal swapping mechanisms.

Having said this, given an application that is both long-lived and memory intensive, migrating memory pages to re-establish a NUMA-friendly configuration may be worth the price.³ Consider, for example, a long lived application with various threads that have terminated and new threads that have been created but reside on another node. Data is now no longer local to the threads that need it and sub-optimal access requests now dominate. Application-specific knowledge of a thread’s lifetime and data needs can be used to determine whether an explicit migration is in order.

Finally, the API may provide functions for obtaining page residency or for examining memory access behavior under the current configuration. Such tools may provide the means to implement a monitoring scheme that makes explicit migration adjustments when memory accesses within the NUMA context fall below a defined threshold.

Summary

NUMA, or **Non-Uniform Memory Access**, is a shared memory architecture that describes the placement of main memory modules with respect to processors in a multiprocessor system. The advantage of the NUMA architecture as a *hierarchical* shared memory scheme is its potential to improve *average case access time* through the introduction of fast, local memory. To realize the potential of NUMA systems, however, careful *data placement* is needed. The more data can effectively be placed in memory local to the processor that needs it, the more overall access time will benefit from the architecture.

In this brief technical paper, we have described various strategies and considerations for ensuring optimal data placement within a NUMA-based system. In particular, we have discussed the role of processor affinity, memory allocation strategies that use implicit operating system page placement policies, and the use of the system API's for assigning and migrating memory pages using explicit directives.

References

- ¹ Drepper, Ulrich. "What Every Programmer Should Know About Memory". November 2007.
- ² Intel® 64 and IA-32 Architectures Optimization Reference Manual. See Section 8.8 on "Affinities and Managing Shared Platform Resources". March 2009.
- ³ Lameter, Christoph. "Local and Remote Memory: Memory in a Linux/NUMA System". June 2006.

Author Bio

David E. Ott is a Senior Software Engineer with Intel's Software Solutions Group. He joined Intel in 2005 as a middleware systems engineer for the Technology and Manufacturing Group. Currently, David focuses on power and virtualization aspects of enterprise server platforms. David holds M.S. and Ph.D. degrees in Computer Science from the University of North Carolina at Chapel Hill.