



Performance Characteristics of Explicit Superpage Support

Mel Gorman, Patrick Healy

► To cite this version:

Mel Gorman, Patrick Healy. Performance Characteristics of Explicit Superpage Support. WIOSCA 2010 - Sixth Annual Workshop on the Interaction between Operating Systems and Computer Architecture, Jun 2010, Saint Malo, France. inria-00493770

HAL Id: inria-00493770

<https://inria.hal.science/inria-00493770>

Submitted on 21 Jun 2010

HAL is a multi-disciplinary open access archive for the deposit and dissemination of scientific research documents, whether they are published or not. The documents may come from teaching and research institutions in France or abroad, or from public or private research centers.

L'archive ouverte pluridisciplinaire **HAL**, est destinée au dépôt et à la diffusion de documents scientifiques de niveau recherche, publiés ou non, émanant des établissements d'enseignement et de recherche français ou étrangers, des laboratoires publics ou privés.

Performance Characteristics of Explicit Superpage Support

Mel Gorman

IBM / University of Limerick
mel@csn.ul.ie

Patrick Healy

University of Limerick
patrick.healy@ul.ie

Abstract

Many modern processors support more than one page size. In the 1990s the larger pages, called superpages, were identified as one means of reducing the time spent servicing Translation Lookaside Buffer (TLB) misses by increasing TLB reach. Transparent usage of superpages has seen limited support due to architectural limitations, the cost of monitoring and implementing promotion/demotion, the uncertainty of whether superpages will be a performance boost and the decreasing cost of TLB misses due to hardware innovations. As significant modifications are required to transparently support superpages, the perception is that the cost of transparency will exceed the benefits for real workloads.

This paper describes how processes can explicitly request memory be backed by superpages that is cross-platform, incurs no measurable cost and is suitable for use in a general operating system. By not impacting base page performance, a baseline metric is established that alternative superpage implementations can compare against. A reservation scheme for superpages is used at `mmap()` time that guarantees faults without depending on pre-faulting, the fragmentation state of the system or demotion strategies. It is described how to back different regions of memory using explicit superpage support without application modification and present an evaluation of an implementation running a range of workloads.

Categories and Subject Descriptors D.4 [Storage Management]: Virtual Memory

General Terms Algorithms, Measurement, Performance

Keywords superpage, fragmentation, replacement policy, TLB

1. Introduction

In the 1990s, it was observed that the Translation Lookaside Buffer (TLB) reach was shrinking as a percentage of physical memory and the relative amount of program time spent servicing TLB misses was increasing. Superpages were identified as one means of reducing TLB misses. Preferably it would be fully transparent to the application but there are significant obstacles that prevent transparent superpage support being properly adopted in mainstream operating systems for general applications. This paper begins by describing architectural limitations when translating addresses for multiple page sizes that have hindered transparent superpage support. As superpages typically require physically contiguous pages, the paper

summarises some earlier work on controlling external fragmentation and page reclaim [GOR08] as they are important pre-requisites for supporting superpages [TALLURI92].

This paper then progresses this earlier work by describing a superpage reservation scheme that can be applied at `mmap()` time to guarantee future faults. This eliminates the need for pre-faulting or demotion strategies that are problematic on a number of architectures. These are prerequisites that must be addressed for reliable Explicit Superpage Support; a mechanism that exposes interfaces directly allowing loaders, libraries or applications to directly request superpages. It is shown that explicit superpage support can back memory regions without direct application modification through the use of libraries.

Over the course of the research, superpage support was implemented on Linux and merged into the mainline releases. This paper presents an evaluation of an implementation based on Linux using AMD® Phenom (X86-64) and POWER® (PPC64) processors. By tackling the fundamental problems of superpage support and enabling usage without application modification, this implementation represents a comparison point for transparent superpage support implementations on the same platform. There is an expectation that similar mechanisms could be implemented on other platforms.

2. Limitations of Transparent Support

There are architectural limitations to consider when implementing superpage support that is cross-platform. Limiting the number of architecture-specific paths reduces maintenance and test costs so the common architectural features must be supported before considering architecture-specific options. This section considers three architectures that make up 99.6% of the Top 500 Supercomputers [MSSD08] list as being a reasonable cross-section of mainstream architectures. Understanding the limitations of two of these architectures in particular illustrate why fully transparent superpage support be problematic for cross-platform operating systems.

2.1 POWER®

POWER® occupies 3 of the top 10 spots in the Top 500 list, has a 12% share overall and occupies the top spot as part of a hybrid system [MSSD08]. The current Power Architecture specification allows multiple page sizes but initially it only allowed two; a base and an implementation-specific large page size between 8KB and 256MB [POWER03]. With POWER4, the page sizes were 4K and 16MB with POWER5+ adding 64K and 16GB superpages. Unfortunately, in general¹ the Memory Management Unit (MMU) requires the same superpage size be used within segments of either 256MB or 1TB virtual address range, thus hindering transparent superpage support. Promotion or demotion would require modifying the entire segment at once with is prohibitively expensive and makes transparent superpage support on POWER impractical.

¹ The specification allows 4K and 64K pages to mix within a segment but support is not universal.

Permission to make digital or hard copies of all or part of this work for personal or classroom use is granted without fee provided that copies are not made or distributed for profit or commercial advantage and that copies bear this notice and the full citation on the first page. To copy otherwise, to republish, to post on servers or to redistribute to lists, requires prior specific permission and/or a fee.

WIOSCA '10 June 19, 2010, Saint-Malo, France.
Copyright © 2010 ACM [to be supplied]...\$5.00

2.2 Itanium® (IA-64)

Itanium maintains a 1.8% share of the Top 500 list and is used on very large single-image machines such as the SGI Altrix 4700 supporting up to 1024 processors on a single machine. The Itanium supports 11 different page sizes but the Itanium port of Linux has similar problems to POWER with respect to address translation that hinder transparent support.

Two mutually exclusive hardware-defined formats usable by page table walkers exist for Itanium called short- and long-format Virtual Hashed Page Table (VHPT); SF-VHPT and LF-VHPT respectively. SF-VHPT has a smaller cache-footprint but can only use one page size per region. As the 64-bit address space is divided into eight regions, this is a significant restriction as multiple pages can only be used within a region if the hardware page table walker is disabled - a significant penalty. The LF-VHPT has larger PTEs and is less cache friendly but potentially has a smaller TLB footprint and has no restrictions on superpage placement.

Linux favours SF-VHPT for two reasons that may be important to other platforms. Firstly, TLB-and-context-switch-intensive workloads are relatively rare and Linux has a policy of minimising cache usage. Secondly, the “hardware-independent” representation of page tables used by Linux is a 4-level hierarchical page table format which, if mapped in the correct place, can be used directly by the SF-VHPT hardware walker. LF-VHPT can still be used on Linux [WIENARD08] but with more expensive TLB reloads. As this can only be potentially offset by superpage usage, the move to the long-format is currently unjustified on that system. This severely limits the possibility of transparent support on the Itanium platform.

2.3 X86 Variants (Intel® EM64T, AMD® X86-64)

X86 variants dominate the Top 500 list with a share of 73.8% for Intel EM64T and 12% for AMD X86-64. 2MB or 4MB superpages are supported depending on the address translation used. In 32-bit mode, 4MB superpages are typically available unless Physical Address Extension (PAE) is enabled. If PAE or 64-bit mode is used, the superpage size is generally 2MB but some recent AMD processors also add a 1GB superpage size.

Due to how addresses are translated, there are limits on how many superpage sizes can be supported but no significant limit on where the superpages are placed. The difficulties for transparent support are two-fold, one hardware and one software. The TLB characteristics differ significantly between processor implementations. The Intel Pentium D, including Xeon variants, have fully-associative TLBs. Others such as the Core Duo have set-associative TLBs with different numbers of entries for different page sizes. The AMD Phenom 9950 processor reports multiple TLBs with the L1 DTLB fully-associative and L2 DTLB set-associative. This degree of variation significantly complicates cost/benefit calculations regarding superpages.

The software problem is specific to Linux “open-coding” 4-level page table traversal without checking the underlying page size instead of using an abstracted API. Transparent support is possible but implementing it incurs a performance cost. Before that cost is incurred, it must be established what the limits of explicit superpage support are.

3. Related Work

This section summarises some earlier work [GOR08] that is a prerequisite for explicit superpage support. Architectural limitations generally require that superpages be naturally aligned and physically contiguous forcing the system to handle external fragmentation. Memory defragmentation is not a complete solution as wired pages exist that cannot move and may be scattered throughout the

physical address space. Being able to move arbitrary pages is costly due to locking requirements and copying, both of which vary between architectures.

3.1 Page Allocation

Grouping Pages By Mobility [GOR08] divides the physical address space into superpage-sized arenas. Different arena types exist based on a page’s ability to move, be it by migration or reclamation noting that the number of arenas for each type changes based on demand. Each arena has its own free-lists so the search cost is similar to that of the standard allocator. The placement policy does not guarantee that a specific number of superpages can be allocated so the system’s memory may also be partitioned. Arenas for wired pages are only ever created on one side of the partition guaranteeing that one partition can be filled with superpages by reclaiming or moving pages. On Non-Uniform Memory Access (NUMA) systems, the partition for wired pages is evenly divided on all nodes to avoid kernel accesses being frequently remote. On batch systems with workloads of varying requirements, the partition would be sized to the largest estimated number of superpages required. If superpages are not allocated, movable base pages still use the partition.

3.2 Page Reclaim

On allocation failure, Age-Biased Arena Reclaim [GOR08] uses the normal page reclaim algorithm to select a reference page and reclaim the naturally aligned arena if all pages are free or reclaimable as likely will be the case due to the placement policy. This minimises disturbance of the normal page reclaim decisions made by the operating system

4. Superpage Reservation

Systems normally ensure successful faults by recording the number of pages required for existing mappings and refusing to create new mappings that exceed available memory and swap space. This is insufficient for superpages as contiguity requirements can fail a page fault due to external fragmentation. Faulting superpage-backed mappings at `mmap()` time would increase the cost of `mmap()` and perform poorly with Non-Uniform Memory Access (NUMA). Optimally on NUMA, CPUs always access pages local to the CPU without copying data between nodes. Pre-faulting places pages local to the thread calling `mmap()` and not necessarily local to the threads accessing the data, incurring costly cross-node accesses.

A better alternative is to reserve the superpages necessary to guarantee future faults at `mmap()` time. Counters track the current state of the system, some of which are tracked on a per-NUMA-node basis so that NUMA policies can be applied correctly.

When creating a new mapping, faulting a page within an existing mapping or unmapping an existing mapping, it is determined how many pages have already been allocated or reserved. For example, pages for shared mappings may already have been allocated or a reserve created by another process. Searching the page cache or page tables for references to allocated pages would be unsuitably expensive. Instead, each mapping is treated as a file. Shared mappings of data use one file whereas private mappings have their own unique file. Each file is associated with one or more `struct file_region` stored on a linked list and defined as follows.

```
struct file_region {
    struct list_head link;
    long from, to;
};
```

The size of the file in superpages is the initial required reservation to succeed all future faults. Each `file_region` on the list implies that to - from reserved pages have been used.

Due to the fact that shared and private mappings are accounted differently, there are differences in the run-time behaviour of an application using shared or private mappings backed by superpages.

4.1 Shared Mapping Accounting

Shared mappings reservations for the same data are based on the largest mapping created by any process i.e. the size of the file. The `file_region` is associated with one shared structure such as an inode. On fault, the page cache is searched and the page inserted into the processes page tables if it exists, otherwise the page is allocated, inserted into both page cache before updating the reservations, `file_region` structures and inserting into the page table.

4.2 Private Mapping Accounting

Private mappings always reserve superpages for the entire mapping as even existing pages require new pages for a Copy-On-Write (COW) fault. Read-only private mappings could avoid reserving until `mprotect()` was called but applications may not correctly handle `mprotect()` errors if the reserve could not be made. Due to the per-mapping nature of the reservation, the `file_region` is linked to the virtual area management structure (`struct vm_area_struct` in Linux).

`fork()` presents a difficulty. Shared mappings can use the same pages and reservations between parent and child but it is unknown how many superpages a child will need for private mappings. Duplicating reservations on `fork()` is expensive, may be impossible and impacts applications calling `exec()` after `fork()`.

An alternative is to only guarantee faults within the address space that created the mapping. A child process's mappings have no reservations but on fault, a superpage allocation without a reservation is attempted. If successful, the fault completes or else the child gets killed. A COW fault from the parent similarly attempts to allocate a new superpage. On allocation failure, the page is unmapped from any children and the mappings are flagged. The parent safely uses the original page but child gets killed if it faults within the flagged mapping. Hence, a programming model based on superpage should not assume reliable superpage behaviour for private mappings in child processes. Either private mappings should be `madvise()` with `MADV_DONTFORK` or child processes should take care not to access the mappings. This model guarantees the address space owner can successfully fault without depending on in-kernel demotion strategies that cannot be implemented for all architectures due to limitations on superpage placement. In the event a child is killed, a message is logged stating that "PID N killed to to inadequate hugepage pool" to make it clear what occurred.

5. Explicit Superpage Support

A simple interface for superpages would be similar to `brk()` or `malloc()` but it is difficult to create a shared mapping without extending the interface with unique semantics. The System V shared memory interface exists but application error can leak memory and private mappings are impossible to create. This section describes an interface for an application to map superpages without transparent support.

5.1 RAM-Based Filesystem

The Virtual File System (VFS) API for the creation, mapping and unmapping of files, both shared and private is well understood by developers. A RAM-based filesystem can create files that are backed by superpages when memory mapped. The filesystem API

does not allow the passing of additional information but each mount of the filesystem can support different page sizes. Alternatively, an `fadvise()`-like interface could specify the page size before `mmap()` is called. The filesystem is a low-level interface and more suitable for use by libraries than directly by applications.

5.2 System V Shared Memory

The semantics of shared memory created and attached with the System V interfaces `shmget()` and `shmat()` are similarly well understood. By passing a new flag, `SHM_HUGETLB`, the kernel is instructed to back the region with superpages. Internally, files of the size requested by `shmget()` can be created on a private instance of the filesystem. `shmat()` maps the file similar to a normal file mapped as `mmap(MAP_SHARED)`, thus sharing code between the filesystem and shared memory implementations. An error should be returned if the size is not superpage-aligned or sufficient superpages, leaving the decision on whether to use base pages to the library or application.

5.3 Anonymous mmap() Mappings

Memory mappings created with `mmap()` take a number of flags such as `MAP_PRIVATE`, `MAP_ANONYMOUS` and others described in the manual page. Similar to the previous interfaces, `mmap()` can be extended to support a `MAP_HUGETLB` flag which automatically backs the mapping with superpages.

The difficulty with this interface is that setting the page size after `mmap()` is not an option due to possible architectural limitations and setting a default page size to be used for `mmap()` can race when threads are in use. If the architecture does not limit the placement of superpages, then `madvise()` can be extended after `mmap()` and before the pages are faulted. Otherwise, using `MAP_HUGETLB` implies that a default superpage size is used.

5.4 Explicit Programming API

Three functions are defined that are useful to applications or libraries with superpage-awareness. The naming reflects the Linux implementation and uses the term "huge page" for superpages.

`hugelbfs.unlinked_fd()` creates a file on the RAM-based filesystem for the requested superpage size and unlinks, rendering it invisible to other processes. This is useful when a process needs complete control over the size of the mapping placement within the address space.

`get_huge_pages()` takes a length and flags parameter. The length must be superpage-aligned and return a buffer to a superpage-aligned length of memory or NULL on failure. It is suitable for use in the implementation of custom allocators. It is not suitable for use as a drop-in replacement for `malloc()` as it is a raw interface to superpages that does not cache-color buffers.

`get_hugepage_region()` similarly takes a length and flags parameter. The length need not be aligned and wastage is used to cache-color the buffer. On allocation failure, a buffer backed by base pages may be returned if requested by the caller. This is suitable as a drop-in replacement requiring that an application convert call-sites that create large buffers to use `get_hugepage_region()` instead.

6. Backing Memory Sections with Superpages

This section describes how different types of mappings of a process can be backed by superpages using explicit superpage support.

6.1 Heap

Allocations increasing the size of a heap currently use either `brk()` or `mmap()` to create memory regions. A custom allocator should use `get_hugepage_region()` and divide the superpage up into

buffers as appropriate. This requires source modification which may be undesirable but it is a relatively straight-forward exercise.

If the allocator provide hooks for applications to create more heap then `get_huge_pages()` can be used by pre-loading and registering the hook before the application starts. GNU libc uses `sbrk()` for small allocations and `mmap()` for blocks larger than `MMAP_THRESHOLD` which is 128kB by default. A `morecore` hook allows registry of a custom function that allocate more heap.

6.2 Mapping Text/Data/BSS

The effectiveness of backing data section with superpages depends on the compiler and the application. If the data section contains pre-initialised data, then backing it with superpages will be effective. If the data section is zeros but stored on disk, then the copy will be effective. If it uses `malloc()` with a hook, then it can be backed similar to the heap. Otherwise, the loader must be modified.

Text and data are mapped with different permissions and the protection granularity depends on the page size. Applications should be relinked with the sections superpage-aligned which requires a suitable linker tool chain. Otherwise, architectural limitations determine if sections can be partially backed by superpages.

On `mmap()` of a section, the kernel can check the page cache. If the data exists but is stored in an unaligned base page, then the data should be copied to a superpage and the existing page cache entry replaced so that only one copy exists. If a superpage cannot be allocated, then the base pages are used. As the checks are made before the insertion in the process page tables, architectural limitations with address translation are avoided. Alternatively the loader can use `get_hugepage_region()` or similar to allocate a suitably large region, copy the data from disk and remap it to a suitable place. It is implementation-specific as to whether text is shared between processes or not.

6.3 Stack

The stack area of memory has similar requirements to `brk()` but is not as easily handled with hooks like `malloc()` or copied and remapped like text and data sections. The maximum stack size is unknown at application start time making it inherently unreliable to back with superpages without a demotion strategy. Consequently, explicit superpage support requires the stack be a fixed size.

7. Evaluation

Explicit superpage support as described in this paper has been implemented over the course of a number of kernel releases. This evaluation was based on Linux Kernel 2.6.31, downloaded from <http://www.kernel.org> and configured using the distribution's kernel configuration file. `libhugetlbfs` 2.7 was used to transparently access superpages without application modification. The distribution used was Debian GNU/Linux Lenny current as of October 2009. This distribution is freely available for independent verification and is commonly used as a server operating system. The compiler used was gcc 4.3.2.

Performance results are based on test-runs without profiling as it can distort results due to increased interrupt handling from the PMU and the profiling daemon gathering data. Profile information was gathered during a separate using `OProfile`². Due to limitations of the PMU, DTLB and cache misses data were collected during separate runs. In all cases, the relative performance of base pages and superpages matched so data collected during profiling should be indicative of what happened during no-profile runs.

During the base page runs, no superpages were used. During the superpage runs, the maximum number of superpages that was possible to use were allocated. Depending on the workload, this varied

from negligible amounts to almost the size of physical memory. As the underlying implementation does not support the paging of superpages, there are no eviction rates to report but to make the comparison fair, it was confirmed that there was no meaningful paging activity while the tests were running with base pages.

The two machines were chosen based on the parts or models being readily available and are described in Table 1. Each workload was first run with base pages, then with superpages. They were rerun with profiling before moving onto the next workload. All superpages were released for each base page run and reallocated for superpage runs. All tests ran without rebooting without any failures to allocate superpages.

7.1 STREAM (Memory Throughput)

STREAM [MCCALPIN07] is a synthetic memory bandwidth benchmark that measures the performance of four long vector operations: Copy, Scale, Add and Triad. It can be used to calculate the number of floating point operations performed during the benchmark to estimate the cost of the "average" memory access. Simplistically, more bandwidth is better.

The C version of the benchmark was selected and used three statically allocated arrays for calculations. Modified versions of the benchmark using `malloc()` and `get_hugepage_region()` were found to have similar performance characteristics.

The benchmark has two parameters: `N`, the size of the array and `OFFSET`, the number of elements padding the end of the array. A range of values for `N` were used to generate workloads between 128K and 2GB. For each size of `N` chosen, the benchmark was run 10 times and an average taken. The benchmark is sensitive to cache placement and optimal layout varies between architectures as noted on the benchmark author's homepage. Where the standard deviation of 10 iterations exceeded 5% of the throughput, `OFFSET` was increased to add one cache-line of padding between the arrays and the benchmark for that value of `N` reran. High standard deviation were only observed when the total working set was around the size of the L1, L2 or all caches combined.

The benchmark avoids data re-use, be it in registers or in the cache. Hence, benefits from superpages would be due to fewer faults, a slight reduction in TLB misses as fewer TLB entries are needed for the working set and an increase in available cache as less translation information needs to be stored.

To use superpages, the benchmark was first compiled with the `libhugetlbfs ld` wrapper to align the text and data sections to a superpage boundary [LIBHTLB09]. It was then loaded with `hugectl --text --data --no-preload`.

Figure 1 shows the comparison for the different STREAM operations when using base pages and superpages. On the test machines, there were massive reductions in DTLB misses. In the case of PPC64, it is reported as a 100% reduction but due to limitations to OProfile, TLB misses were not reduced to zero but there were so few that OProfile did not sample a miss.

X86-64 gained 3.5% in the Add operation but was offset by a 3.31% regression in the Triad operation. Again, DTLB misses were significantly reduced and the time spent servicing misses was reduced from 1.84% to close to 0%. Cache misses were also reduced overall. The differences may again be potentially explained by cache coloring with Triad suffering significant conflicts.

On STREAM, data is not often re-used so increased cache capacity does not necessarily help STREAM unless the data is pre-fetched and cache conflicts are not an issue. As cache misses were not significantly reduced on X86-64, it implies that data was not being accurately pre-fetched or there were significant cache conflicts. Further investigation showed that there were 60 times more cache misses than TLB misses implying that cache conflicts

² <http://oprofile.sourceforge.net/news/>

CPU	AMD Phenom® 9950 Quad-Core	CPU	PPC970MP, altivec supported
CPU Frequency	1.3GHz	CPU Frequency	2.5GHz
# Physical CPUs	1 (4 cores)	# Physical CPUs	2 x dual core (4 cores in all)
L1 Cache	64K Data, 64K Instruction per core	L1 Cache	32K Data, 64K Instruction per core
L2 Cache	512K Unified per core	L2 Cache	1024K Unified per core
L3 Cache	2048K Unified Shared per chip	L3 Cache	N/a
Main Memory	8 GB	Main Memory	10GB
Mainboard	Gigabyte GA-MA78GM-S2H	Mainboard	Specific to the machine model
Superpage Size	2MB (1GB available but unused)	Superpage Size	16MB
Est. TLB Miss Cost	51 cycles	Est. TLB Miss Cost	563 cycles
Machine Model	Custom built	Machine Model	Terrasoft YDL Powerstation
	X86-64		PPC64

Table 1. Machine Configurations Used For Evaluation

were a greater problem than TLB misses for this workload in the chosen configuration.

In contrast, PPC64 has a very high TLB miss cost and the reduction in misses resulted in large performance gains of between 11.27% for the Copy operation to 16.57% for Triad. Further, an estimated 18.23% of time was spent servicing TLB misses for the STREAM benchmark. As cache misses only occurred 3-4 times more than TLB misses, it made TLB misses a more dominant factor than the other test machines. Cache misses were also significantly reduced, possibly helped by the fact that the arrays were rarely aligned to a superpage boundary.

The results show that TLB misses are not always a significant bottleneck for workloads with poor cache locality. In the situation where data structures are aligned to a superpage size, superpages can in fact impair performance due to increased cache conflicts depending on the CPU implementation. For the STREAM benchmark in particular, there is a potential for between 1.8% and 5% of performance to be gained by using superpages on X86-64 machines but care should be taken in future tests to avoid aligning arrays to the superpage size. In the PPC64 case, significant performance increases can be achieved with superpages due to the high cost of a TLB miss.

7.2 SysBench (Database Workload)

OnLine Transaction Processing (OLTP) is a general class of workload where clients perform a sequence of operations whose end result must appear to be an indivisible operation. TPC-C³ is considered an industry standard for the evaluation of OLTP but requires significant capital investment and is extremely complex to setup. Sysbench⁴ is a system performance benchmark comprising file IO, scheduler, memory allocation, threading and OLTP benchmarks. The setup requirements are less complicated and supports database back-ends for MySQL, Postgres and Oracle. Postgres⁵ was used for this experiment on the grounds that it uses a shared memory segment similar to Oracle making it a meaningful comparison with a commercial database server. Both Sysbench 0.4.8 and Postgres 8.3.4 were built from source.

Postgres was configured to use a 756MB shared buffer, an effective cache of 150MB, a maximum of 6*CPUs were allowed to connect and the “update_process_title” turned off. Options that check-out, fsync, log or synchronise were turned off to avoid interference from IO. **pg_ctl** was invoked with `hugesctl --shm pg_ctl` when backing the shared memory segments superpages.

The system was configured to allow the postgres user to use superpages with `shmget()` (`hugetlb_shm_group`). Postgres uses

System V shared memory but it has no superpage awareness. To transparently use superpages, **pg_ctl** was invoked with `hugesctl --shm pg_ctl` to automatically back shared memory with superpages.

The Sysbench client did not use superpages as it would be expected that a database administrator has limited control of the clients. The table size was 10 million rows, read-only to avoid IO and the test type was “complex” described making each operation by the client a database transaction. Tests were run varying the number of clients accessing the database from one to four times the number of CPU cores in the system, a total of 16 clients for both systems.

Figure 2 shows the comparison of the number of transactions per second when using base pages and superpages. Overall, X86-64 shows a 2.81% improvement and 4.57% on PPC64 with reliable reductions in TLB, although not massive reductions on X86-64. Broadly speaking, the performance gains were close to predicted maximums of 2.3% for X86-64 and 3.63% for PPC64. Superpages performed slightly better than predicted due to the slight reduction in cache misses and remaining differences can be accounted for by small performance differences in performance when profiling. The difference between predictions and reality are so slight that there is a high degree of confidence that all potential performance improvements due to superpages are being achieved.

Anecdotal evidence in private communication with performance analysts presume that database workloads benefit from superpages by between 2% and 7% and these results would appear to correlate.

7.3 SPECcpu 2006 v1.1 (Computational)

SpecCPU 2006 v1.1⁶ is a standardised CPU-intensive benchmark used in evaluations for HPC that also stresses the memory subsystem. A `--reportable` run was made comprising “test”, “train” and three “ref” sets of input data. Three sets of runs compare base pages, superpages backing just the heap and superpages backing text, data and the heap. Only base tuning was used with no special compile options other than what was required to compile the tests. The benchmark suite classifies the run as “not publishable” on the grounds it did not recognise the `-O2 -m64` compiler switches but this is not considered significant.

Figure 3 shows the comparison of the time taken to complete the Integer portion of the benchmark when using base pages and superpages.

The difficulty in analysing the performance of this benchmark was that profiling introduced an unusually large amount of overhead in comparison to other benchmarks in this study. In comparison to other workloads, the number of PMU events recorded was

³ <http://www.tpc.org/tpcc/>

⁴ <http://sysbench.sourceforge.net/>

⁵ <http://www.postgresql.org/>

⁶ <http://www.spec.org/cpu2006/>

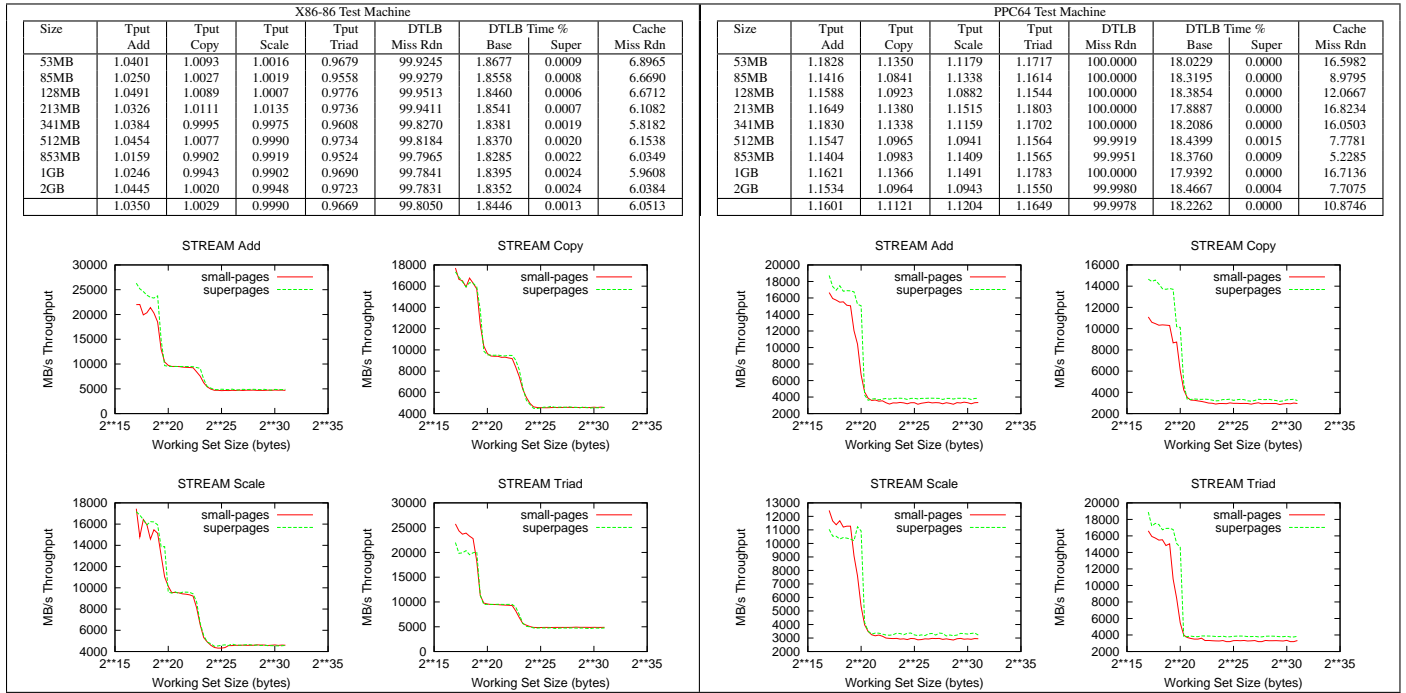


Figure 1. STREAM Comparison with Superpages

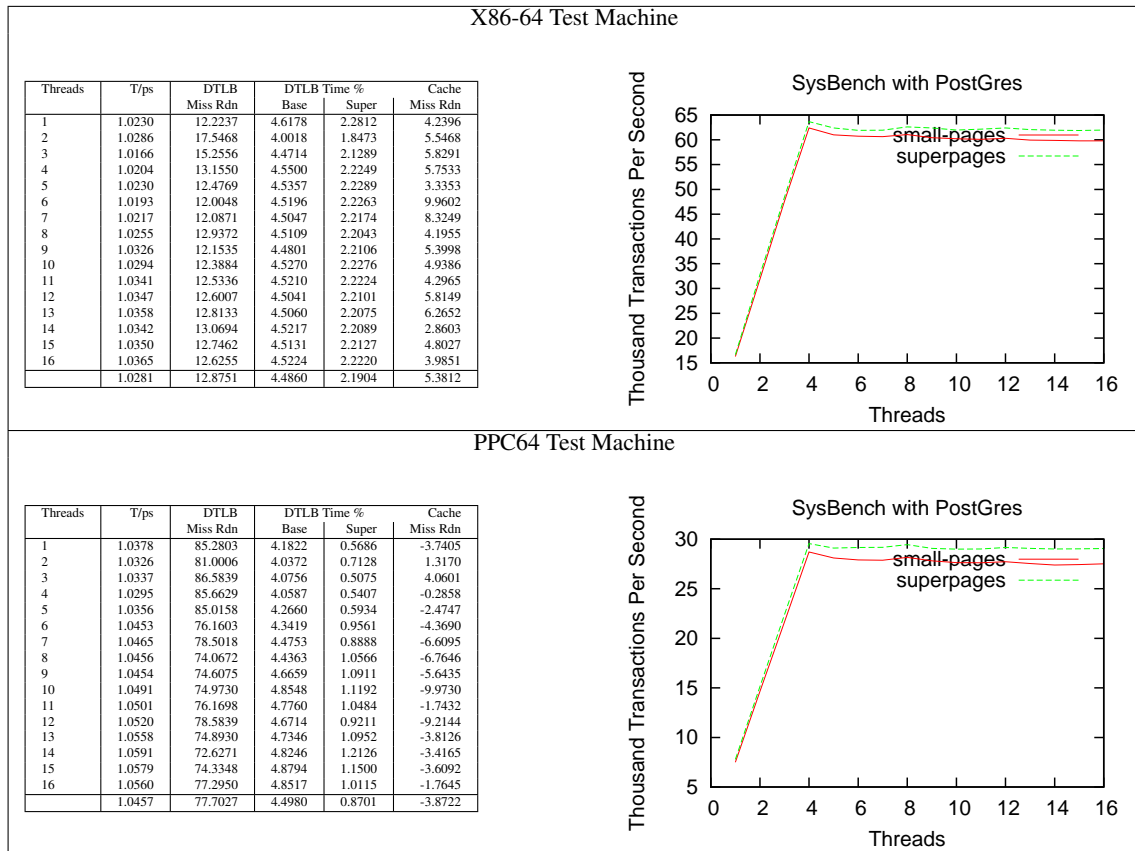


Figure 2. Sysbench Comparison with Superpages

far higher making the overhead of interrupts was correspondingly higher and a greater amount of CPU time was dedicated to the profiler than on other benchmarks. Unfortunately, reducing the granularity impaired the accuracy of the TLB measurements and could not be otherwise worked around on the available hardware. In all case where the actual performance exceeded the predicted performance, the difference was found in the profiling overhead.

X86-64 showed negligible difference in the number of TLB misses due to the erratic nature of the reference pattern for many of the tests. However, as the cost of a superpage TLB miss is often lower on that architecture due to the different number of pagetable levels, there can still be a performance gain in many cases. PPC64 in some cases reduced the number of TLB misses due to taking advantage of a larger superpage size but they were in the minority.

Despite the difficulties in the exact analysis of the reason for the performance gain, X86-64 gained 10.40% and PPC64 gained 15.14%.

Figure 4 shows the comparison of the time taken to complete the Floating Point portion of the benchmark when using base pages and superpages.

On X86-64 there was a negligible difference in the number of TLB misses reduced and gains were due to a lower cost of a superpage TLB miss. Similarly, the profiling overhead was unavoidable and made an accurate analysis of the TLB behaviour problematic but where actual performance gains differed from expected gains, a significant portion of the difference was accounted for by the profiling overhead. Overall, X86-64 gained 5.19% and PPC64 gained 8.20%.

7.4 SPECjvm 2008 (Java)

Java is used in an increasing number of scenarios, including real time systems and it dominates in the execution of business-logic related applications. Particularly within application servers, the Java Virtual Machine (JVM) uses large quantities of virtual address space that can benefit from being backed by superpages. SpecJVM 2008⁷ is a benchmark suite for Java Runtime Environments (JRE). According to the documentation, the intention is to reflect the performance of the processor and memory system with a low dependence on file or network I/O. Crucially for HPC, it includes SCIMark⁸ which is a Java benchmark for scientific and numerical computing.

The 64-bit versions of IBM Java Standard Edition Version 6 SP 3 were used for both X86-64 and PPC64. Installation on PPC64 required `libstc++-5` from outside of the distribution as the distribution uses `libstc++-6` and the backwards compatibility library was 32-bit only. The JVM was configured to use a maximum of 756MB for a heap. Unlike the other benchmarks, the JVM is superpage-aware and uses superpage-backed shared memory segments when `-Xlp` is specified to the JVM. Superpages were allocated in a static pool as required by the JVM. As tuning parameters were used, the `--peak` parameter was specified as required by the run rules.

Figure 5 shows the comparison of the number of transactions per second when using base pages and superpages. Overall, X86-64 shows a 7.57% performance improvement with 4.54% on PPC64 with reliable reductions in TLB misses. There are a wide range of results in the test categories with `scimarkLarge` suffering despite a prediction that tests would gain by between 2% and 21.2% on the three machines. An examination of the sub-tests that make up the `scimarkLarge` benchmark showed significantly different TLB behaviour and performance benefits with a mix of gains and losses complicating a general analysis from the context of superpages. In

general, the differences in the sub-tests account for the apparent difference between predicted and actual performance.

The disparity of the results show that superpages can help JVM workloads but by how much or little depends on the exact workload being executed within the JVM. As enabling superpages on JVMs is generally so trivial, the best recommendation is to test the workload with and without superpages measuring both absolute performance and profiling CPU and TLB miss counts.

8. Conclusions

This paper noted that transparent superpage support did not take architectural limitations into account and overly simplified cost functions. Explicit support is usable by 99.6% of the architectures making up the Top 500 supercomputer list, was implemented for the operating system running on 87.8% of them and does not necessarily require application modification. Evaluation was performed on two significantly different machines.

Average performance improvements due to explicit superpage use were very roughly in the range of 2% to 10% on X86-64 and 4% to 15% on PPC64 running a variety of workloads without imposing additional overhead on applications using exclusively base pages. PPC64 typically benefited to a greater degree than X86-64 which can be mostly accounted for by two facts. The first is that the cost of a TLB miss on the target machine was very high - 563 cycles versus 51 cycles for a miss on the X86-64 based machine. The second is that the larger superpage size on PPC64 as well as the larger number of TLB entries meant that the TLB reach was increased by a far greater degree on PPC64 than on X86-64.

It was not reported in the figures but superpage use reduced the number of page faults incurred by the system by a factor related to the difference in size between base and superpages. The significance of this is mitigated by the fact that the workloads faulted the bulk of their data early in the process lifetime and did not page heavily. I/O rates are affected during application start-up as superpages are being allocated but as this principally affected application initialisation, it was not found to have made any meaningful performance difference to the target workload.

The performance improvements are less dramatic than predicted in early literature but those estimations were based on software TLB miss handlers and in-order processors. The returns have diminished due to faster hardware pagetable walkers and speculative address translation. However, bottlenecks still exist and superpages can reduce TLB miss rates as well as cache miss rates due to less translation information being stored. This may be of particular importance to virtualised workloads where translation an address in a guest is significantly more expensive than native translation.

Explicit superpage provides an important tool for predictably mitigating the cost of address translation on a variety of platforms. It establishes a performance baseline for continued research on transparent superpage page and should have comparable performance to specialised OS kernels that enable superpage usage on a per-application basis.

Acknowledgments

The design and implementation of superpages support in Linux has a long history and would not have reached its current state without the sterling work of the Linux community, particularly Andy Whitcroft, William Irwin, Adam Litke, David Gibson, Nishanth Aravamudan, Peter Zijlstra, Kenneth Chen, Hugh Dickins, Nick Piggin, Steve Fox, Jon Tollefson, Andi Kleen and Eric Munson. Parts of the material presented in this work are supported by the Defense Advanced Research Projects Agency under its Agreement No. HR0011-07-9-0002.

⁷ <http://www.specbench.org/jvm2008/>

⁸ <http://math.nist.gov/scimark2/>

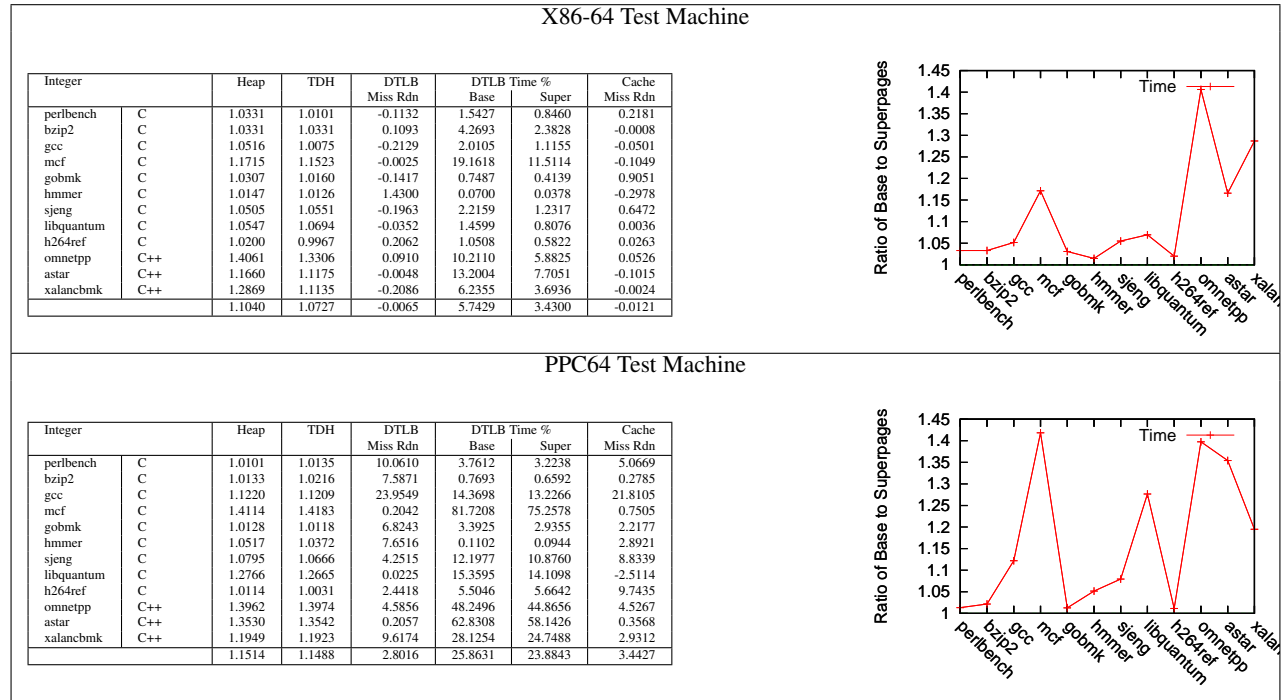


Figure 3. SPECcpu Integer Comparison with Superpages

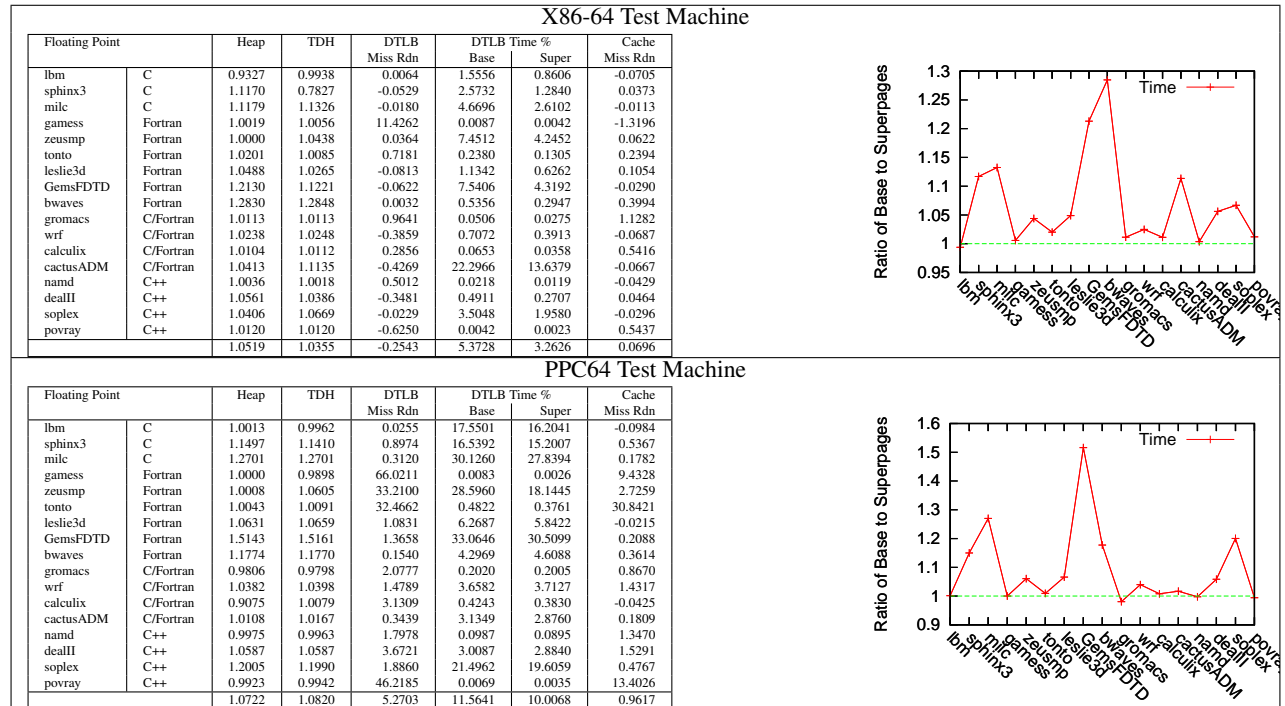


Figure 4. SPECcpu Floating Point Comparison with Superpages

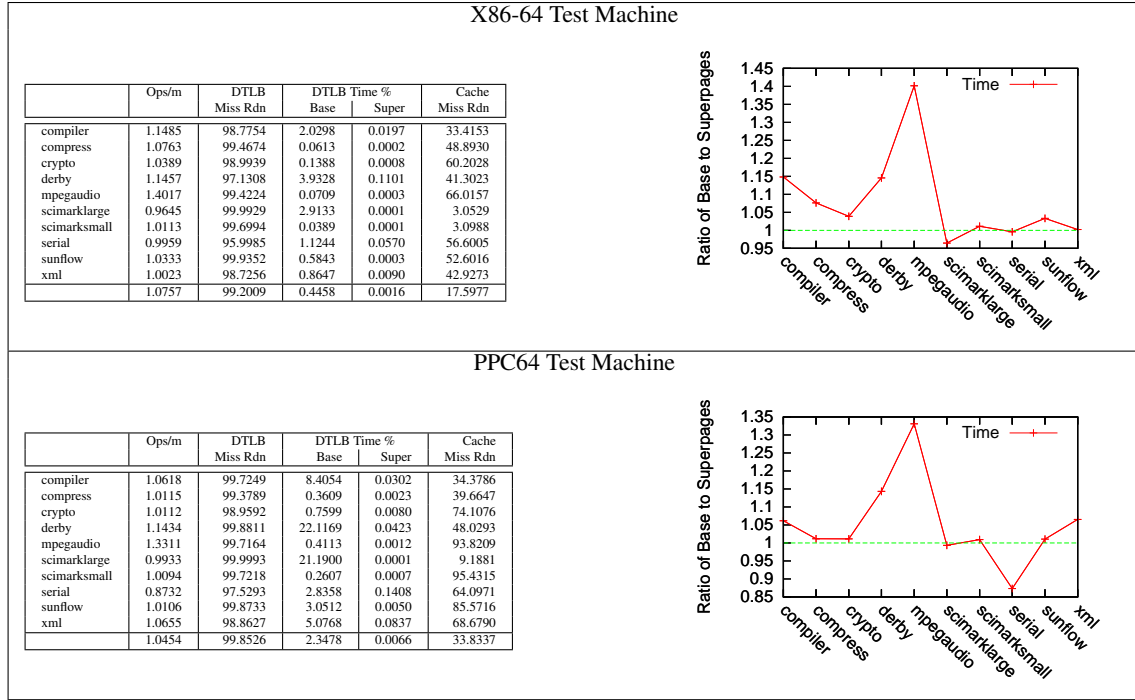


Figure 5. SPECjvm Comparison with Superpages

References

- [POWER03] Joe Wetzel, Ed Silha, Cathy May, Brad Frey. PowerPC Operating Environment Architecture Book III v2.01
- [MSSD08] Hans Meuer, Erich Strohmaier, Horst Simon and Jack Donarra. Top 500 Supercomputers List for November 2008. In <http://www.top500.org>
- [GOR08] Mel Gorman and Patrick Healy. Supporting superpage allocation without additional hardware support In *ISMM '08: Proceedings of the 7th international symposium on Memory management*,
- [TALLURI92] Madhusudhan Talluri and Shing Kong and Mark D. Hill and David A. Patterson and Sun Microsystems Laboratories. Trade-offs in supporting two page sizes In *Proceedings of the 19th Annual International Symposium on Computer Architecture*,
- [NAV04] Juan E. Navarro. Transparent operating system support for superpages PhD Thesis, Rice University, Houston, TX, USA, 2004.
- [CASCVAL05] Calin Cascaval and Evelyn Duesterwald and Peter F. Sweeney and Robert W. Wisniewski. Multiple Page Size Modeling and Optimization In *Proc. 14th International Conference on Parallel Architecture and Compilation Techniques*.
- [CWH03] Matthew Chapman, Ian Wienand, and Gernot Heiser. Itanium Page Tables and TLB Technical Report UNSW-CSE-TR-0307, School of Computer Science and Engineering, University of NSW, May 2003.
- [WIENARD08] Ian Wienard. Transparent Large-Page Support for Itanium Linux Masters Thesis, University of NSW, July 2008.
- [MCCALPIN07] John D. McCalpin. STREAM: Sustainable Memory Bandwidth in High Performance Computers In a continually updated technical report. <http://www.cs.virginia.edu/stream/>
- [GANAPATHY98] Narayanan Ganapathy and Curt Schimmel. General purpose operating system support for multiple page sizes In *Proceedings of the USENIX Conference. USENIX*
- [LIBHTLB09] Various Authors. libhugetlbfs 2.7 HOWTO In the libhugetlbfs source tarball. <http://sourceforge.net/projects/libhugetlbfs>
- [GOR08a] Mel Gorman, Using the Direct Hugepage Allocation API with STREAM In author's home site, <http://www.csn.ul.ie/~mel/docs/stream-api>