

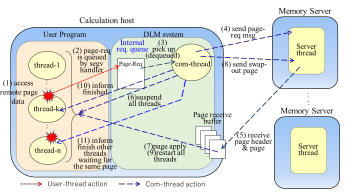
Efficient Swap Protocol for Remote Memory Paging in Out-of-Core Multi-Thread Applications

Hiroko Midorikawa, Kenji Kitagawa, Hikari Ohura, JST CREST & Seikei University, Tokyo Japan
midori@st.seikei.ac.jp http://www.ci.seikei.ac.jp/midori/paper

Abstract

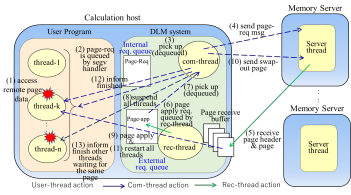
A new page swap protocol is proposed for user-level remote memory paging systems to accelerate the performance of out-of-core processing with multi-thread user programs and libraries written in OpenMP and pthread. The original swap protocol has a bottleneck in efficient page swapping, which is requested by multiple threads in a user program, because all the MPI communications to the memory servers and page swap managements are allocated to one system thread. However, this protocol has a benefit that it is widely available anywhere even if MPI thread-support-level is not so high. The new protocol uses two independent system threads: one for page swapping, and the other for managing user thread requests. Although it requires the highest MPI thread-support-level (multiple), which is usually considered to degrade the MPI communication performance compared to that in lower MPI thread-support-level, the new protocol achieves higher performance than the original protocol in micro benchmark, Stream benchmark, matrix multiplication, stencil computation, and 3-Dimensional FFT.

Original page swap implementation (base: r3-protocol)



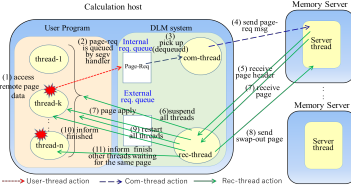
- Centralized control system by com-thread**
- No mutual exclusions for accessing the DLM internal system data.
 - Widely available anywhere, MPI thread-support-level is low. (FUNNELED, SERIALIZED).
 - Bottleneck in efficient page swapping. All requests from user threads are managed by a single com-thread, one by one.

Revised page swap implementation (revised 1: r58-protocol)



- Introducing receiver thread to the centralized control system (r3-protocol)**
- No mutual exclusions for accessing the DLM internal system data.
 - Concurrent MPI communications require MPI thread-support-level MULTIPLE.
 - Rec-thread dedicates receiving swap-in pages in receiver-buffers.
 - Com-thread manages sending page-requests, page-applying, and sending swap-out pages.

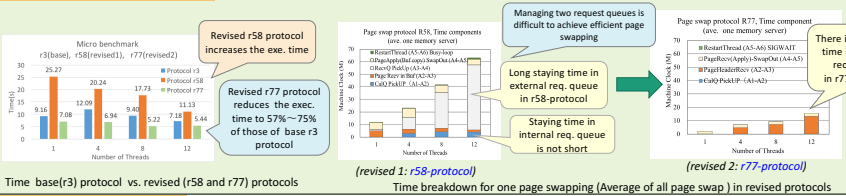
Revised page swap implementation (revised 2: r77-protocol)



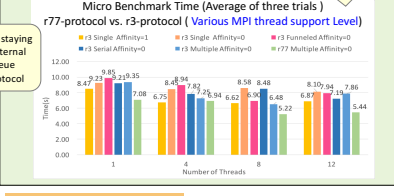
- Two system threads manage their works independently, without page-buffers.**
- No mutual exclusions are required to access the DLM internal system data.
 - Two MPI Communicating threads: requires MPI thread-support-level MULTIPLE.
 - Two system thread run independently. Page-buffers are removed.
 - Rec-thread manages all page swapping (receive, apply, send).
 - Com-thread manages only user requests.

Performance evaluations : base protocol vs. revised protocols

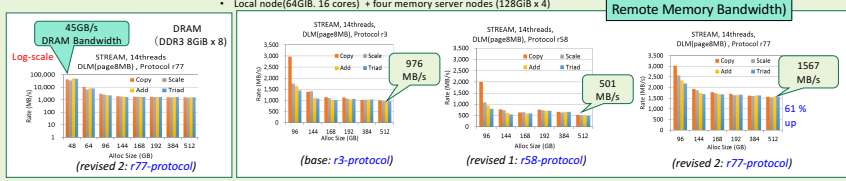
Micro Benchmark (Fig.3.)



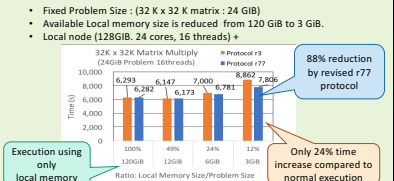
Micro Benchmark (Fig.3.)



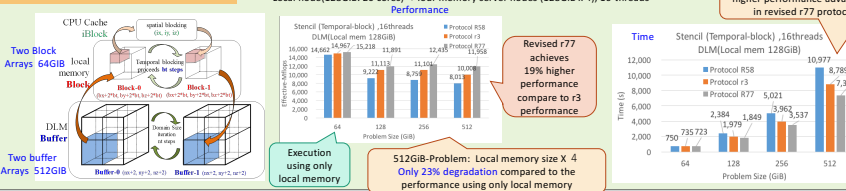
STREAM Benchmark [3]



Matrix Multiplication



7-point Stencil computation with Temporal blocking [4]

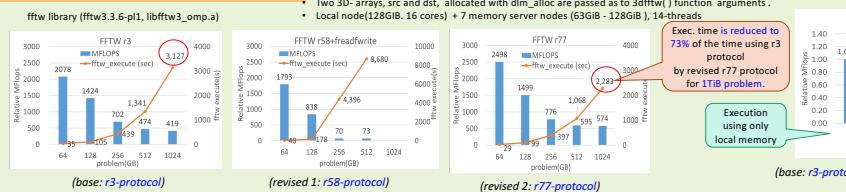


Result Summary

In all the cases, the proposed protocol (r77) achieved a better performance by using two system threads that are separately responsible for page swap and user thread request management and communicate to the memory servers independently. Comparison between revised r77 protocol and base r3 protocol:

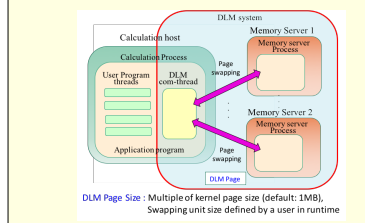
- Micro Benchmark: Exec. time 43% - 25% reduction
- Stream Benchmark: Remote Memory bandwidth 61% increases
- Stencil Computing: Exec. time of 512GiB-problem Local memory ratio 25% 19% reduction
- Matrix Multiplication: Exec. time of 24GiB-Problem Local memory ratio 25% 12% reduction
- 3D-FFT: Exec. time of 1TiB-Problem Local memory ratio 12.5% 12% reduction

3D Fourier Transform using fftw library [5]

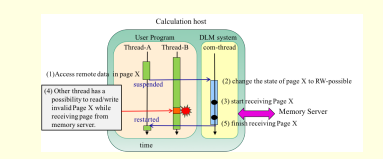


DLM (Distributed Large Memory) [1][2] for Out-Of-Core Multi-thread Program Executions

A user-level remote memory paging system, DLM (Distributed Large Memory), offers virtual large memory using a cluster of distributed node memories [1][2]. It is highly portable to various clusters and is an easy-to-use remote memory for people without any knowledge of MPI. It was designed for users who need to solve large-size problems that require the processing of large amounts of data beyond the capacity of local memories, using existing algorithms and programs originally designed for shared-memory models. They prefer to accept the extra execution time caused by partially using remote memory instead of the local memory, because converting existing code algorithms to parallel MPI programs is not an easy task and requires substantial costs. The DLM supports multi-thread programs and libraries written in OpenMP and pthread[1].



- DLM System :**
- Virtually provides a single large memory using distributed memories in a cluster for multi-thread programs designed for shared-memory-models, such as OpenMP and pthread.
 - A Transparent system to use applications : The segv-signal-handler invoked by each thread in an application fetches a remote page implicitly.
 - Little modification for existing programs
 - No need of MPI programming knowledge for using remote memories



- Fig.2 Difficulties in page swapping by multiple threads**
- Data inconsistency is generated when page fetch and data access are occurred simultaneously.
 - Temporal suspension of all user threads in a program during copying a remote page to the user address space is employed. pthread_kill(pthread_t, sig);
 - pthread_create() is hooked by the DLM system to register thread IDs of user threads which are generated and terminated dynamically in program execution.

The Programming interface for DLM

Convert malloc to **dim_alloc** and **dim_a[...]** allocations

```
#define MAX 1000
dim_a = (int (*)[MAX]) dim_alloc (sizeof(int) * MAX * MAX);
```

Use the **dimc translator** for array-data [6]

```
#define ENUM ((unsigned long) (1L<<34)) //16G Elements
int main (int argc, char *argv[])
{
    int i, j;
    for (i = 0; i < MAX; i++)
        a[i][i] = i * i;
    for (i = 0; i < MAX; i++)
        for (j = 0; j < MAX; j++)
            printf("%d", a[i][j]);
    return 0;
}
```

Converted by **dimc translator**

```
int (*dim_a)[MAX];
int main (int argc, char *argv[])
{
    int i, j;
    dim_startup (&argc, &argv);
    dim_a = (int (*)[MAX]) dim_alloc (sizeof(int) * MAX * MAX);
    for (i = 0; i < MAX; i++)
        for (j = 0; j < MAX; j++)
            dim_a[i][j] = i * i;
    for (i = 0; i < MAX; i++)
        for (j = 0; j < MAX; j++)
            printf("%d", dim_a[i][j]);
    return 0;
}
```

Current and Future works

- Protocols designed for DLM cache-model (r3, r77 protocols were designed for DLM flat-model).
- Protocols designed for m-DSM (Distributed Shared Memory for multithread programs) (r58 is modified version of the current protocol used in m-DSM.)

Experimental Environment

| Server Name | Sandy Bridge Server cres12,3,4 | Haswell Server cres5,6,7,8 | Broadwell Server cres9,10 |
|-------------|--|--|--|
| Network | Infiniband single FDRx4 (56Gbps) | | |
| Node CPU | Xeon E5-2687W v3 (3.0GHz) 2 CPU x 8 Core | Xeon E5-2687W v4 (3.0GHz) 2 CPU x 10 Core | Xeon E5-2687W v4 (3.0GHz) 2 CPU x 10 Core |
| Node Memory | DDR4-1600, resist. 8GBx4 (64 GB) cres13: 16GB x 4 or 12 (128 /192 GB) cres2: 8GB x 16 (128 GB) | | |
| L1 cache | 32 KiB | 32 KiB | 32 KiB |
| L2 cache | 256 KiB | 256 KiB | 256 KiB |
| L3 cache | 25.4 GiB | 25.4 GiB | 30 MiB |
| OS | CentOS 7.1.1503 (x86_64, kernel 3.10.95) | CentOS 7.1.1503 (x86_64, kernel 3.10.95) | CentOS 7.2 (x86_64, kernel 3.10.95) |
| Compiler | gcc version 4.8.3-03 | | |
| MPI | MPICH2 2.0.1 | mvapich2.0.1 | mvapich2.0.1 |