

# 遠隔メモリを用いた大容量仮想メモリ DLM における メモリ管理機構の導入

斉藤 和広† 緑川 博子† 甲斐 宗徳†

†成蹊大学工学研究科情報処理専攻 〒180-8633 東京都武蔵野市吉祥寺北町 3-3-1

E-mail: †dm083406@cc.seikei.ac.jp, {midori, kai}@st.seikei.ac.jp

**あらまし** 筆者らはローカル物理メモリサイズに制限されず、クラスタの各ノードの遠隔メモリを集めて仮想的に大容量メモリとする分散大容量メモリシステム DLM を構築、評価してきた。従来の DLM では遠隔メモリ利用のための動的メモリ割り当て機能のみを提供し、メモリ解放と再利用機能は備えていなかった。大容量メモリを必要とするアプリケーションには、始めに必要なデータ領域を確保して最後までそれを用いるものもあるが、実行中に動的割り当てと解放を繰り返し行うものもある。このため、DLM においてもメモリ解放と再利用を行う機構の導入が求められていた。今回、メモリ解放機能と効率的なメモリ利用を実現するための最小限のメモリ管理機構を設計し、DLM へ導入した。本稿ではメモリ管理機構の概要について述べ、Red-Black Tree へのランダムなノード挿入削除を行うマイクロベンチマークと、実応用である遺伝子情報処理におけるクラスタリング処理におけるメモリ割付・解放関数の初期性能評価を行ったので報告する。

**キーワード** クラスタコンピューティング, 大容量メモリ, 遠隔メモリ, 分散メモリ, メモリ管理

## The Dynamic Memory Management for Distributed Large Memory System : DLM

Kazuhiro SAITO<sup>†</sup> Hiroko MIDORIKAWA<sup>†</sup> and Munenori KAI<sup>†</sup>

† Graduate School of Engineering, Seikei University,

3-3-1, Kichijoujikota-machi, Musashino-shi, Tokyo 180-8633, Japan

E-mail: †dm083406@cc.seikei.ac.jp, {midori, kai}@st.seikei.ac.jp

**Abstract** Authors designed and evaluated Distributed Large Memory System: DLM, which gives us very large virtual memory beyond the size of local memory by using remote memory distributed over cluster nodes. Originally, the DLM only supports a function for dynamic memory allocation, but it is required to be equipped with a memory release function for the applications that reuse memory resources in their runtime. We incorporate a simple memory management mechanism into the DLM to support memory release function for efficient memory utilization. In this report, we describe the design and mechanism of the memory management and show an initial performance evaluation for a red-black tree as a micro benchmark and a cluster3, a real bioinformatics application.

**Keyword** Cluster Computing, Large Memory, Remote Memory, Distributed Memory, Memory Management

### 1. はじめに

筆者らはローカルの物理メモリサイズに制限されることなく、クラスタの各ノードの遠隔メモリを利用し仮想的に大容量のメモリ空間を提供するシステム、DLM(Distributed Large Memory)を構築、評価してきた[1]-[4]。

従来の DLM には、大容量仮想メモリ利用の方法として遠隔メモリ動的割り当てと遠隔スワップの二つの機能を提供していた。物理計算や遺伝子分析等の大容量データを用いるアプリケーションでは、利用する領域を計算の始めに確保しこれを最後まで利用する形式

を取っているものがあるが、実行中に動的割り当てと解放を繰り返し行うものもあるため、DLM にもメモリ資源再利用のための遠隔メモリ解放機能とそのためのメモリ管理機構を試作導入した。

これまで DLM は、カーネルの変更を最小限にしたユーザレベルソフトウェアとして設計・評価してきた。遠隔メモリアクセス用デバイスを用いた OS カーネルスワップによる方式と比べ、ユーザレベルソフトウェアを用いた DLM がより高速かつ高安定性が得られることを示した[1]-[3]。今回新しいメモリ管理機構もユーザレベルソフトウェアとして構築した。

本稿では、第2節でDLMのシステム構成とユーザインタフェースを、第3節で新しく導入したDLMメモリ管理機構を述べ、第4節でその動作確認と初期性能評価を報告する。

## 2. DLM

現在DLMには2つのバージョンがあり、DLMのコンパイラを用いることで遠隔メモリサーバ起動をユーザから隠蔽したDLM-S(DLM-Single Client)[1]-[3]と、明示的に常駐のサーバを立ち上げることでマルチクライアント型にしたDLM-M(DLM-Multi Client)[3]を構築・評価してきた。本稿での実装・評価はDLM-Mで行ったので、このシステム構成とインタフェースを述べる。

### 2.1. DLM-Mのシステム構成

DLMはクラスタ環境における複数台のコンピュータを、あたかも一台のマシン上のメモリ資源であるかのように利用する。そのため、図1のようにクラスタノードを計算ノードとメモリサーバの二つに分類する。なお、図1はマルチクライアントバージョンであるDLM-Mの図である[3]。

計算ノードはユーザが実際にプログラムを実行するノードで、ユーザプログラムを処理する計算スレッドとメモリサーバとの通信を行う通信スレッドからなる。メモリサーバは計算ノードからのリクエストを受け、メモリを提供するノードである。DLMでのメモリ管理はDLMページ単位で行っている。DLMページサイズはOSのページサイズの倍数であればシステム構築時に自由に設定することができる。DLMではDLMページ管理表によりメモリを管理しており、このエントリを操作することで動的割り当てを可能としている。なお、計算ノードは全ノードの利用DLMページを管理しているが、メモリサーバは自身が保持するDLMページのみを管理している。また、ノード間の通信にはTCP/IPを用い、基本的にDLMページ単位でデータの送受信を行う。

計算ノードで遠隔メモリが必要になったとき、DLMページ単位での遠隔スワップによってメモリサーバが保持しているデータを利用する。

### 2.2. ユーザインタフェース

DLM-MではDLM-Sと異なり、ライブラリ関数のみを提供し、C言語で利用することができる。図2はDLM仮想メモリの利用例である。プログラムの流れは、dlm\_startup()によるDLMの初期設定から始まる。その後、dlm\_alloc()によりDLMの仮想メモリを利用したい変数に動的に割り当て、dlm\_free()を用いることで指定変数を解放することができる。dlm\_shutdown()はDLMのサービスを終了するとき用いる。実行時には、図

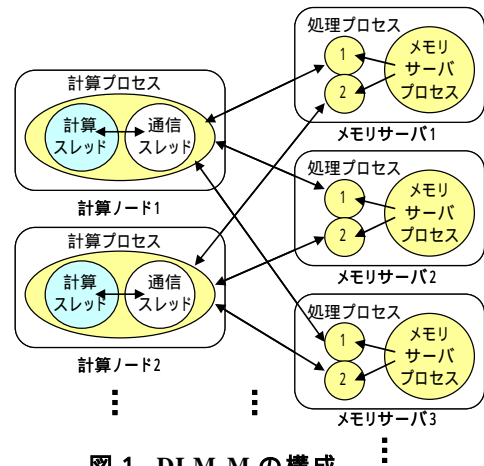


図1 DLM-Mの構成

```
#include <dlmm.h>
#define N 1<<30 /* 1G (*4B=4GB) */
int main(int argc, char* argv[]){
    int *a;
    dlm_startup(&argc, &argv); /* DLM-Mの初期化 */
    a=(int*)dlm_alloc(sizeof(int)*N); /* 動的割り当て */
    . . . /* 計算 */
    dlm_free(a); /* 解放 */
    . . .
    dlm_shutdown(); /* DLM-Mの終了 */
    return 0;
}
```

図2 DLM-Mの利用例

```
calhost 2048 // 2GB 計算ノード
memhost1 2048 // 2GB メモリサーバ1
memhost2 4096 // 4GB メモリサーバ2
:
```

図3 ホスト指定ファイル

3のようなDLMで利用する計算ノードとメモリサーバのホスト名と、それぞれで利用可能なメモリ量を記述するホスト指定ファイルをオプション-fで指定する。また-nで利用台数を指定することもできる。DLM-Mでは、実行前にメモリサーバとするノード上でdlmm\_serverを起動しておく必要がある[4]。

## 3. 遠隔メモリ管理機構

DLMの仮想メモリ解放機能のための管理機構と、これを利用したdlm\_alloc関数と、更にdlm\_free関数について述べる。DLMはユーザレベルソフトウェアでの実装を行っているため、カーネルのページフレーム管理機構とは別の最小限のデータ構造をユーザレベルソフトウェアで設計する。

### 3.1. 計算ノードのデータ構造

従来の DLM では大容量仮想メモリの動的割り当てのみを提供してきたので、DLM ページ単位でのメモリ管理だけで対応することができた。しかしメモリ解放機能を実装するには `d1m_alloc` による割り当て単位での管理が必要となる。また解放された資源のフラグメンテーションを回避し効率的な再利用するためには、空き領域の管理が必要となる。

`d1m_alloc` による割り当て単位での管理のためのデータ構造は図4のような単方向リスト構造になっていて、これを変数リストと呼ぶ。これは DLM 仮想メモリ空間上で先頭からアドレス順に連結される。変数リストにおける各ノード(変数ノード)には図5のように変数のメモリ情報等を持つ。また探索の高速化のため、DLM ページ管理表エントリからも各 DLM ページ内の先頭となる変数ノードへリンクしている。図4の DLM 仮想メモリ空間からの矢印、図5の DLM ページ管理表エントリからの矢印がその連結を示している。

空き領域管理には、図6のようにサイズ別に双方向連結リストを用いていて、これを隙間リストと呼ぶ。サイズ別隙間管理表の  $2^k$  のエントリには、 $2^{k-1} + 1 \sim 2^k B$  の範囲のサイズの隙間がサイズ順に連結される。管理される空き領域は以下の二種類がある。

- ・ DLM 仮想メモリ空間の先頭と変数の間の隙間
- ・ 変数と変数の隙間

隙間リストにおける各ノード(隙間ノード)にはその隙間情報等があり、更に変数削除時に隙間を生成するときの前後の隙間との結合を容易にするために、図7のように変数ノードに DLM 仮想メモリ上ですぐ前にある隙間ノードを示すポインタ `preg` が存在する。

### 3.2. メモリサーバのデータ構造

計算ノードとメモリサーバは DLM ページ単位でメモリのやり取りをするため、メモリサーバと計算ノード間でアドレスを管理する必要がない。メモリサーバは、メモリ割り当て要求を受け取ると必要な数のページを格納するためのスロットと呼ばれる領域を用意してページを保持する。メモリ解放要求を受けると、そのページで利用していたスロットは `munmap` せずに、再利用できるように図8のような空スロットリストに登録する。次に割り当て要求が来たときには、その空きスロットリストに利用可能なスロットがあれば、それを利用して DLM ページを格納する。この場合はメモリサーバは DLM ページ管理表を更新するだけで処理を完了する。

### 3.3. DLM メモリ動的割り当て

DLM ライブラリ関数の `d1m_alloc(size_t size)` を用いることで、DLM 仮想メモリへの動的割り当てをすることができる。従来はこの関数が呼ばれた時点での最終

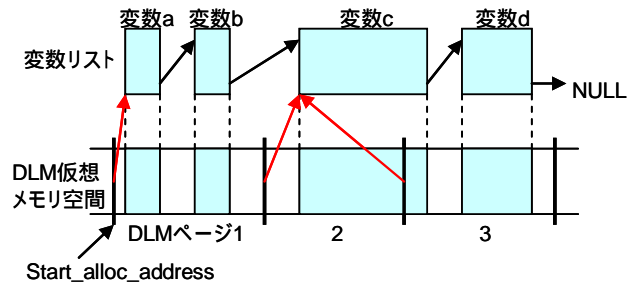


図4 変数リストとDLM 仮想メモリ空間の関係

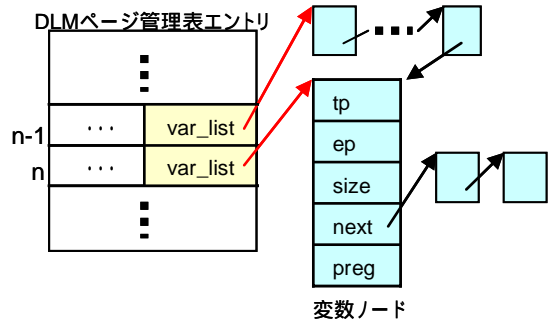


図5 DLM ページ管理表から変数リストへのリンク

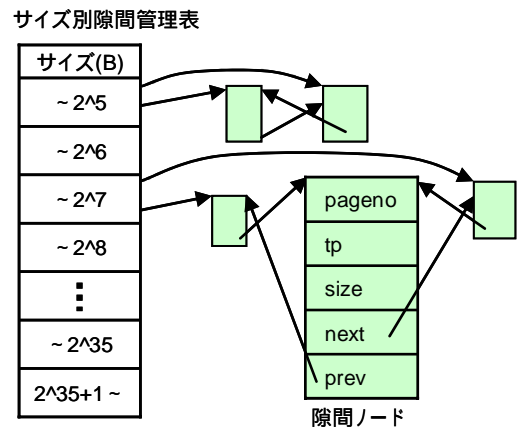


図6 計算ノードの隙間リスト

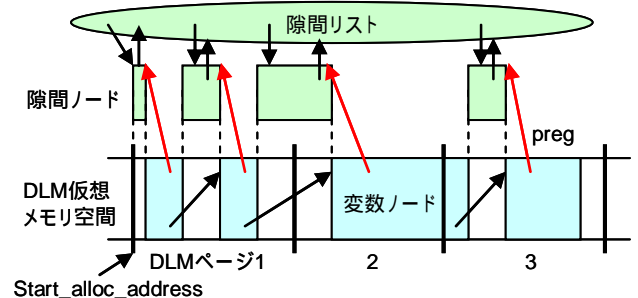


図7 変数ノードから隙間ノードへのリンク

割り当てアドレスから順に割り当てるだけであったが、変数管理と隙間管理が追加されたことにより、大幅に構造が変わった。

まず、隙間が一つも存在しないか、引数 `size` 以上の隙間が存在しない場合を考える。この場合は従来の

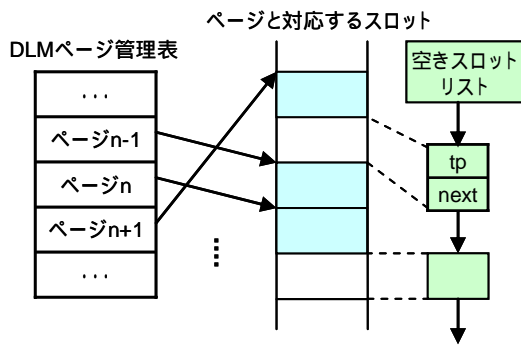


図8 メモリサーバのDLMページ管理とスロット管理

dlm\_alloc()と同様に必要なDLMページ数分を適切なノードに割り当てる。size分全てを割り当てたらDLMページ表を更新し、変数リストに登録する。

次に利用可能な隙間が存在する場合を考える。始めにsizeより大きい最小サイズの隙間ノードを探索し、発見した隙間ノードの先頭アドレスから順にDLMページ単位で割り当てていく。このとき以下のパターンが考えられる。

そのDLMページにおいて既にどこかのノードに割り当てられている場合、そのDLMページ管理表を更新するのみ。

そのDLMページがどこにも割り当てられていなくて計算ノードに空きがある場合、その分をmmapしDLMページ管理表を更新する。

そのDLMページが割り当てられてなく計算ノードに空きがない場合、空きのあるメモリサーバに割り当て要求を出し、すぐにDLMページ管理表を更新する。メモリサーバは要求を受け取ると、空きスロットが存在すれば空きスロットリストから、なければmmapでスロットを作り指定DLMページを保存し、DLMページ管理表を更新する。

これらをsize分完了したら、この変数情報を変数リストに挿入する。最後に隙間ノードを削除し、この隙間ノードを全て使い切らなかった場合は残った隙間を新しく隙間ノードとして登録する。そしてこのとき作成した変数ノードの次の変数ノードのpregをこの新しい隙間ノードへのポインタとして登録する。また使い切った場合にはこのpregをNULLにする。

### 3.4. DLMメモリ解放

DLMライブラリ関数のdlm\_free(void\* p)を用いることで、DLM仮想メモリに割り当てたデータを解放し、再利用可能にすることができる。

この関数が呼ばれると、まず引数pの変数情報を変数リストから取得する。次に、最初のDLMページか

ら順にDLMページ単位で解放していく。このとき以下のパターンが考えられる

そのDLMページにおいて他に変数が残っている場合、必要に応じてそのDLMページ管理表を更新するのみ。

そのDLMページに他に変数がなく、計算ノードが保持していた場合、その分をmunmapしDLMページ管理表を更新する。

そのDLMページに他に変数がなく、メモリサーバが保持していた場合、メモリサーバに解放要求を出しDLM管理表を更新する。メモリサーバはその要求を受け取ると指定DLMページの指定しているスロットを一つずつ空きスロットリストに繋ぎ、DLMページ管理表を更新する。

次に解放する変数の前後の変数ノードの情報から解放する変数の前後の隙間を取得し、その変数が開放されることで出来る隙間とそれらの隙間を結合する。結合後、後ろの変数のpregを更新する。またこのとき、DLM仮想メモリ空間上で最後尾の変数だった場合はその前にある隙間ノードを削除するだけである。そして最後に解放した変数ノードを削除する。

## 4. 初期性能評価

今回導入した遠隔メモリ管理機構の初期性能を評価した。評価にはマイクロベンチマークとしてSTAMP[4]のライブラリを用い、また実際のアプリケーションに近い例として、遺伝子分析で用いられるCluster3.0[5]を用いた。なお、評価は表1のマシンを用いて計算ノード1台とメモリサーバ1台で行い、DLMページサイズは128KBに設定した。

### 4.1. マイクロベンチマーク

STAMP[4]はマルチスレッドによるトランザクション処理のメモリの性能を調べるためのベンチマーク集である。今回は動作確認と初期評価のためにこの中のライブラリであるRed-Black Treeを実装したrbtree.cの逐次版を用いて、以下のような簡易のトランザクション処理を行った(main.c)。

50万個のランダム整数を挿入(insert)

1000個のランダム整数を削除

1000個のランダム整数を挿入

とを100回繰り返す(transaction)

残っているノードを全削除(release)

このプログラムの特徴は、変数と隙間の数が非常に多くできることと、隙間のサイズが一定サイズに多く集中することである。このとき、乱数は種を固定したときのlong型の0~1,000,000で、全てのノードとランダム整数はmallocで動的に確保される。ノードのサイズは48Bなので、動的メモリの最大使用量は約28MB

表 1 実行環境

ホスト	サーバマシン g1	クラスタマシン hp0
マシン	HP ML150 G3	HP ML150 G2 x 7 Nodes
CPU	Xeon E5310 1.6GHz QuadCore x 2CPU	Xeon 2.8GHz x 2CPU HyperThread
メモリ	8GB	1GB
L2キャッシュ	4MB/CPU	1MB/CPU
OS	Linux kernel2.6.23.17-88.fc7 x86_64	Linux kernel2.6.20-1.2320.fc5 x86_64
コンパイラ	gcc version 4.1.1 20070105	gcc version 4.1.1 20070105
NIC	NC7781 OnBoard Gigabit NIC	Broadcom 5721 PCI-Express Gigabit NIC
ネットワーク	1GbitEthernet	

表 2 main と main\_dlmm のメモリ操作別実行結果

	rbtree	local	alloc(s)	count	free(s)	count	Total(s)
insert	main		0.13235	500002	0.01979	106494	1.420
	main_dlmm	100%	5.77808		1.27921		8.490
		74%	6.53484		1.63273		494.349
		56%	7.38121	2.04381	1858.074		
transaction	main		0.01463	47419	0.00937	39894	0.186
	main_dlmm	100%	0.70120		0.69541		1.824
		74%	1.34586		1.58248		770.246
		56%	1.44040	1.75564	1483.255		
release	main		0	0	0.31215	401033	0.425
	main_dlmm	100%	0		90.82537		91.070
		74%	0		90.30602		91.088
		56%	0	99.96436	376.969		

となる。

次に ,malloc と free をそれぞれ dlm\_alloc と dlm\_free に代えたもの(main\_dlmm.c)を用意し,同様に上記の処理を行った。このとき,本応用でのメモリ使用量に対する計算ノードが保持するローカルメモリ量の割合(ローカル率)を 100%, 74%, 56% に変えて実行した。そしてこれらの実行時の (insert) の時間, (transaction) の時間, (release) の時間をそれぞれ計測した。

表 2 がその結果である。alloc 列が dlm\_alloc(malloc) 関数の実行時間, free 列が dlm\_free(free)関数の実行時間, count それぞれの実行回数, Total が Red-Black Tree 操作の時間も含めた各処理の合計時間となっている。なお insert 処理においての free は, 挿入のために一度ノードを動的に割り当てるが,すでに値が存在する場合はそれを挿入せずにすぐ free するよう rbtree.c で実装されているからである。

このプログラムは, 変数の数は全体で最大 401,033 個, DLM ページ 1 つの中では最大 2,730 個で, 各 DLM ページ内の変数リスト長が非常に長くなる。変数探索ではその変数がある DLM ページの先頭の変数からリストを辿るため, 一回の変数探索の平均メモリアクセス回数が 1365 回と非常に多くなる。また, 隙間数は release 時に最大で約十万個にまで達した。更に多くの隙間のサイズが同じため, 隙間管理表の 1 エントリに集中してしまう。すなわち, Red-Black Tree は alloc と free の繰り返し頻度が高く, 変数・隙間リストの探索回数も多いため, メモリ管理上非常に過酷な例である。このような状況においては malloc と free に比べ dlm\_alloc と dlm\_free 共に非常に低い性能結果となった。しかし alloc time や free time においては例えローカル率が低くとも, main\_dlmm.c の 74% や 56% の Total 実行時間のような急激な速度低下は見られなかった。

次に, main\_dlmm.c の実行時間に対する通信, メモリ・変数・隙間操作の時間の割合を dlm\_alloc と



図 9 main\_dlmm における dlm\_alloc の時間の割合

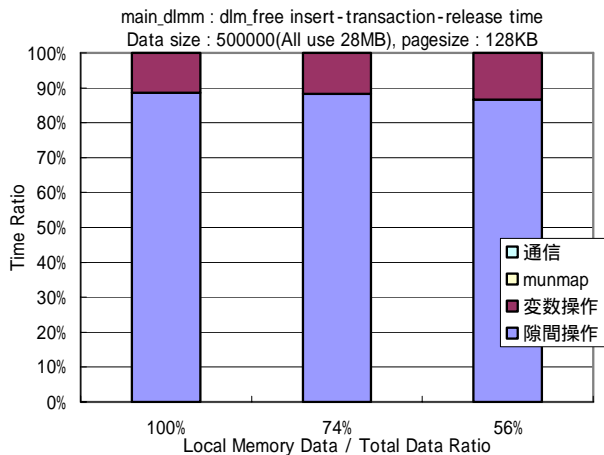


図 10 main\_dlmm における dlm\_free の時間の割合

dlm\_free それぞれでグラフ化したものが図 9 と図 10 である。図 9 から, dlm\_alloc におけるオーバーヘッドとなっているのが変数操作, 図 10 から dlm\_free では隙間操作であることがわかった。これは, 先に述べたような, 変数・隙間リストの探索が大きな原因となっていることの裏づけとなったといえる。

#### 4.2. Cluster 3.0

遺伝子分析のための様々なクラスタリング手法を提供するアプリケーションである Cluster3.0[5]の C 言語コマンドライン版を用いて性能評価を行った。本評



表 3 Cluster3.0 の実行時間

Cluster 3.0		Time(s)				
	local	alloc	realloc	calloc	dml free	Total
Original		0.017208	0.001193	0.000012	0.009185	9.637763
DLM version	100%	0.199455	0.112264	0.000022	0.148814	10.93558
	96%	0.204349	0.112475	0.000022	0.151385	13.99166
	77%	0.225566	0.112512	0.000022	0.16416	195.533
	58%	0.228844	0.112608	0.000022	0.16443	495.0805
	38%	0.237418	0.1123	0.000022	0.164922	906.9881
19%	0.244479	0.11241	0.000022	0.1665	1201.194	
count		24773	4937	1	24774	

価では、このアプリケーションのデフォルトの階層クラスタリングと、1.2MB の遺伝子情報入力ファイル demo.txt を用いた。なお、プログラム中に出てくるすべての動的割り当て・解放は dlm\_alloc と dlm\_free に変換した。

表 3 は実行時間と、各メモリ操作関数で要した時間とそれらの呼び出された回数、さらに全処理に要した時間を表にしたものである。このうち、Original が DLM 用に変換する前の元の Cluster3.0 のプログラムの実行時間で、DLM version は計算ノードのローカル率別に実行した。count はそれぞれのメモリ操作関数が呼ばれた回数である。Red-Black Tree での実行時と同様、ローカル率の割合が低い場合でも大きな性能低下がなかった。また全体の実行時間に対する比も 1% 未満と少なかった。これは、利用された変数の数が 12433 個と Red-Black Tree に比べ非常に少ない、これは、一つの変数のサイズが大きい (~ 39,456B) ので DLM ページ内部の変数ノード数が少ないことと、最大で隙間が 4 つしかできなかったことに起因する。

図 1 1、図 1 2 は図 9 と図 1 0 同様に操作毎の時間の割合をグラフにしたものである。Red-Black Tree に比べ変数・隙間操作ともに比率は少ないが、やはりこれらがオーバーヘッドとなっていることがわかった。

## 5. おわりに

本稿では、遠隔メモリを用いた大容量の DLM 仮想メモリ空間における動的メモリの管理機構を導入し、その初期評価を行った。現段階では Red-Black Tree のような非常に過酷に alloc と free を繰り返すプログラムにおいては未だ課題が残る結果となったが、一方 Cluster3.0 では全体の実行時間に対するメモリ操作時間の比が 1% 以下と非常に少なく、DLM のメモリ管理機構が Red-Black Tree 実行時ほどの大きなオーバーヘッドにはならないことが評価からわかった。今後の展望として、今回大きなオーバーヘッドとなった変数・隙間管理のリスト構造について、データ構造の再設計も考慮にいったチューニングを試みたいと考えている。

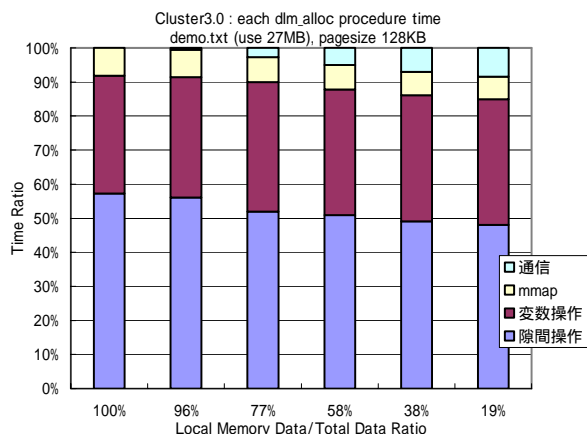


図 1 1 Cluster3.0 における dlm\_alloc の時間の割合

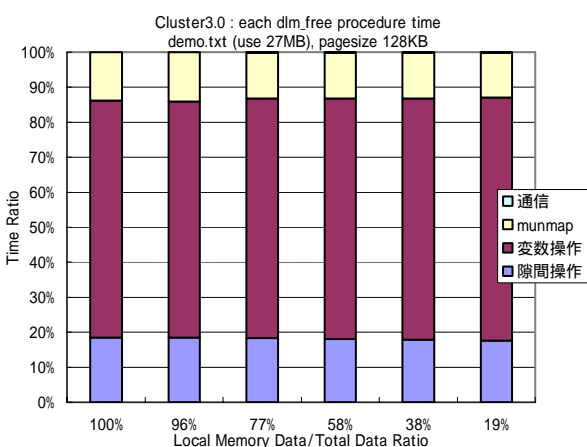


図 1 2 Cluster3.0 における dlm\_free の時間の割合

## 参考文献

- [1] DLM の設計と DLM コンパイラの構築”, 電子情報通信学会研究報告 CPSY 研究会報告, 信学技法 Vol.102, P.29-34, No.398, Dec.2007
- [2] 緑川, 黒川, 姫野, “遠隔メモリスワップのユーザーレベルソフトウェア DLM と性能評価”, 電子情報通信学会研究報告 CPSY 研究会報告, Aug.2008
- [3] 緑川, 黒川, 姫野, “遠隔メモリを利用する分散大容量メモリスシステム DLM の設計と 10Gb Ethernet における初期評価”, 情報処理学会論文誌: コンピューティングシステム, Vol.1 No.3(ACS1), 2008-12
- [4] 齊藤, 緑川, 甲斐, “マルチクライアント向け分散型大容量メモリスシステム DLM-M の設計と実装”, 情報科学技術フォーラム FIT2008, FIT 論文集, C-003, Sep.2008
- [5] STAMP (Stanford Transactional Applications for Multi-Processing)  
<http://stamp.stanford.edu/>
- [6] Cluster 3.0  
<http://bonsai.ims.u-tokyo.ac.jp/~mdehoon/software/cluster/>