

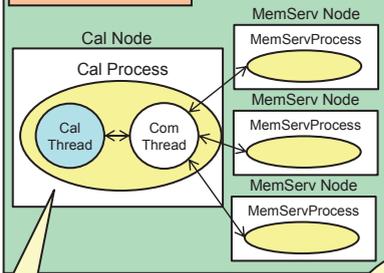
遠隔メモリ利用で大容量データ処理を可能にする 逐次プログラムのためのCコンパイラ

吉村 礎, 緑川 博子(成蹊大)

分散大容量メモリシステム (Distributed Large Memory System)

ネットワークで結ばれたコンピュータの物理メモリを通信によって利用し、逐次処理用に**仮想的に大容量のメモリ空間**を提供するシステム。

DLMシステム構成



DLMの遠隔Swap

計算ノード(Cal Node)で利用可能なメモリサイズを超えるメモリが必要な場合

→遠隔メモリ(MemServ)にデータを割り付ける
メモリサーバ上(MemServ Node)に割り付けたデータへアクセスする場合

→欲しいデータを持つメモリサーバ(MemServNode)からそのデータを含むDLMページを計算ノード(Cal Node)へ持ってきて(スワップイン)、不必要と思われるDLMページをメモリサーバ(MemServNode)へ送る(スワップアウト)。

DLM設定ファイル

MachineFile	MemoryFile
hp0	100
hp1	600
...	...

MachineFile
使用ホストを指定
MemoryFile
それぞれのサーバで
利用可能メモリ量を指定
(単位はMB)

プログラム実行コマンド例

```
mpirun -np 2 -machinefile machinefile01 / sample -- -f memfile01
```

プロセッサ2個使用

MachineFile指定

MemoryFile指定

Cal Processで
逐次プログラムが
走る

DLM コンパイラ

DLMシステムをユーザに使用しやすくするため開発。

定義

通常データ：ローカルメモリのみ使用する変数
DLMデータ：遠隔メモリにも展開可能な変数
DLMプログラム：DLMデータが宣言されているプログラム

DLMのAPIの例

静的データ宣言に適用
`dlim int a[2000];`
動的データ宣言に適用
`a = (int*)dlim_alloc(2000*sizeof(int));`

通常の変数宣言に「dlim」を
付加するだけ。
malloc関数をdlim_alloc関数に
書き換えるだけ。

プログラム例

通常のCプログラムに、DLMシステムを使用したい変数に「dlim」と付加するだけで、実行可能。DLMプログラムをDLMコンパイラにかかるとC言語プログラムに変換される。

通常のCプログラム

```
#define MAX 1000
int a[MAX][MAX];

int main (int argc, char *argv[])
{
    int i, j;
    for (i = 0; i < MAX; i++)
        for (j = 0; j < MAX; j++)
            a[i][j] = func();
    for (i = 0; i < MAX; i++)
        for (j = 0; j < MAX; j++)
            printf("%d", a[i][j]);
    return 0;
}

int func() {
    double b[MAX] = {1.0};
    int ans, i;
    for (i = 0; i < MAX; i++) ans += b[i];
    return ans;
}
```

DLMプログラム

```
#define MAX 1000
dlim int a[MAX][MAX];

int main (int argc, char *argv[])
{
    int i, j;
    for (i = 0; i < MAX; i++)
        for (j = 0; j < MAX; j++)
            a[i][j] = func();
    for (i = 0; i < MAX; i++)
        for (j = 0; j < MAX; j++)
            printf("%d", a[i][j]);
    return 0;
}

int func() {
    dlim double b[MAX] = {1.0};
    int ans, i;
    for (i = 0; i < MAX; i++) ans += b[i];
    return ans;
}
```

グローバル変数
指定可能

ローカル変数
指定可能

変換後のC言語プログラム

```
int (*__dlim_a_0)[1000];

int main (int argc, char *argv[])
{
    int i, j;
    dlim_startup(&argc, &argv);
    __dlim_a_0 = (int(*)[1000])dlim_alloc(1000*1000*sizeof(int));
    for (i = 0; i < 1000; i++)
        for (j = 0; j < 1000; j++)
            __dlim_a_0[i][j] = func();
    for (i = 0; i < 1000; i++)
        for (j = 0; j < 1000; j++)
            printf("%d", __dlim_a_0[i][j]);
    dlim_shutdown();
    return 0;
}

int func() {
    double (*__dlim_b_1);
    int ans, i;
    __dlim_b_1 = (double(*)[1000])dlim_alloc(1000*sizeof(double));
    for (i = 0; i < 1000; i++) ans += __dlim_b_1[i];
    dlim_free(__dlim_b_1);
    return ans;
}
```

DLMコンパイラを
使用しない場合は
このプログラムを
ユーザが書くことで
DLMシステムを
使用できる。

従来との特徴比較

変数	従来	DLM
静的変数	静的データ領域	ヒープ領域 (遠隔メモリにも動的に展開)
関数内部自動変数	スタック領域	
動的メモリ割り当て	ヒープ領域	

既存のプログラムで
ローカルメモリサイズに
影響をされず、
データの大規模化が
可能！

利点

稼働実験

HimenoBenchmarkとNPBBenchmarkで変換・実行を確認。

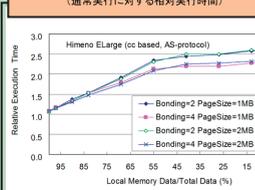
HimenoBenchmark

```
dlim float p[MIMAX][MJMAX][MKMAX];
dlim float a[MIMAX][MJMAX][MKMAX];
c[3][MIMAX][MJMAX][MKMAX];
dlim float bnd[MIMAX][MJMAX][MKMAX];
dlim float wrk1[MIMAX][MJMAX][MKMAX];
wrk2[MIMAX][MJMAX][MKMAX];
```

たった
4か所
変更する
だけで
使用可！

Himeno Benchmark におけるDLMの性能

(通常実行に対する相対実行時間)



95%遠隔
メモリを
使用しても
通常の
約2.5倍しか
遅くならない

実験環境：東京大学情報基盤センター, T2K (ha8000)
Network : 40Gbps, 20Gbps, Memory : 20GB/node x 2nodes

DLMコンパイラ構成

