

OS に中粒度の実時間性を付加する仕組みとその応用 または、ある日の T 研究室 # 4 もしくは、bit の没原稿 # 1

Presented by ^{Tada} 多田 ^{Yoshikatsu} 好克

コンピュータサイエンス誌 “ bad ” 2000 年 7 月号 Vol.32, No.7.

1 プロローグ

ある日の T 研究室¹。二日酔いの頭をだまし
ながら、T 教授がメールを読んでいる。

T : なんぢゃ、これは？

B : 何ですか、それは？

T : おお、B くんか。なんぢゃ、その綺麗な
恰好は ...

A : お早うございます。

T : ありゃ。A さんまで美しいお姿で ... 今
日はいったい何があるのぢゃ？

A : 卒業式ですよ、先生。修士の 2 年間、本
当にお世話になりました。

T : そうか。卒業式か。うんうん、君たちは
本当にお世話したぞ。

A : 何を言ってるんですか。

T : 社会人学生の D くんの恰好は、いつもと
あまり変わらないなあ。で、C くんは？

B : 卒業旅行と称して海外に行っています。

T : 何！彼は卒業までに 1 度は bit に出てみ
たいと言っていたのに、運の悪い奴ぢゃ。

A , B : えーっ、また bit の原稿なんですか？

T : うん。多田くんから変なプログラムが送
られて来てのぉ。

B : 面白そうですねえ。

A : あっ、私、卒業式に出てきま〜す。

T : うむ。今回は紙面も少ないし、君たちは
付き合う必要はない。ただし、夕方から
の呑み会は忘れないようにな。

A , B , D : はーい。

¹T 研究室については bit の 1998 年 10 月号「スレ
ドの作り方」の記事を参照のこと。

2 *et* と *MRSA* について

というわけで、T 教授は多田くんから送ら
れてきたメールの中身を吟味することにした。
メールには、tar して gzip して uuencode し
たプログラムと以下のような説明とが含まれ
ていた。

T 先生 :

多田です。

OS の外部で Unix システムとは独立に走
行するカーネルモードのスレッドを作りま
した。で、*et*(external thread) と命名しまし
た。可愛いでしょ。*et* を実現する仕組みは
MRSA(moderate real-time system adapter)
と言います。耐性黄色葡萄球菌とは何の関係
もありません。フリビ²の 2.2.5 で動きます。
遊んでみてください。なお、詳細は文献 [1] を
見てくださいね。

- - Y o s h i T a d a

うーむ、相変わらずだのぉ。とつぶやきな
がら、T 教授はメールで送られてきたプログ
ラムを展開し始めた。

Readme によると、*MRSA* は FreeBSD
2.2.5 に標準の lkm (loadable kernel module)
を利用して実現しているらしい。Unix の全
システムと対等な立場で *et* を走らせ、メタス

²FreeBSD のこと。ちなみに、ポケビはポケットビス
ケッツ、トラギはトランジスタ技術。

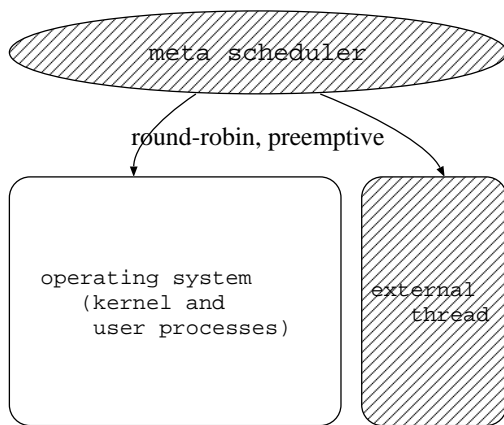


図 1: MRSA と *et* の概念図.

スケジューラによって Unix と *et* とを交互に実行する (図 1 参照)。 *et* は Unix の時分割スケジューリングの対象外となり、20m 秒³程度の中粒度の周期性を持って CPU を占有した処理を遂行できる。また、メタスケジューラは CPU 横取りをするので、*et* を使ってサーバのような無限ループの処理を記述することも可能である。

図の斜線部分が MRSA だな。その一部分で *et* が動いているという ... うーむ、分かったような分からないような。

3 *et* の利用者インタフェース

et の利用者は、*et* のプログラムを書くことになる。*et* のプログラムは Unix のカーネルモードで動作する。そのため、通常の C ライブラリ (libc.a) を呼び出すことはできない。一方、カーネルのデータにアクセスしたり、カーネルプログラムの関数を呼び出したりすることはできる。このことを除けば、*et* のプログラムは通常の C のプログラムと変わらない。

ただし、プログラムの実行は main() ではなく et_main() という関数から始まる。この関数が終了すれば *et* の実行も終了する。また、et_exit() を呼び出して、*et* の実行を陽に終えることもできる。

Unix と *et* とを交互に動かす間隔は、記号

³フリビのカーネルの記号定数を変更すれば、数 m 秒程度までの、より細かい周期性も実現できる。

定数 UNIX_SLICE と ET_SLICE を定義して指定する。この単位はチック (tick) と呼ばれる時間間隔であり、現在のフリビでは 10m 秒である。ただし、Unix のカーネルを書き換えれば 1 m 秒程度まで短くすることが可能である。

et の実行が ET_SLICE チックを使い切ると、CPU 横取りが起こり、Unix の実行が再開される。そして UNIX_SLICE チック後には、再び *et* の実行が再開する。*et* が et_yield() を呼び出すと *et* の実行は ET_SLICE を使い切る前に一時中断する。この場合も UNIX_SLICE チックが経過すると、再び *et* の実行が再開する。

et_lock() は *et* に対する CPU 横取りを禁止する。ただし、*et* が制御を確保し続けると Unix カーネルに不都合が生じる。et_lock() を呼び出した場合は、短期間の内に et_unlock() を呼び出し、CPU 横取りを可能にしなければならない。

上述の関数の他に MRSA では Unix との同期や入出力のための関数を用意している。まず、p_sleep() と p_wakeup() は Unix カーネルと同期を取るための関数である。これらの関数には、Unix の sleep(), wakeup() と同様に、同期を取る事象を識別するための適切な変数の番地を引数として与える。

p_putchar(), p_printf() は *et* から Unix コンソールへの出力関数である。また、p_write(), p_writec(), p_read(), p_readc() は Unix のプロセスと情報のやり取りをする。Unix プロセス側ではデバイス/dev/etio をオープンし、read, write することによりデータの授受ができる。ちなみに、p_wrtcec(), p_readc() は 1 文字の受渡しをする。なお、これらの関数はランデブーを基本とし、相手がデータの読み書きをするまでブロックする。

4 *et* の応用プログラム

ふーっ。パイプをくゆらしながら、T 教授は長いため息をついた。面白そうだが、いったい何に使うのだろう。usr という名のディレ

クトリにいくつかの応用プログラムがあったので、それを make することにした。

make してみると、proc_mod.o, map_mod.o, io_mod.o というオブジェクトが作られた。これらは、それぞれ、Unix のプロセスの状態を可視化する *et* プログラム、あるプロセスの仮想アドレスのマッピング情報を可視化する *et* プログラム、そして、プロセスから直接入出力装置を制御するための *et* プログラムである。

なるほど、*et* は Unix とは独立に作動するので、Unix の各種状態をプロンプ効果なしに観測できるわけだな。そうすると、*et* を使って Unix カーネルの各種パラメータを動的にチューニングすることもできるのか。io_mod.o は入出力装置を直接制御するという魂胆だな。デバイスドライバを書かなくても、*et* ではある程度の実時間性を実現できるので、好都合というわけだ。

ふーっと、もう一度パイプの煙を漂わせ、T 教授は proc_mod.o を実行させることにした。

proc_mod.o は lkm のモジュールぢゃな。ということは、modload コマンドを使ってこれをカーネル内に押し込んでと ...

```
% modload -u -o /tmp/et_mod
-e et_init proc_mod.o
```

これで、*et* が動き始めたわけだな。次に *et* からの情報を取り出すプロセス proc を動かせば良い ... proc は /dev/etio をオープンして read する通常の Unix プログラムだから、そのまま

```
% proc
```

と打てば良いのぢゃな。

おお、こ、これは ..(と、画面には突然図2のウィンドウが現れた。)

画面には上のようなウィンドウが現れ、個々の小さな四角が緑や黄、赤に点滅している。色は Unix プロセスの状態を表し、緑が実行中、黄が実行可能、赤がサスペンド状態とのこと。また、図で黒く見えるのは実は青色で、事象を待って休眠中のプロセスらしい。

うーむ、カラーで見えない bit の読者が可哀相ぢゃ。美味しそうにパイプを吸いながら、

T 教授はつぶやいた。*et* が Unix カーネル内の情報を調べ、それを書き出す。Unix プロセスはデバイスドライバから情報を読み出し、Xウィンドウに描画するというわけぢゃな。どれ、少しプログラムを見ておくれ。

```
struct dat {
    int    d_pid;
    char   d_stat;
} d_buf[MAXBUF];
```

というデータ構造に個々のプロセスの状態を集めるのが、以下の *et* プログラムぢゃな。

```
void et_main() {
    struct proc    *p;
    int            i;
    while (1) {
        i = 0;
        for (p=allproc.lh_first; p!=0;
             p=p->p_list.le_next) {
            d_buf[i].d_pid = p->p_pid;
            d_buf[i].d_stat = p->p_stat;
            if (i++ >= MAXBUF)
                break;
        }
        p_write(d_buf,
                sizeof(struct dat) * i);
    }
}
```

単にプロセスのリストを嘗めて d_buf に情報を集め、それをデバイスドライバに書き出しているだけだな。プログラムは無限ループだが、p_write() がブロックするし、CPU 横取りも起きるから問題無いということか ...

で、Unix プロセスの方は、我々の良く知っている普通のプログラムだな。

```
#define BUFSIZE sizeof(struct dat)

main() {
    int fd;
    int i, j;
    if ((fd=open("/dev/etio",
                 O_RDONLY)) == -1)
```

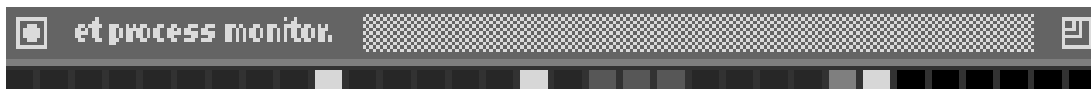


図 2: Unix プロセスの実行状態の可視化

```
exit(1);
while (1) {
    i = read(fd, d_buf,
            BUFSIZE * MAXBUF);
    i /= BUFSIZE;
    /* d_buf の値を使った描画 */
};
close(fd); /* 実行されない */
}
```

これまた単純だなあ。多田くんの書くプログラムはいつもこうじゃ。要するに/dev/etioからd_bufに情報を読み出す。で、その情報を使って描画してるだけだ。このプログラムも無限ループだが、readがブロックするのでetと同期が取れるのだな。

ps コマンドを使うと、ps コマンドを実行しているプロセスが常に実行中になるが、etを使って情報の収集をやるとUnixの自然な姿が観測できるのぢゃな。それに、描画はUnixのプロセスでやるから、プログラムも簡単と。こりゃ、なかなか便利ぢゃな。

仮想アドレスのマッピング情報を見るプログラムも似ているな。etではカーネル関数が利用できるの、マップの情報を得るのも簡単なこと。100m秒の間隔で情報を取得しているが、もし、描画処理が間に合わなくても入出力がランデブーなのでOKということか...

入出力装置の制御プログラムは便利そうだなあ。Windowsにはソフトウェアでタイミングを取る入出力装置も多いが、etを使えばUnixでも似たようなことができそうぢゃ。etを使ってロボットの動きを制御し、かつ、動作のプランニングにはUnixプロセスを使うといった応用にも使えるかもしれない。うーむ。もう少し付き合ってみるか。

5 MRSA の仕組み

T教授は、etを実現しているMRSAのプログラムを調べることにした。Readmeにはフリビのlkmを利用していると書いてあったが、それならばモジュールの取込み時に初期化をやっているはずだな。どれどれ...

```
et_load(struct lkm_table *lkmtpl,
        int cmd) {
    extern int etio_begin();
    int err;
    err = etio_begin();
    stack = (unsigned int)
        kmem_alloc(kmem_map, s_size);
    sp = (int *)((stack+s_size) & ~3);
    *--sp = (int)et_exit;
    et_buf[R_ESP] = (int)--sp;
    et_buf[R_EIP] = (int)et_main;
    timeout(to_et, 0, UNIX_SLICE);
    return(err);
}
```

ふむ。思った通りぢゃ。etio_begin()によってデバイスドライバの準備をして、後はetのコンテキスト作りぢゃな。カーネル関数のkmem_alloc()を使ってetのスタックを確保し、et_bufにetの初期コンテキストを格納すると...多田くんのことだから、おそらく、et_bufとunix_bufという2つのコンテキスト格納領域を使って、etとUnixを交互に動かすつもりだろう。してみると、この2つのバッファはjmp_buf型で、コンテキスト切替えにはsetjmp()とlongjmp()を使うのだな[2]。多田くんも進歩がないなあ⁴。しかし、最後のtimeout()は何だ？これはカーネル関数で、確かマニュアルがあったはずだぞ。

と、T教授はフリビのマニュアルのtimeout(9)を読み始めた。ちなみに、Unixマニユ

⁴放っといてくれ。

アルのセクション9と言うのは耳慣れないかもしれないが、intro(9)によるとカーネル内の関数のインタフェースなどがまとめられているらしい。

何々、第3引数で指定したチックが経過した後、第1引数の関数を呼び出してくれるのぢゃな。第2引数は関数への引数だが、今は使わないから関係ないと。ということは、UNIX_SLICE チック後に `to_et()` という関数を呼び出せと言っているのだな。どうせ、`to_et()` は `unix_buf` に Unix のコンテキストを保存して、`et_buf` から `et` のコンテキストを回復するのぢゃらう。一応、確認しておくか ...

```
void to_et() {
    if (setjmp(unix_buf)) {
        /* just return */
    } else {
        timeout(to_unix, 0, ET_SLICE);
        longjmp(et_buf);
    }
}
```

ふっふっふっ、凶星だな。なるほど、今度はここで `ET_SLICE` 後に `to_unix()` を呼び出すように設定する。`to_unix()` は `to_et()` と同じようなものだから見るまでもなかるう。

どれ、呑みに行くかな。Aさん達はまだかな？

6 RT-Linux との比較

Aさん、Bくん、Dくんが卒業式から帰ってくるのを待つ間に、T教授の元にもう1通のメールが届いた。

T先生：

多田です。

MRSA と良く似たシステムに Real-Time Linux と言うのがあります。有名なシステムなので、Web で捜せばすぐに見つかると思います。文献 [3][4] もあります。bit に記事を書くなら、比較しておいてください。

- - Yoshi Tada

比較しておいてくださいだお！どこにそんな暇があると言うのぢゃ。と思ったが、取り敢えず卒業式が終わるまでは暇か ... どれ、少し調べてみよう。

何々、RT-Linux も *MRSA* と同じように Unix のカーネル外でカーネルモードのスレッドを走らせるのか。違いはどこにあるのぢゃ？ふむふむ、RT-Linux では制御の獲得をタイム割込みで実現しているのか。*MRSA* よりも粒度の細かな制御ができるということか。しかし、プログラミングは大変そうぢゃお。Unix カーネルとの排他制御も考えなければならぬぞ。

だいたい、以下のような比較で良いかな(表1参照)。どれ、では、この表を多田くんにメールしておこう。

さて、これで今日の仕事は終わりぢゃ。無事 bit の原稿も書けたし、そろそろAさん達も帰って来る頃だし、ぼちぼち呑みに行くことにするかな。

7 エピローグ

T教授が呑みに行ってしまったので、*et* と *MRSA* に関する話はこれで終わり。最後に、話を少しまとめて、この記事も終わることにする。

まず、*MRSA* を利用することにより、稼働中のカーネルに新しい機能を追加できる。また、あらかじめカーネルをモジュール化して作っておくことにより、必要な機能を追加したり不要な部分を取り外したりして、用途に合ったカーネルを動的に構築することも可能となる。

このようなカーネルの発展性は、マイクロカーネルとサーバ群による構成でも実現されてきた。しかし、サーバ間の通信にコンテキスト切替えが必要なことから、十分な処理速度を確保しにくい [5]。他方 *lkm* を使った方法では、単層カーネル (monolithic kernel) に機能を追加するので速度の問題はない。しかし、

表 1: *MRSA* と Real-Time Linux との比較 .

比較項目	<i>MRSA</i>	RT-Linux
制御の獲得法	timeout ルーチン	タイマ割込み
最短起動周期	20 m秒*	数 10 μ 秒
外部事象への対応	20 m秒* + 数 100 μ 秒	< 20 μ 秒
CPU 横取り	あり	なし
カーネル関数	利用可	利用不可
カーネルとの排他制御	不要	必要
デバイスドライバ	ランデブー	ブロック不可
プログラミング	比較的容易	難しい
割込み処理	オーバヘッドなし	かなりのオーバヘッド
システムの導入	簡単	ドライバの書換えが必要

*) Unix カーネルの簡単な変更により、数 m 秒程度まで改善可能。

モジュールのプログラムに虫があると、OS 全体が脱線転覆するという問題がある [6]。

また、lkm ではシステムコールやファイル操作、デバイスドライバなど、制御の取得点が決っており、自発的に処理を開始するプログラムを記述しにくかった。*MRSA* はこの欠点を補っていると考えることができる。既に見てきたように、*et* のプログラミングにおいては無限ループを書くことも可能である。

et のプログラムは、短い時間間隔で周期的に実行される。いわゆる実時間処理ではないが、動画や音声のマルチメディア処理、ロボットの制御などに利用できるのではないかと考えている。中粒度の実時間処理を実行できる一方で、我々が馴れ親しんだ Unix の機能をそのまま使えることのメリットは大きい。プログラムも比較的簡単である。

今後は、*MRSA* と *et* を利用したシステムの記述に挑戦したいと考えている。

参考文献

- [1] 多田好克, 中村嘉志: *MRSA*: Unix 上の中粒度実時間処理系, 信学技報, CPSY99-123, pp. 1-8 (2000年3月) .
- [2] 多田好克, 寺田実: 移植性・拡張性に優れ

たCのコールチンライブラリ実現法, 信学会論文誌, J73-D-1, 12, pp. 961-970 (1990年12月) .

- [3] Barabanov, M., and V. Yodaiken: Introducing Real-Time Linux, Linux Journal, pp. 19-23 (Feb., 1997).
- [4] 船木陸議: RT-Linux の構造と実装の詳細, インターフェース, 25, 11, pp. 88-100, CQ 出版社 (1999年11月) .
- [5] Liedtke, Jochen: Toward Real Microkernels, CACM, 39, 9, pp. 70-77 (Sep., 1996). (谷口秀夫訳: 真のマイクロカーネルへ向けて, bit, 29, 8, pp. 11-21, 共立出版 (1997年9月) .)
- [6] 鈴鹿倫之, 中村嘉志, 多田好克: カーネル拡張のための効率的な開発環境, 情報処理学会第41回プログラミング・シンポジウム報告集, pp. 57-64 (2000年1月) .