

東海大学の清水です。

99年5月号の原稿になります。よろしくお願ひします。
COMPAQから新しい発表があったことを受けて当初予定から
内容を変更しました。

また、できれば

<http://www.jp.debian.org/debian/dists/slink/main/{disks-alpha,binary-alpha,source}>

<http://www.jp.debian.org/debian-jp/dists/slink-jp/main/{binary-alpha,source}>

の内容をCDROMに載せて頂けますか？

性能解析ツール lprobe の紹介

清水尚彦 nshimizu@et.u-tokai.ac.jp

最近、学科の連絡なども電子メールを中心に連絡を取り合うようになってきて
ようやく世間並みにペーパーレスになりつつあるのですが、困ったことに
MSWordやら一太郎やら果てはMacintosh binaryの添付ファイルが本文になっている
メールが来ることも多くなり文句を言っても他に方法を知らないといわれたりします。
こういった場合に平文テキストを強く主張してもなかなか受け入れてもらえず
困ってしまいます。普段PC+Unixで生活していて自分の作る書類にはあまり困って
いないだけにWinodwsを使わないほうが悪いといわんばかりの風潮には反発を覚えて
しまいますね。とはいっても何とかしなくては仕事に支障が出てしまいます。
対策としては

1. Linux+WineでMSWordを動かす。(もしくはWordが読めるLinux用のワープロを導入する)
2. 日常の環境をMSWindowsに変更する
3. 二台のPCを切り替えて使う

などがありますが、どれもスマートではありません。1はだいぶましですが、仕事に
使うほどの安定度が確保できるか評価が難しいところではあります。

こんな時に大型コンピュータの世界ではVMを導入して同時に両方のOSを動作させる
のですが、i386アーキテクチャは自己仮想化をサポートしないのでVMの導入は困難です。
Alphaでも、ハードウェアリソースに書き込みのみのレジスタがあったりするのでそのまま
仮想化するのにはオーバヘッドが大きそうですが、やってみる価値は十分にありそうですね。

私はすぐにでも必要なので、もう少し現実的な方法としてInterixをx86のWinodwsNTの
上にのせてそこに普段使うアプリケーションを乗せてしまうという方法を検討中です。

本当はWinFrameを使ってリモートのNTのアプリケーションを使った方がスマ
ートですが、WinFrameはとんでもない値段が付いているので個人では手が出ま
せん。NTをメインとするとしても

普段使っているTeXやtgifやjvimなどを捨てたくはない

ですよ。TeXやjvimはNTでも動作するものがあるので、残りはtgifです。これを
ソースからNTで動かすように変更するのは大工事になるのですが、Interixなら
かなり簡単に移植できるのではないかと密かに期待しています。

tgifははまだですが少しだけNTを使いはじめました。しかし思った以上に安定度が悪く
一日に何回かはリセットボタンを押すはめになっています。本当にこんな

ひどいものをみんな使っているのだろうか？とも不思議ですが、私の

システムだけが不安定なのだろうか？もっとも固まるのはNetscapeを使っているときが
ほとんどのような気がするのでマイクロソフト社の陰謀かもしれません :-)

InterixはAlphaでも動作しますが、NTなのであくまでもアドレス空間は

32ビットでしか動きません。(もちろん、プログラムは64ビットレジスタを使えます。

ここでの32ビットアドレスとは上位32ビットが0のアドレスのことです。)そこで、
数GBのメモリが必要な人がAlphaのマシンにNTを入れてもあまり嬉しいことはありません。

Linuxの2.2.Xのカーネルが出てきました。

Linux/Alphaのインストールという目でみると2.0.Xに比べカーネル作成時にGenericという
選択肢が増えたのが大きい変化です。これはマザーボードを

カーネルが自動認識するものでマザーボードに強く依存した従

来のカーネルに比べてインストール時のトラブルが減る可能性があります。

(もっともMILOはまだボード別になっているのでまだ意識は必要です。SRMからMILO
を使わずにインストールすればGenericだけですむかもしれません。)

Genericは2.1の系列にも入っていたようでRedHatのサイトにあるカーネル(2.1.124)には
Genericのものも置いてあります。

OpenBSDやNetBSDでは(Digital Unixでも)自動認識していたのでやればできる
ことだったはずですが、これまでなされていなかったのが不思議ですね。

年明け早々にCOMPAQからAlphaコミュニティに大きなプレゼントが来ました。それはlprobeと呼ばれるパッケージです。このパッケージはAlphaの性能を測定するための特殊なカウンタを使うためのフロントエンドになっていて簡単な設定で様々な性能上の特性をとらえる事が出来るようになります。

<ftp://metalab.unc.edu:/pub/Linux/devel>

から入手できます。
ファイルは幾つかあり、

lprobe_suite-4.0-1.alpha.rpm
lprobe_suite-4.0-1.src.rpm
lprobe_suite-4.0-driver.o.gz
lprobe_suite-4.0.lsm
lprobe_suite-4.0.tgz

のようになっています。
今回はこのlprobeの使い方と何が出来るのかについてまとめてみたいと思います。

1. Multia(UDB)のトラブルシューティングとDebian 2.1のインストール

後述のlprobeと

2.2.Xのカーネルを試すためにトラブル中のMultia(UDB)を復活させようとWEBサイトを探していたところ(時期的に現役のマシンはAS200 4/100を含めて学生が常時使っていて私が自由に使えるAlphaのマシンは少ないのです)

<http://brouhaha.com/~eric/computers/udb.html>

にMultiaの情報だけでなくマザーボードのマニュアルも置いてありました。その中のTrouble Shootingの項目から私のトラブルに該当するものがありました。現象としてはブートのモードを切り替える時にそのまま固まってしまうというもので、私はBIOSのEEPROMイメージが壊れたものだと思っていたのですが、マニュアルによるとバックアップ電池が切れた場合の現象とのことでした。^^;

(この程度の問題ならユーザーズマニュアルに載せておいて欲しいなあ)

さっそくパソコンショップに出かけて「これと同じ型の電池はありますか?」と尋ねたところMacintosh用のバックアップ電池と同型だそうで案外簡単に手に入りました。

ということで、無事Multia君も生き返り新しいカーネルを試せると

思ったのですが、思わぬ所に伏兵がいました。

実は長いこと使わなかったうちにMultia君のフロッピードライブの調子がすっかりおかしくなってしまったのです。

フロッピーからデータを読み出そうとするとウンウンうなってI/Oエラーとなってしまう ;_;

それでも何回かやっているうちにMIL0だけは起動できたのでHDDにあった昔のパーティションに新しいMIL0を入れることができました。

(昔のパーティションからブートすればよいではないかと気が付いた人は鋭いですね。実はMIL0のファイル名にちょっと細工していたのですっかりファイル名を忘れてしまってブートまでこぎ着けなかったのです。MIL0が立ち上がってファイル名が確認できたのでとりあえず古いバージョンのLinux/Alphaは動くようになりました。)

RedHatのgenericカーネルも何回かに一回は立ち上がりペンギンの絵が画面に出て来ました。ところが、ramdiskのフロッピーはどう頑張っても読めないのです。

ということで、今回は残念ながら時間切れです。次回にSRMコンソールからディスクレスブートをするようにしてMultia君の再生に再挑戦します。

というはずだったのですが、フロッピーが読めないのも不便なのでいろいろ試してみました。古いパーティションのカーネルを使って一旦Linuxを立ちあげた後ddコマンドをシェルのループ機能で連続発行してフロッピードライブを動作させ続けているとしばらくするとドライブの音が静かになってきたと思ったらちゃんと読めるセクタが増えているではありませんか!

喜んでMIL0からカーネルを読み込んでramdiskのイメージを読み込ませました。調べてみるとどうにかこの弱ったドライブで読めるのはソニーのディスクだけでマクセルや富士フィルム、DYSANなどのメーカー製はぜんぜん読めない

ようです。(Multia以外のハードでは問題なく読み書きできるのでこれらのメーカーが悪いわけではありません。)

フロッピーが読めたのでRedHat 5.2のインストールに挑戦です。
RedHatのミラーサイトからダウンロードしたramdiskを使ってMIL0から

```
boot sda1:generic.gz root=/dev/fd0 rw load_ramdisk=1 prompt_ramdisk=1
```

としてインストールを開始します。ここでsda1としているのは前述のとおり私は一旦古いシステムを立ちあげたのでその時にgeneric.gzをMIL0のパーティションに入れていたからこう指定しています。カーネルをフロッピーに入れている人は「fd0:generic.gz」に読み替えてください。
カーネルが読み出されるとramdiskのフロッピーを要求してきますので入れ替えてキーを何か押すとramdiskの読み出しが始まり、インストール用のシステムが立ち上がります。私のMultiaにはCDROMがついていないので、ftpを使ってインストールを試みました。ところが、IPアドレスを直接指定してもDHCPを指定してもどちらもエラーになってネットワークを有効にできません。Alt-F3でインストールの途中を確認するとifconfigはちゃんとできていそうなのになぜかスクリプトがエラーを検出します。今回は時間がないのでRedHat 5.2のトライアルはこれでおしまいになりました。落ち着いたらゆっくり原因を探ってみたいと思います。

次にDebianプロジェクトから出ているslinkバージョン(Debian 2.1)のインストールを試してみました。執筆時点ではまだ安定版にはなっていませんがフリーズ版になっていてリリースへの秒読み段階です。
この雑誌が販売されるころには大きく事情が変わっているかもしれませんが、御了承ください。

CDROMがあれば簡単ですが、slinkのアルファ版のCDROMなんてそうそう簡単には手に入らないのでFTPによるインストールを行います。
インストーラがNFSで配布ファイルを含むディレクトリをマウントできるようにもう一台のマシン(x86でOK)を設定しておいてください。
Alphaを持っている方はそれ一台しかない場合は少ないと思われるので一台だけの場合について言及しませんが、もし単独でインストールしたい場合、もしくはNFSのサーバーの設定が難しい場合にはフロッピーベースのインストールを行ってください。

Debianではパッケージの管理を行うツールが使われており最初にインストールするのはベースパッケージと呼ばれるものとカーネル・モジュールだけです。Debianのミラーサイト(www.jp.debian.orgなど)の

```
debian/dists/slink/main/disks-alpha/1998-11-15
```

から自分のマシンに合わせてMIL0とカーネルの入ったファイルとroot1440.binとbase2_1.tgzを取ってきましょう。
*.binというファイルは1.4MBの大きさですが、base2_1.tgzは8MB以上あります。
jensenだけは執筆時点でupdateが入っていますが、他のマシンはすべてこのディレクトリの下から取ってきます。
新たにインストールをする場合には必ずupdateのディレクトリを確認して自分のマシン用のupdateファイルが出ているか確認してください。

base*.binというファイルがたくさんありますが、こちらはフロッピーベースのインストールで使うものですから今回は使いません。
ここには次のようなファイル/ディレクトリがならんでいます。

.	book1	eb164	root1440.bin
..	cabriolet	eb64p	ruffian
alcor	config.book1	eb66	sable
avanti	config.generic	eb66p	sable-g
base14-1.bin	config.jensen	jensen	sx164
base14-2.bin	config.miata	lx164	takara
base14-3.bin	config.miata-s	miata	xl
base14-4.bin	config.mikasa	mikasa	xlt
base14-5.bin	config.noritake	noname	
base14-6.bin	config.sable	noritake	
base2_1.tgz	config.sable-g	pc164	

各マシン毎のディレクトリには大体次の二つのファイルが入っています。

```
drv1440.bin resc1440.bin
```

drv1440.binはARC コンソールかAlphaBIOS用のMIL0をロードするためのもので

```
linload.exe  
milo
```

のファイルが入っています。
一方、resc1440.binにはカーネルが入っています。
ブートの方法には大きく分けて三つあります。

1. jensen, noritake, sable, mikasa, book1の各マシンではカーネルとともにaboutというブートマネージャが入っていてMIL0を使わずにSRMからブートさせます。ブートのコマンドは次のようにします。

```
>>>boot dva0 -file linux -flags "root=/dev/fd0 rw ramdisk_size=5120  
load_ramdisk=1 prompt_ramdisk=1"
```

2. ARCもしくはAlphaBIOSのマシンでは前述のようにlinload.exeからMIL0を起動させますので、BIOSの設定画面でフロッピーディスクからlinload.exeを起動してOSファイルとしてmiloを読み出すようにします。

3. それ以外のSRMを持つマシンはresc1440.binのフロッピーにMIL0とカーネルが両方とも入っていてSRMからフロッピーをブートするとMIL0が起動します。

```
>>>boot dva0
```

MIL0を起動した場合にはMIL0からLinuxをさらに起動する必要があります。
resc1440.binのフロッピーディスクをドライブにセットして次のようにブートさせます。

```
MIL0>boot fd0:linux root=/dev/fd0 rw ramdisk_size=5120 load_ramdisk=1  
prompt_ramdisk=1
```

ramdiskのフロッピーを入れるように指示して来ますのでroot1440.binをセットしてキーを押します。

ARCやAlphaBIOSのマシンでハードディスクからブートしたい場合にはインストール時にMIL0をブートするための小規模(3MBほど)なFATのファイルシステムを作っておきます。
また、SRMを使う場合には「必ず」BSDディスクラベルによってパーティションの設定をしてください。BSDdisklabelはシェルプロンプトからfdiskを手動で起動して「b」コマンドで行います。

ベースシステムさえ展開できれば残りのインストールはx86のインストールとほとんど同じです。
今回インストールしてみてもつくづく感じたのですが、回線が細い場合にはftpベースでのインストールは大変です。とても家から電話回線でftpインストールする気にはならないですね。安定板になったらCDROMが出回ると思われるのでこの雑誌が出版されるころには簡単にCDROMが手に入ると思います。

2. lprobeの概要と利用方法

前述の通り、COMPAQからGPLにしたがってリリースされたlprobeはプロセッサの性能を様々な角度から測定するための有力なツールになります。
このようなツールを一般に公開して頂いたCOMPAQ社に拍手 \(_)/

さて、このツールを使うにはカーネルのバージョンが2.1.132より大きいか2.2.xである必要があります。
テストはSRMでは行われていますが、MIL0からブートしたマシンでは十分には行われていないようです。もっともMIL0に問題があればMIL0の担当者が直す予定だと書かれていますので最新のMIL0を使えば問題ないかもしれません。
前章のMultiaのインストールは実はlprobeを使うために行ったのです。残念ながら誌面の都合で今回はリリースノートから概要の説明に留めます。具体的なインストールや利用方法は機を改めて書くつもりです。
RedHatを使っている人はrpmとコンパイル済みのカーネルモジュールが提供されるのでカーネルのバージョンさえ気を付けておけば簡単にインストールできます。

2.1 性能カウンタと lprobe

まずはこのツールがどんなものかを説明します。
CPUの設計は既存のプログラム(RISCでは特に有名なベンチマークに的を絞って)の定量的な測定を行って設計上のさまざまなトレードオフを決定しています。測定方法の一つはアドレステレースで、これは全部の命令をトレースすることで新しいプロセッサのシミュレーションデータを採取するものです。ところが、この方法には大きな問題があります。というのはトレースを取るにはすべての命令に割り込みをかけて必要な情報を採取するのですが、この時プログラムの動的な振る舞い(特にキャッシュや割り込み)が変わってしまうのです。そこで、実動プログラムの特性を動的に捕らえるためにハードウェアモニタもしくは性能カウンタと呼ぶ機構が導入されることがあります。これはシステムのいろいろな事象をカウントするカウンタになっていてプログラムが動作したときの性能やイベントをカウントしてプログラムの振る舞いを推定するために使われるものです。

この性能カウンタはCPUの設計だけでなくソフトウェアのチューニングにも威力を発揮します。たとえばキャッシュミスや低減するようなアルゴリズムを開発してもどうも性能が上がらないなあというような場合に性能カウンタを用いてキャッシュミスがどの程度減ったのか定量的な評価が可能になります。また、ついでに他の性能要素を観測することでキャッシュミスを対策したのは見当違いだったなんてことがわかる可能性もあります。パイプライン動作や分岐予測、キャッシュの振る舞いや変換バッファのミスなど性能に関するパラメータならたいしては性能カウンタで取得できるように作ります。Alphaにはプロセッサの実装に依存しますが複数の性能カウンタがあります。たとえば今現役で売られている21164(A)のシリーズは3本のカウンタを持っています。そこで、3つまでのパラメータは同時に測定できることとなります。これ以上のパラメータは同時には取得できないので同じプログラムを測定パラメータを換えて何回か流して全体の振る舞いを推定します。実はこのカウンタは一般に使うものではないのでハードウェアの低減のために少しビットを減らしています。その代わり所定の値を超えたときに割り込みを発生することができます。メモリ上にデータを集計するカウンタやヒストグラムを作成しておき割り込みルーチンのなかでこれらの値を更新することでカウンタの長さを超えた長時間の測定が可能になります。

lprobeはこの便利な性能カウンタをユーザーに開放するものです。カウンタの存在すら企業秘密にする会社が多い中で性能カウンタの制御・解析ツールを一般に公開するのは大英断だと思います。

2.2 採取イベント

21164(A)の世代のチップで採取できるイベントを示します。
見れば分かるようなイベントのほうが役に立つものが多いのですが、それぞれのイベントの詳細はハードウェアリファレンスマニュアルを参照してください。

イベント名	カウンタ番号	説明
bcache_misses	2	外部キャッシュのミス数
bcache_hits	1	外部キャッシュのヒット数
bcache_victims	1	外部キャッシュビクティムエントリへのヒット数
branch_instructions	1	分岐命令数
branch_mispredicts	2	分岐予測失敗数
conditional_branches	1	条件分岐命令数
cycles	0	サイクル数
dcache_accesses	1	一次キャッシュアクセス数
dcache_misses	2	一次キャッシュミス
dtb_misses	2	データアドレス変換バッファ(DTB)ミス
dual_issue_cycles	1	2命令実行のサイクル数
external	2	外部事象(ボードによってサポート)
float_operations	1	浮動小数点命令数
icache_accesses	1	命令キャッシュ参照数
icache_misses	2	命令キャッシュミス
integer_operations	1	整数演算命令
itb_misses	2	命令アドレス変換バッファミス
jsr_ret_instructions	1	サブルーチンリターン命令
ldu_replays	2	キャッシュミスを起こしたロードの結果を使う命令の再実行
load_instructions	1	ロード命令数
load_locked	2	ロック付きロード
loads_merged_in_maf	2	MAFにマージされたロード
long_stalls	2	12サイクル以上命令発行が中断されたストール
mem_barrier	2	メモリバリア命令

PC_mispredicts	2	PC 予測失敗
pipeline_dry	1	命令の供給がなくパイプラインが止まるサイクル数
pipeline_frozen	0	資源競合による命令発行の中断
quad_issue_cycles	1	4 命令同時発行数
replay_trap	1	リプレイトラップ
scache_access	1	二次キャッシュ参照
scache_misses	2	二次キャッシュミス
scache_reads	1	二次キャッシュ読み出し要求
scache_read_misses	2	二次キャッシュ読み出しミス
scache_shared_write	2	二次キャッシュへのシェアード書き込み
scache_victims	1	二次キャッシュビクティムヒット
scache_write	2/1	二次キャッシュへの書き込み
scache_write_misses	2	二次キャッシュへの書き込みミス
single_issue_cycles	1	一命令のみ発行できたサイクル
split_issue_cycles	1	命令スロットの部分的命令発行
store_instructions	1	ストア命令数
system_invalidates	2	システムインバリデイト
system_read_requests	1	システム読み出し要求
system_requests	1	システム要求
triple_issue_cycles	1	3 命令発行
wb_maf_full_replays	2	MAF もしくは WB フルによるリプレイトラップ

2.3 実行コマンドと実行時オプション

すぐ使ってみたい人のために実行コマンドと実行時オプションを簡単に列挙します。動作の詳細はドキュメントを参考に実際に動かしてみてください。

lprobe は次のシンタックスで実行します。

```
lprobe [-option [-option..]] [event[.freq][.skip] [event[.freq][.skip]...]]
```

実行時オプションは次のようなものがあります。

```
-buffer_count [number]
-collect
-delay_start [seconds]
-duration [seconds]
-help
-input [file]
-interval [number]
-method [sampling-method *1]
-mode
-nocollect
-output [file]
-pc_range [address1-address2].granularity
```

*1 sampling-method count, sample, mode, ipl, address, total

イベント名は表にあげたものを使います。割り込み頻度を freq で制御しますが、比較的長いプログラムの評価をするときにはデフォルトの low で OK でしょう。もう少し精度の高い測定をしたいときには freq の部分に high を指定してください。ただし、割り込み頻度があがると CPU のタスクスイッチが増えてプログラムの振る舞いが変わる可能性もあります。

使い方の例

```
$lprobe -method count -duration 300 dcache_misses
```

これは 300 秒の間におきるキャッシュミスをカウントして出力します。mode を指定すると CPU のモードごとのデータをヒストグラムとして表示するのでカーネルとユーザーのデータを分けて把握できます。

3. おわりに

Debian-JP プロジェクトのおかげで比較的簡単に最新の日本語環境もそろえられる可能性があります。現在 slink.jp にはすでにかなりの日本語ソフトが入っています。パッケージ管理の自動化や性能関係のライブラリの積極的な取り込みなど注目すべき配布パッケージになってきていると思います。

今回は細い大学の回線から ftp でインストールしてみたのですが、やはり CDROM でパッケージを持っていないとインストールは大変ですね。MSword の件は悩ましいですがきちんとインストールできたら Debian-JP をメインマシンに入れて使ってみたいと思っています。

連載の読者の方から使っていない EB164 のマシンを譲っていただけることになりました。誌面を借りてお礼申し上げます。
これで遅れ馳せながら EV5 の世代のマシンをようやく入手できます。
性能解析のツール lprobe も手に入ったので数値計算関連の性能についてさらに突っ込んだ解析を含んだ話を展開できるだろうと思っています。