

LinuxWorld Expo/Tokyo 2005

Workshop 配布資料

トラブルシューティング

実践で役立つ障害解析のポイント

2005/06/03

MIRACLE LINUX Corporation

Professional Service

Daisuke Tonoki



- **トラブルシューティングのコツ**
- **逆引きトラブルシューティング**

トラブルシューティングのコツ



- MIRACLE LINUX のサポートメニューの紹介
- 無償
 - 無償インストール・サポート
- 年間契約
 - アップデート・サポート
 - プロダクト・サポート
 - エンタープライズ・サポート
- 1ショット
 - インシデント・サポート
 - ダンプ解析サポート



MIRACLEサポート専用のmcinfoコマンド(スクリプト)

- 一度の実行でサポートに必要な情報を収集します。
- 自動的にインストールされます。
- 収集する情報の一例
 - MIRACLE LINUXのバージョン (/etc/miraclelinux-release)
 - 起動時のメッセージ、デバイスの初期化処理など (dmesg)
 - CPUの種類と個数 (/proc/cpuinfo)
 - 実メモリーとSWAPの状態 (/proc/meminfo,/proc/swap)
 - マウントしているデバイス (df)
 - ディスクのドライブ割り当て (fdisk -l)
 - PCIデバイスのリスト (lspci)
 - ロードされているモジュールのリスト (lsmod)
 - インストールされているRPMのリスト (rpm -qa)
 - 最近のsyslog (messages*)



- 障害の情報
 - 現象の内容、何ができて、何ができないのか
 - コンソールに表示されるエラーメッセージ
 - エラーログ `/var/log/アプリケーション名/*log*`
 - Oracleのアラートログ `alert_${SID}.log`

- 現象再現手順
 - どのような状態で、なにをすれば発生するのか

- 障害発生環境の情報
 - ハードウェアの構成
 - O.S.の環境(Version,Kernel,glibc.etc...)
 - アプリケーションの種類、バージョン

逆引きトラブルシューティング



- インストールできない
- SANを接続すると起動できない
- DATドライブが利用できない
- フルリストアしたが起動できない
- CronでのOracleバックアップができない
- サーバにアクセスできない(1)
- サーバにアクセスできない(2)
- LKCDのダンプが取得できない

インストールできない



現象

- ・MIRACLE LINUX をインストール中に以下のメッセージが出てインストールができません。

「エラーが発生しました-新規ファイルシステムを作成する有効なデバイスが見つかりません。
ハードウェアをチェックして、
この問題の原因を調査してください」

その他の確認事項

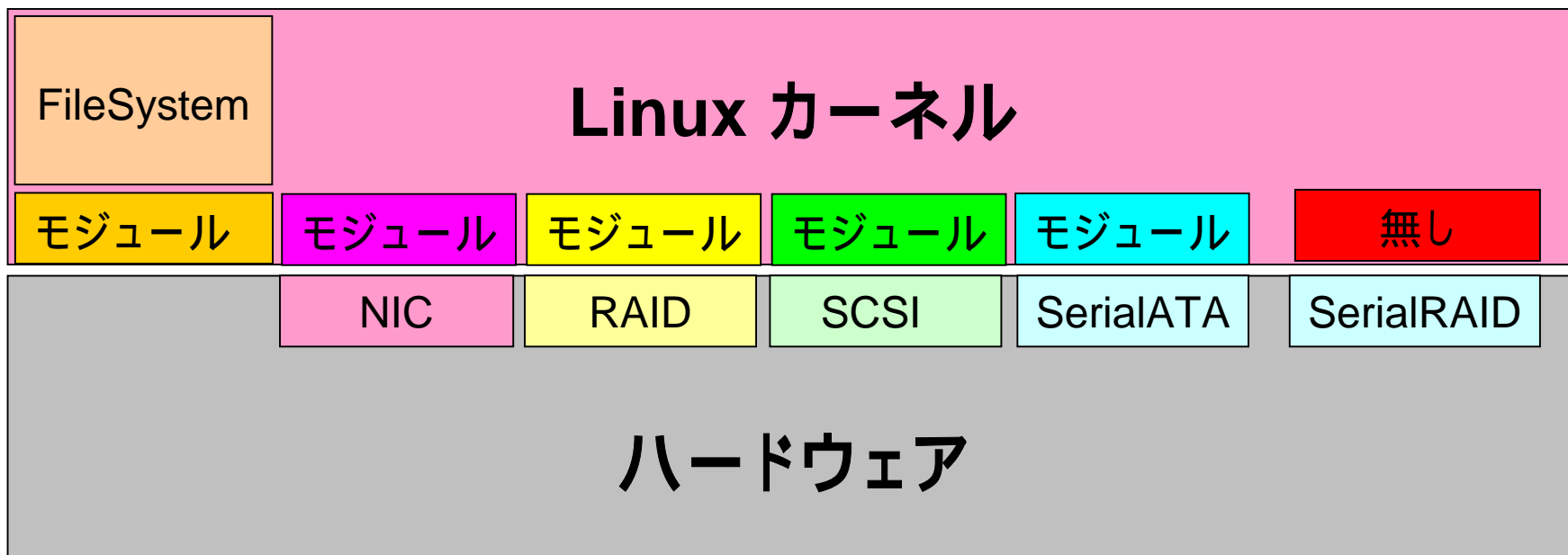
- ・サーバ、利用オプションの詳細な型番確認。
- ・BIOS,RAID構成の確認。
- ・ハードウェアのエラー確認。

インストールできない



カーネル + ハードウェアの関係

LinuxO.S.が特定のハードウェアで動作するかは、LinuxO.S.のカーネルがそのハードウェアに対応するドライバを含んでいるかどうかによって依存します。



インストールできない



原因はこれ

- ・LinuxではサポートされていないRAIDコントローラを利用していました。
- ・今回のケースではICH6RというSerialATAコントローラ内蔵のチップです。ICH5Rなども同様です。
- ・ハードメーカーのページではLinux利用不可と案内されていました。

これで解決

- ・BIOSでオンボードRAID機能をOFFにしました。
- ・Linux標準のSoftwareRAIDでインストールしました。

インストールできない



参考資料SoftwareRAID

現在ICH6RなどIntelチップに内蔵されているオンボードRAIDコントローラが利用できるサーバが増えてきています。

一見ハードウェアRAIDの様ですが、これはドライバがソフトウェア部分を担当するソフトウェアRAIDとなっており、これがドライバがクローズドである理由(RAID実装部がRAIDコントローラーのファームではなく、ソフトウェアとしてドライバに存在してる)となっています。

Iostressテストの実行時間比較

Linux標準・ソフトウェアRAID

ML30 + Linux SoftwareRAID(raid1)

Average Time:55.81

Fastest Time:54.12

Slowest Time:57.60

クローズドソース・オンボードRAID

???? + オンボードRAID(raid1)

Average Time:73.94

Fastest Time:66.56

Slowest Time:78.02

FCストレージを接続すると起動できない



現象

- ・O.S.のインストールが終了し、RACの設定を開始するためFCストレージに接続したが、サーバが起動しなくなってしまった。

その他確認事項

- ・FCストレージの接続ケーブルを外すと起動できる。

ケース1

- ・起動中にカーネルパニックになる。

ケース2

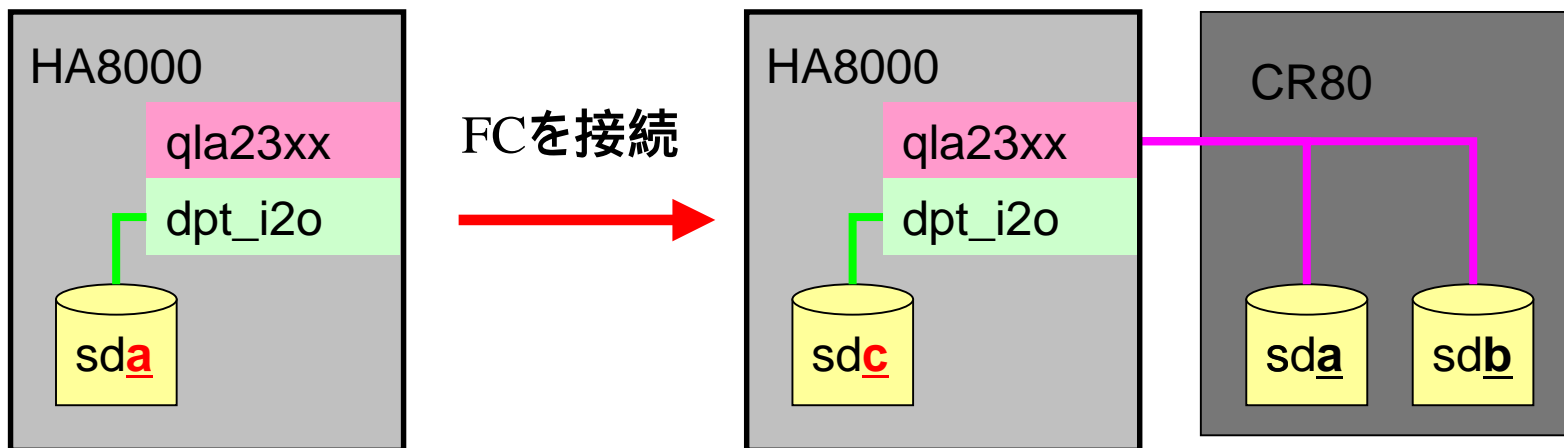
- ・grub (ブートローダー) のメニューも表示されない。

FCストレージを接続すると起動できない



原因はこれ(ケース1)

- ・FCストレージを接続しない状態では内蔵RAIDコントローラーのディスクが sda として認識されています。
- ・コントローラーの順番ディスクの順番よりデバイス名が割り振られるため LABEL を利用していないケースではこのような問題が発生する可能性が高いです。



LABELはパーティションに名前をつけることができ、デバイス名の変化に対応できます。

FCストレージを接続すると起動できない



これで解決(ケース1)

・`/etc/modules.conf`を編集してLinuxが起動する時にドライバが読み込まれる順番を入れ替えました。

```
alias scsi_hostadapter qla23xx  
alias scsi_hostadapter1 megaraid
```



```
alias scsi_hostadapter megaraid  
alias scsi_hostadapter1 qla23xx
```

・このような変更を行った場合は、`initrd`の再作成も必要です。

```
# mkinitrd -f /boot/inird-2.4.21-20.19AXsmp.img 2.4.21-20.19AXsmp
```

FCストレージを接続すると起動できない



原因はこれ(ケース2)

- ・BIOSによるブート設定にて、内蔵RAIDコントローラよりFCコントローラの優先度が高く設定されていたためでした。

```
-Removable Disk
-CD-ROM
-SCSI
- QLogic2300 ←ここ
- Adaptec RAID
```

これで解決(ケース2)

- ・F2プロンプトより呼び出されるBIOS設定にて、内蔵RAIDコントローラの優先度を高く設定しました。

```
-SCSI
- Adaptec RAID ←ここ
- QLogic2300
```

DATが利用できない



現象

- ・O.S.のセットアップが終了し、DATへバックアップを取得しようとしたが、以下のようなエラーが出てDATが利用できません。

```
# tar cf /dev/nst0 /home
tar: /dev/nst0: open 不能: そのようなデバイスはありません
tar: エラーを回復できません: 直ちに終了します。
```

その他確認事項

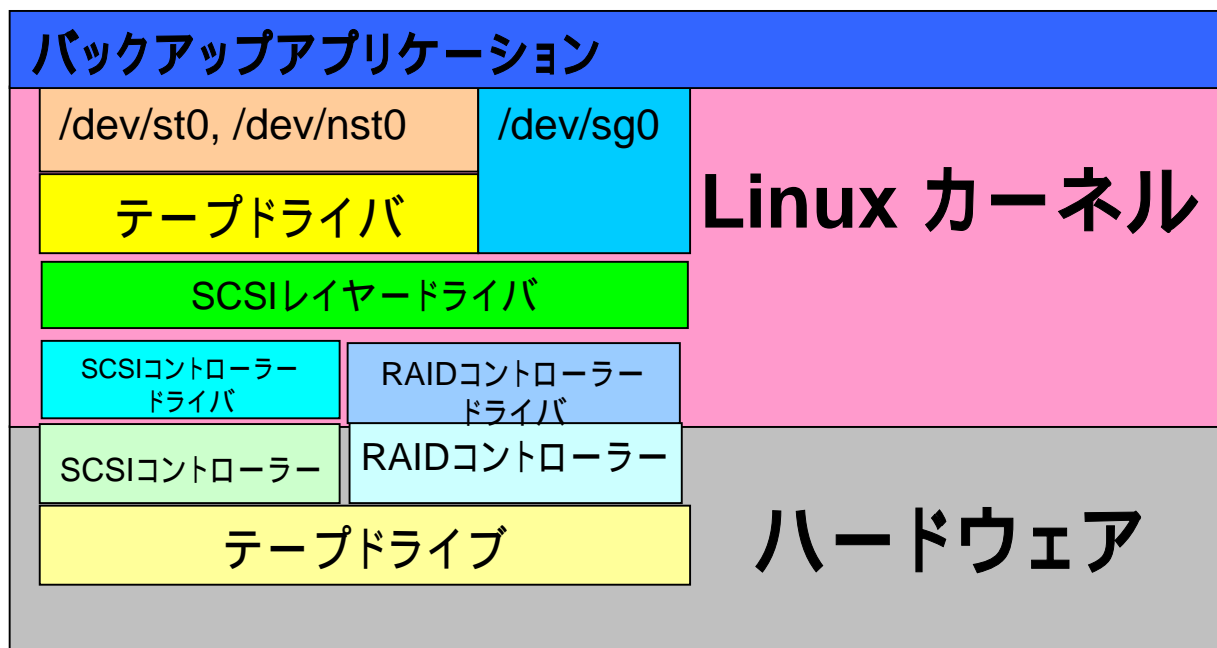
- ・lspciで接続されているSCSIコントローラーの調査。
- ・lsmodでロードされているドライバの調査。

DATが利用できない



LinuxでのDAT(TapeDrive)の利用のイメージ

- テープドライバはst.oで、カーネルに内蔵しています。
- 起動時に自動的に認識してドライバがロードされます。
- 利用するデバイスは、 /dev/st0 , /dev/nst0 , /dev/sg0 など。



DATが利用できない



原因はこれ

- ・テープドライブ接続専用のSCSIコントローラーは増設されておらず、RAIDコントローラーにDATが接続されていました。

これで解決

ケース1

通常のRAIDコントローラーはテープドライブ接続をサポートしていないため、テープドライブ用のSCSIコントローラーを増設しました。

ケース2

HP社のccissドライバを利用するRAIDコントローラーでは、テープドライブ接続をサポートしているため以下のコマンドを実行しました。

```
# echo "engage scsi" > /proc/driver/cciss/cciss0
```

参考URL by HP(COMPAQ)

http://www1.jpn.hp.com/products/software/oe/linux/mainstream/support/doc/option/array/tape_cciss.html

リストア後に起動できない



現象

・dumpコマンドで/, /boot,/homeなどをバックアップ後、リストアのテストを行いましたが、画面に「grub」と表示された時点で停止してしまいます。

その他の確認事項

・リストア後のパーティションをレスキューモードで参照するとデータは正常にリストアされています。

リストア後に起動できない



GRUB (ブートローダー) の動作

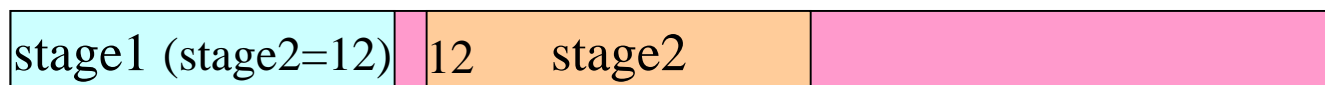
- ・ハードウェアPOST(Power on self test)
- ・MBR(MasterBootRecord)からgrub (stage1,stage2)の読み込み
- ・stage2がgrub.confの読み込み (起動メニューの表示)
- ・カーネルイメージの読み込み (grubがfilesystemにアクセス)
/boot/vmlinuz-\${kernel_version}
- ・必要ならばinitrdのマウント (モジュール読み込みのため)
/boot/initrd-\${kernel_version}.img
- ・rootファイルシステムのマウント
- ・initの実行
- ・各種デーモンの起動

リストア後に起動できない



原因はこれ

- ・MBRに記録されたgrubのstage1にはstage2の物理的な位置が記録されている。
- ・リストアされた/boot/grub/stage2の物理位置が変わっていました。



これで解決

位置がズレている

- ・レスキューモードで起動し、grub-installを実行しました。



grub-installを実行して位置とアドレスを一致させました。

CronでOracleのexpが実行できない



現象

- ・下のようなバックアップスクリプトをcronでスケジュールしましたが、正常に動作しませんでした。

```
#!/bin/sh
/u01/ora8/bin/exp user1/passwd owner=(user1)
file=/home/ora8/backup.dmp
```

その他の確認事項

- ・シェルはoracleユーザーのcrontabに登録されています。
- ・oracleユーザーが、じかに実行した場合は成功します。
- ・dbshut (Oracleの停止スクリプト) はcronでも正常に実行できます。

CronでOracleのexpが実行できない



切り分け

- ・logが残るようにシェルスクリプトを修正しました。

```
#!/bin/sh  
/u01/ora8/bin/exp user1/passwd owner=(user1)  
file=/home/ora8/backup.dmp > /tmp/test.log 2>&1
```

- ・logに記録されたエラーは以下の通りです。

```
Message 206 not found; No message file for product=RDBMS, facility=EXP:  
Release 8.1.7.0.1 - Production on Tue May 6 12:00:00 2003 (c) Copyright  
2000 Oracle Corpor  
Invalid format of Export utility name  
Verify that ORACLE_HOME is properly set  
Export terminated unsuccessfully
```

CronでOracleのexpが実行できない



原因はこれ

- ・oracleユーザーのcrontabに登録されていても、oracleユーザーの環境変数は引き継がれません。
- ・シェルまたはcrontabに環境変数をセットする必要があります。

これで解決

- ・シェルスクリプトにORACLE関連の設定行を追加しました。

```
#!/bin/sh
export ORACLE_HOME=/u01/oracle/product/8.1.7/
export ORACLE_SID=ora8
```

参考情報: dbshutスクリプトには同様の環境変数がセット済みでした。

サーバにアクセスできない(1)



現象

- ・Samba接続でサーバに2GBのファイルを転送するとサーバのネットワークサービスが利用できなくなる。

その他の確認事項

- ・[service network status]、[ifconfig]は正常な値が帰る。
- ・pingには応答がある。
- ・sshや他のネットワークサービスもアクセスできない。
- ・sarのネットワークデータを見るとIPパケット送受信が、現象前にピークとなっている。

```
# sar -n DEV -f /var/log/sa/sa13
```

- ・sarのディスクIOを確認すると直前の負荷は非常に高い。

```
# sar -d -f /var/log/sa/sa13
```


サーバにアクセスできない(1)



原因はこれ

・sarのデータをさらにチェックしていくとpagecache関連の増加とswapの発生が認められ、サーバの反応速度が著しく低下しているものと推測されました。

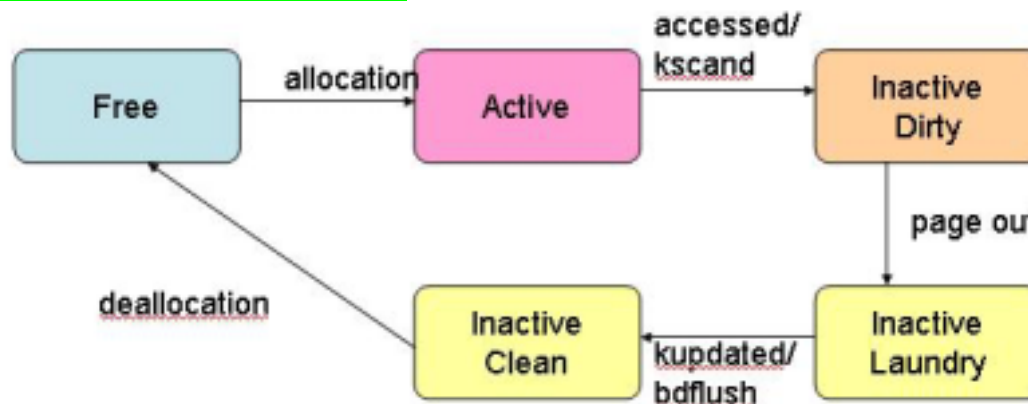
```
# sar -r -f /var/log/sa/sa13
```

00時10分00秒	kmemfree	kmemused	%memused	kmemshrd	kbbuffers	kbcached	kbswpfree	kbswpused	%swpused
01時10分00秒	23568	6037744	99.61	0	2028	5624136	2093220	3252	0.16
02時10分01秒	39784	6021528	99.34	0	2128	5638376	2094988	1484	0.07
03時10分00秒	37984	6023328	99.37	0	1988	5632136	2095292	1180	0.06
04時10分01秒	40240	6021072	99.34	0	57336	5567492	2095292	1180	0.06
05時10分00秒	39928	6021384	99.34	0	1944	5625076	2095296	1176	0.06
06時10分01秒	39328	6021984	99.35	0	1900	5625316	2095296	1176	0.06

サーバにアクセスできない(1)



原因補足ページキャッシュ



[Free],[Active],[InactiveDirty],[InactiveLaundry],[InactiveClean]はページキャッシュの状態を表しています。

1. [Free]は利用可能な状態を表しており、ここからプロセスに配分されたページは[Active]となります。
 2. [Active]はプロセスに配分されている状態を表しており、kscandによりアクセス頻度が低いと判断された場合は[InactiveDirty]状態に移されます。
 3. [InactiveDirty]は一定期間使われていない状態を表しており、メインメモリーからの削除候補となります。
 4. [InactiveLaundry]はディスクへの書き込みを待っている状態で、書き込みが完了した後に再利用可能な[InactiveClean]状態になります。
 5. [InactiveClean]はディスクにsyncされた後の状態を表しており、再利用可能状態です。
- 3,4,5いずれの状態でも一度アクセスがあれば[Active]状態に戻ります。

サーバにアクセスできない(1)



これで解決その1

以下のコマンドを実行してpagecacheパラメーターを変更しました。

```
# echo "10 20 70">/proc/sys/vm/pagecache
```

デフォルト値 1 15 100(単位パーセント)

左から最小値、維持目標値、最大値が実メモリーに占める割合を表しています。注目すべきは右端の最大値で、デフォルトの状態では100%までページキャッシュとして利用されることとなります。OracleなどのBatch系ジョブではディスクデータの再利用が頻繁に行われるため、最大のパフォーマンスを発揮します。ファイルサーバなど頻繁にユーザごとに違うデータが読み書きされるサーバでは、この値を70%程度まで低くするとパフォーマンスが改善できる場合があります。恒久的に設定する方法は/etc/sysctl.confに[vm.pagecache = 1 15 100]と追加します。

サーバにアクセスできない(1)



これで解決その2

以下のコマンドを実行してinactive_clean_percentパラメーターを変更しました。

```
# echo "70">/proc/sys/vm/inactive_clean_percent
```

デフォルト値 30 (単位パーセント)

このパラメーターはすでに使わなくなった[InactiveDirty]状態のページキャッシュの容量に対する、ページキャッシュの再利用可能な[InactiveClean]状態と、再利用可能にするための書き込み待ち[InactiveLaundry]状態の合計容量の割り合いを指定します。

数式で表すと以下の通りです。

$$[\text{InactiveDirty}] \times 30\% < [\text{InactiveClean}] + [\text{InactiveLaundry}]$$

このパラメーターの値を増やすことにより使われないページキャッシュ[InactiveDirty]はより積極的にディスクへ書き込まれ、再利用可能な[InactiveClean]状態の割合が増えます。

恒久的に設定する方法は/etc/sysctl.confに[vm.inactive_clean_percent = 30]と追加します。

サーバにアクセスできない(2)



現象

- ・NFS接続でサーバに1.5GBのファイルを転送するとサーバのネットワーク通信が一切できなくなる。

その他の確認事項

- ・[service network status]、[ifconfig]は正常な値が帰る。
- ・pingには応答がない。
- ・[service network restart] すると回復する。
- ・**sar**のネットワークデータを見るとIPパケット送受信が、すべてが停止している。

```
# sar -n DEV -f /var/log/sa/sa13
```

- ・**sar**のディスクIOを確認すると直前の負荷は低い

```
# sar -d -f /var/log/sa/sa13
```

サーバにアクセスできない(2)



原因はこれ

・パケットの送受信の記録が一切ない状態から、サーバが接続されていたSwitch/Hubとの間で、不正なPAUSEフレームがやり取りされていると推測しました。

これで解決

/etc/modules.confに以下の行を追加しました。

```
options bcm5700 auto_flow_control=0 tx_flow_control=0 rx_flow_control=0
```

Switch/Hub、bcm5700がPAUSEフレームを送信しないよう設定し解決しました。

参考資料はドライバソース添付のREADME

LKCDのdumpが正常に取得できない



現象

- ・サーバがkernel panicを起こすため、LKCDのdumpを取得して調査しているが、dumpファイルが正常に取得できません。

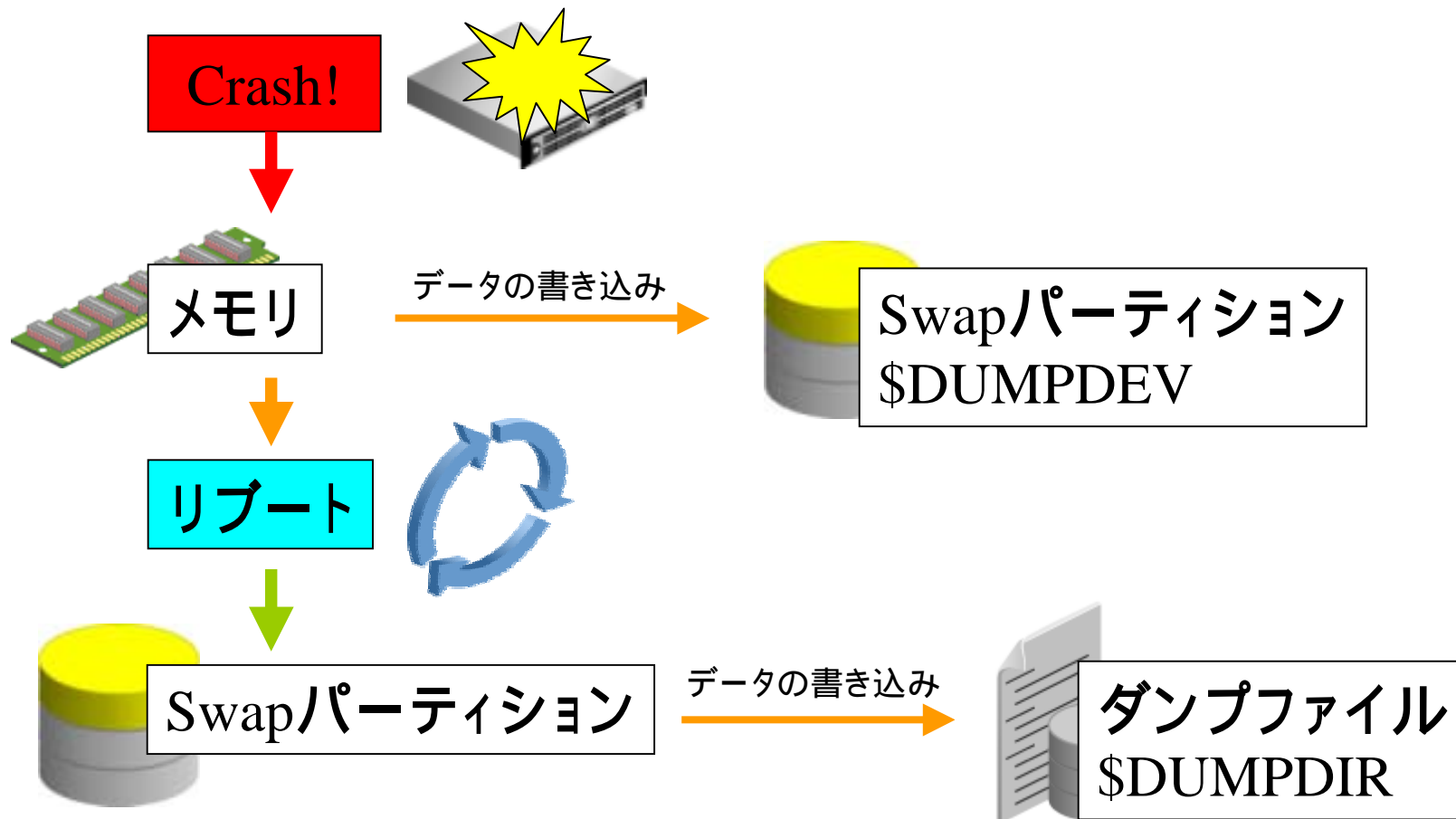
その他の確認事項

- ・サーバが搭載している物理メモリーは8GB。
- ・swapパーティションは2GB。
- ・/varパーティションは1GB。
- ・/varパーティションは現象発生時は100%使用となってしまいます。

LKCDのdumpが正常に取得できない



ダンプファイル取得の流れ



LKCDのdumpが正常に取得できない



原因はこれ

- ・前のページの通りdumpの一次保存先swapパーティション
最終保存先/var/log/dump双方の容量がたりなかった。

これで解決

- ・新たなパーティションを作成し/dev/vmdump(DUMPDEV)に
シンボリックリンクを張りなおした。 通常はswapのリンク
- ・/etc/sysconfig/dumpにあるDUMPDIRパラメーターを、
十分な容量があるパーティション上に変更した。

LKCDのdumpが正常に取得できない



参考資料LKCDとは？

➤ 問題の調査対象

アプリケーションの問題	ログファイル strace coreファイル
パフォーマンスの問題	sar vmstat iostat
O.S.の問題	syslog dmesg /procファイルシステム

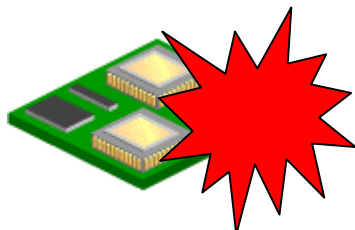
これらのファイルに問題解決の手がかりが残っていない
障害はどのように調査するのか？

LKCDのdumpが正常に取得できない

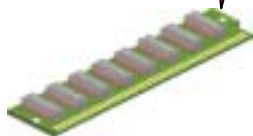


参考資料LKCDとは？

- サーバが停止しても、ダンプが取得できる場合



Linuxが処理中のデータが異常を起こした場合。



Linuxが管理するデータは全てメモリに存在する。



ログに記録を残せない状態なので、強制的にクラッシュダンプを取得する。



ダンプファイルが残るので原因究明が可能対策を講じることができる。



Do the Next, Open your Window

MIRACLE