

AIXCHINA

AIX 问题检测工具和技巧 学习笔记（一）

www.aixchina.com

AIX 中国论坛发表的所有文章版权均属相关权利人所有，受《中华人民共和国著作权法》及其它相关法律的保护。

如出于商业目的使用本资料或有牵涉版权的问题请速与论坛管理员联系。管理员电子邮件：aixchina@21cn.com

AIX Problem Determination Tools and Techniques

IBM Certification Study Guide Test 185

第一章 ICATE 考试科目 185 的学习内容.....	4
1. system dumps	4
2. Crash	4
3. Trace	4
4. 文件系统和性能问题检测工具	4
5. 网络问题检测	5
6. 错误日志和检测	5
7. 其他问题检测工具	5
第二章 客户相关内容	6
1. 确定故障	6
2. 向客户采集信息	6
3. 从机器采集信息	6
第三章 启动问题	7
1. 启动过程	7
2. BIST 和 POST	7
(1) MCA 系统	7
(2) PCI 系统	9
3. 启动第一阶段	9
4. 启动第二阶段	9
(1) LED 551 , 555 , 557	10
(2) LED 552 , 554 , 556	10
(3) LED 518	11
(4) alog	11
5. 启动第三阶段	12
(1) /etc/inittab	12
(2) LED 553	13
(3) LED c31	13
(4) LED 581	13
(5) 错误日志中关于启动的信息	14

第一章 ICATE 考试科目 185 的学习内容

1. system dumps

创建 system dump
理解正确的 system dump 设备
判断 system dump 数据位置
通过 LED 代码确定 system dump 状态
system dump 后采取合适的处理方法
判断 system dump 是否成功
使用 snap 命令

2. Crash

理解 crash 命令的使用和目的
检查 system dump 的状态
用 crash 看堆栈的使用
使用 crash 子命令 stat
进程表里的操作数据
理解 crash 的 stack trace 输出
理解 crash 的进程输出
理解 crash 的 TTY 输出

3. Trace

启动停止 trace
运行 trace
生成 trace 报告
理解 trace 输出
用 trace 调试进程问题

4. 文件系统和性能问题检测工具

用工具检查和修正有问题的文件系统
理解文件系统特性
解决文件系统的 mount 问题
修复有问题的文件系统
vmstat

iostat
filemon

5. 网络问题检测

用问题检测工具找到网络问题
解决网络性能问题
改正 host name 解析问题
NFS mount 问题的原因
路由问题的原因
解决路由问题

6. 错误日志和检测

使用错误日志
理解错误日志
使用检测工具 diag

7. 其他问题检测工具

用 dbx 设置中断点
用 dbx 单步执行程序
用 dbx 带参数执行程序
理解 core 文件
用 core 文件调试程序
读 shell script 脚本程序
调试 shell script 脚本程序

第二章 客户相关内容

1 . 确定故障

一个故障的可能由多个原因导致，比如客户反应不能打印，检查后发现/var 文件系统满了，但扩大/var 后问题可能仍未解决，其他原因包括线缆松动、lpd 守护进程的问题等等都会导致不能打印的现象。因此需要确定客户关心的最终故障，并为之解决。

2 . 向客户采集信息

下面这些问题有助于了解故障的原因：

故障现象是什么？

你是如何发现这个故障的？你是否做了什么不常做的操作才发现这个故障？

其他地方发生了相同的故障吗？是否有多台机器发生了相同的故障？

如果多台机器故障，他们的相同点和不同点在哪里？

最近系统做了什么修改？加了什么新的硬件？软件的配置做了什么修改？

如有做修改，在做修改之前，必要的条件是否都满足？

比如系统微码级别、软件安装的顺序是否正确？硬件的连接、地址设置是否正确？

3 . 从机器采集信息

有用的命令包括 lsdev、lspv、lsvg、lslpp、lsattr、df、mount、errpt，它们可以了解系统的配置情况和错误记录。

如果客户使用 SMIT 或 WSM 管理系统，/smit.log 和/websm.log 包含了这些管理操作的 log 信息。其他象 system dump 和 checkstop 文件常用于解决一些特定的故障。

第三章 启动问题

1 . 启动过程

硬件和软件问题都会导致启动过程挂起，MCA 和 PCI 机器的启动过程有一些重要的不同，启动流程如下：

检查和初始化硬件

MCA 机器包括 BIST (built-in self test) 和 POST (power-on self test) 过程，PCI 机器只有 POST 过程。

加载 BLV 到 RAM 文件系统和移交控制权给 BLV

配置基本的设备

处理/etc/inittab

BLV 内容：

AIX 核心，/unix(软链接到/usr/lib/boot/unix_mp 或 unix_up)就是核心的一个拷贝

Rc.boot 脚本文件

简化的 ODM，只包含基本设备

启动命令，如 cfgmgr、bootinfo

因为 rootvg 现在不可用，所有在内存中创建 RAMFS 文件系统需要的信息都在 BLV 中，在这之后，init 进程加载开始配置基本设备，这是启动的第一阶段。Init 执行 Rc.boot 带参数 1。

启动的第二阶段，尝试激活 rootvg，这是最容易发生问题的一段过程，比如文件系统或 jfslog 的问题等，下一步，控制将交给 rootvg init 进程，RAMFS 释放。

最后，从硬盘加载的 init 进程（不是从 BLV 加载的），执行 rc.boot 带参数 3，配置剩下的设备，最后一步是按照/etc/inittab 的设置运行。这是启动的第三阶段。

2 . BIST 和 POST

(1) MCA 系统

MCA 系统启动时，发生的第一件事就是 BIST，这些测试存在 EPROM 芯片中，主要测试主板上的部件。LED 代码从 100 到 195，代表了硬件的状态。

BIST 后是 POST，POST 测试在加载 boot image 时涉及到的硬件，LED 代码从 200-2E7。这个阶段软硬件问题都会引起停机。

MCA 机器通过钥匙的位置决定是按 Normal 的 bootlist 顺序启动，还是按 Service 的 bootlist 启动。在 Normal 位置启动，init 进程会执行/etc/inittab 中运行级别是 2 的条目。可以用 bootlist 修改 Normal 的启动顺序：

```
#bootlist -m normal hdisk0 hdisk1 rmt0 cd0
```

这个命令表示先在 hdisk0 上查找可用的 BLV，如果没有再查找 hdisk1，等等。

当钥匙放在 service 档，一般用于维护操作时，这时，应用和网络的进程不会启动，bootlist 带 -o 参数还可以列出当前启动顺序。

```
#bootlist -m service -o
```

```
fd0
```

```
cd0
```

```
rmt0
```

```
hdisk2
```

```
ent0
```

AIX4.2 以后可以不用具体指出设备名，而只要指出设备类型，就可以设定 bootlist：

```
#bootlist -m service cd rmt scdisk
```

修改 bootlist 还可以通过 diag 的 Task Selections 的 Display or change Bootlist 来设置，包括设置 normal 和 service。

这一阶段常见错误有：

☞ LED 200

钥匙在 secure 档，要转到 normal 或 service 档。

☞ LED 299

BLV 正在加载，如果 299 过后是 201，且长时间不变，说明要重建 BLV。

☞ 常见的 LED 代码含义：

LED	描述
100-195	BIST 发现硬件问题
200	钥匙在 secure 档
201	如果经过 299，重建 BLV 如果没经过 299，POST 遇到硬件问题
221， 721， 221-229， 223-229， 225-229 233-235	NVRAM 中的 bootlist 不存在，或 bootlist 的设备中没有 botimage，或 bootlist 的设备不可用

☞ 如何重建 BLV

当 BLV 不能访问时，首先通过 diag 检查硬件问题，比如线缆松动。

下一步才通过外部介质，如 AIX CD-ROM 启动维护模式，用 Access this Volume Group 菜单，再用 bosboot -ad /dev/hdisk0 重建 BLV。

另外，镜像 rootvg 并不能在镜像盘上自动建 BLV，这种情况下要建 BLV：

通过外部介质，如 AIX CD-ROM 启动维护模式

Start Maintenance for System Recovery

Access a Root Volume Group

Access this Volume Group and start a shell

Access this Volume Group and start a shell before mounting file system 用于文件系统或 jfslog 有问题的情况下。

(2) PCI 系统

PCI 机器和 MCA 机器有很大的不同，只有 POST 过程，没有 BIST 过程。而且也没有钥匙，新的 PCI 机器使用逻辑开关，通过功能键实现。在早期的 PCI 机器里没有 diag 功能。

☞ 改变 PCI 机器的 bootlist

所有的 PCI 机器有 SMS 菜单，启动到控制台初始化时，通过 F1 (图形控制台) 或 1 (字符控制台) 键，进入 SMS 菜单，其中有 boot 菜单，新的 PCI 机器还有 mutilboot 菜单。

☞ Normal 和 Service 启动顺序

个别机器只有 normal，没有 service 启动顺序，如 7248-43P，可以通过修改 normal 顺序达到维护目的。用命令

```
#bootlist -m service -o
```

```
0514-220 bootlist: Invalid mode(service) for this model
```

可以判断机器是否有 service 模式

所有的 PCI 机器有默认的 bootlist，新的 PCI 机器通过按 F5 键用默认 bootlist 进入单用户模式，或 standalone 的 diag 功能。早期的 PCI 机器没有这样的功能。通过将电池放电 30 秒可以恢复默认值。

F6 键则使用客户自设的 service bootlist。

☞ LED 代码

不同的机器有不同的 LED 代码，要通过机器的 service guide 来查找。

3 . 启动第一阶段

系统在创建 RAMFS，启动 init 进程后，rootvg 还没有激活，从这以后 MCA 和 PCI 机器的启动过程是相同的。

第一阶段流程：

☞ PID 1 – init

☞ Rc.boot 1

Init 进程执行了启动脚本 rc.boot 1，在这阶段，restbase 命令从 BLV 拷贝了简化的 ODM 到 RAMFS，如果失败，LED 出现代码 548。

☞ Cfgmgr -f

从简化的 ODM 中读取 Config_Rules 类，将属性 phase 标记为 1 的基本设备配置好，为激活 rootvg 做好准备。

☞ Bootinfo -b

最后这个命令决定最近的启动设备是什么，这时 LED 显示代码为 511。

4 . 启动第二阶段

第二阶段的流程：

☞ ipl_varyon

激活 rootvg，如果失败，LED 显示代码 552，554，556。

☞ mount /dev/hd4

根文件系统/dev/hd4 mount 到 RAMFS 的临时 mount 点/mnt，如果失败，LED 显示 555，557。

☞ mount /usr

- ☞ `mount /var;copycore;umount /var`
下一步 `mount /usr` 和 `/var` 文件系统，如果失败 LED 显示代码 518，`mount /var` 可以让系统从默认的 dump 设备 `/dev/hd6` 中拷贝 dump 到默认的目录 `/var/adm/ras`。
- ☞ `swapon /dev/hd6`
激活主交换区。
- ☞ 将 `/dev` 从 RAMFS 拷贝到硬盘，`mergedev`
- ☞ 将 ODM 从 RAMFS 拷贝到硬盘，`cp CU* /mnt/etc/objrepos`
- ☞ `umount /usr`，`umount /dev/hd4`
- ☞ `mount -f /`，`mount /usr`，`mount /var`
根文件系统从临时 mount 点移到正确位置后，重新 `mount /usr` 和 `/var`。
- ☞ 拷贝启动信息到 `alog`
因控制台现在不能用，输出信息都记录在 `alog` 中。

(1) LED 551 , 555 , 557

- ☞ 原因：
损坏的文件系统
损坏的 JFSlog 设备
损坏的 rootvg 硬盘
- ☞ 处理：
用介质启动，进入维护菜单，选 `Access a Volume Group and start a shell before mounting filesystems`：
检查文件系统
`#fsck -y /dev/hd1`
`#fsck -y /dev/hd2`
`#fsck -y /dev/hd3`
`#fsck -y /dev/hd4`
`#fsck -y /dev/hd9var`
重建 log 设备
`#/usr/sbin/logform /dev/hd8`
如果 BLV 损坏，重建 BLV，修改 `bootlist`
`#bosboot -a -d /dev/hdisk0`
`#bootlist -m normal hdisk0`

(2) LED 552 , 554 , 556

- ☞ 原因：
损坏的文件系统
损坏的 JFSlog 设备
错误的 IPL 设备记录或 magic 号（magic 号代表设备类型）
BLV 中 ODM 数据库损坏
Rootvg 中的硬盘未激活
损坏的 superblock

☞ 处理：

用介质启动，进入维护菜单，选 Access a Volume Group and start a shell before mounting filesystems：

当 fsck 一个文件系统发现第 8 块有问题时，最简单的修复是，删除这个文件系统，重建，再从备份带上恢复，但是/dev/hd4 不能重建，只能重装 AIX。

当需要修复 ODM，首先将系统配置保存的备份目录中：

```
#/usr/sbin/mount /dev/hd4 /mnt
#/usr/sbin/mount /dev/hd2 /usr
#/usr/bin/mkdir /mnt/etc/objrepos/bak
#/usr/bin/cp /mnt/etc/objrepos/Cu* /mnt/etc/objrepos/bak
#/usr/bin/cp /etc/objrepos/Cu* /mnt/etc/objrepos
#/usr/sbin/umount all
#exit
```

然后，检查 BLV 在硬盘上的位置是否正确：

```
#slsv -m hd5
```

将干净的 ODM 数据库保存到 BLV：

```
#savebase -d /dev/hdisk0
```

最后，重建 BLV：

```
#bosboot -ad /dev/hdisk0
#shutdown -Fr
```

如果是文件系统的 superblock 损坏，fsck 不能自动修复时，可以用命令：

```
#dd count=1 bs=4k skip=31 seek=1 if=/dev/hd4 of=/dev/hd4
```

(3) LED 518

Display Value 518

Remote mount of the / (root) and /usr file systems during network boot did not complete successfully.

启动过程中的518并不是上述错误描述的情况，仍然采用前面介绍的方法修复/usr文件系统。

(4) alog

启动过程中的信息可以通过alog命令查看：

```
# alog -ot boot
***** no stderr *****

-----
Time: 12          LEDS: 0x538
invoking top level program -- "/usr/lib/methods/definet > /dev/null
2>&1;opt=`/u
sr/sbin/lstatr -E -l inet0 -a bootup_option -F value`
```

```

if [ $opt = "no" ];then nf=/etc/rc.net
else nf=/etc/rc.bsdnet
fi;$nf -2;x=$?;test $x -ne 0&&echo $nf failed. Check for invalid
command
s >&2;exit $x"
Time: 21          LEDS: 0x539
return code = 0
***** no stdout *****

```

5 . 启动第三阶段

从 rootvg 中加载的 init 进程，在第三阶段完成以下工作

- ☞ /etc/inittab:/sbin/rc.boot 3
inittab 文件中有
brc:sysinit:/sbin/rc.boot 3 >/dev/console 2>&1 # Phase 3 of system boot
启动 rc.boot 3
- ☞ mount /tmp
- ☞ syncvg rootvg &
LED 显示代码 553。
- ☞ Normal boot:cfgmgr -p2
Cfgmgr 从 ODM 中读 Config_rules 文件，配置 phase=2 的设备。
- ☞ Service boot:cfgmgr -p3
Cfgmgr 从 ODM 中读 Config_rules 文件，配置 phase=3 的设备。
- ☞ Cfgcon
配置控制台，信息将送往控制台，同时也记录/var/adm/ras/conslog，这时 LED 代码的含义如下：
c31: 选择控制台
c32: 控制台是LFT终端
c33: 控制台是tty
c34: 控制台是硬盘上的文件
- ☞ Rc.dt boot
- ☞ Savebase
最后将根文件系统中的 ODM 保存到 BLV 中的 ODM。

(1) /etc/inittab

文件的第一行定义了默认的运行级别，如下：

```
inti:2:initdefault:
```

如果没有这行，启动时，系统会提示输入运行级别。

每行的格式如下：

```
Identifier:RunLevel:Action:Command
```

RunLevel 为空表示任何运行级别都要执行这条命令。

(2) LED 553

553 是由于/etc/inittab 不能读引起的。处理方法先检查/dev/hd3 和/dev/hd4 空间够不够，检查/etc/inittab 文件是否被破坏。/etc/inittab 对文件的格式非常敏感，因此建议用以下命令编辑：

```
mkitab
chitab
```

(3) LED c31

c31 不一定是一个错误码，从 CDROM 或 TAPE 启动时，此时要求选择控制台。

(4) LED 581

581 也不一定就是代表错误，系统这时在配置 TCP/IP 运行/etc/rc.net 配置卡、接口和主机名。但是，但执行/etc/rc.net 时，可能由于系统或网络的问题，会挂起，这个时间从 3 分钟到无限期都有可能。

处理（下面的方法前提是没有用 NIS 和 DNS）：

1. 从 service 模式启动
2. 将/etc/rc.net 改名为/etc/rc.net.save
3. 重启

如果正常了，说明问题可能是：

1. 网卡的硬件故障，用 diag 检查
2. 默认路由不正确
3. 网络不可访问，检查网关、名字服务器、NIS master 等
4. IP 地址或掩码不对，用 iptrace 和 ipreport 命令检查
5. ODM 错误，删去并重建网络设备
6. IP 地址解析问题，检查 named、ypbind/ypserv、/etc/hosts
7. 配置文件的结尾行有多余的空格，用 vi 的 set list 检查
8. LPP 安装配置不正确，重装 LPP
9. DNS 中用到 ATMLE 时，可以在/etc/netsvc.conf 中加上 host=local,bind 或修改/etc/rc.net 文件如下：

```
#####
# Part III - Miscellaneous Commands.
#####
# Set the hostid and uname to `hostname`, where hostname has been
# set via ODM in Part I, or directly in Part II.
# (Note it is not required that hostname, hostid and uname all be
# the same).
export NSORDER="local"      <<=====NEW LINE ADDED HERE
/usr/sbin/hostid `hostname` >> $LOGFILE 2>&1
/bin/uname -S `hostname` | sed 's/\..*$/^' >> $LOGFILE 2>&1
unset NSORDER               <<=====NEW LINE ADDED HERE
```

#####

(5) 错误日志中关于启动的信息

错误日志中记录了启动的信息。

```
# errpt
```

```
IDENTIFIER TIMESTAMP T C RESOURCE_NAME DESCRIPTION
```

```
499B30CC 0711125600 T H ent1 ETHERNET DOWN
```

```
1104AA28 0711125200 T S SYSPROC SYSTEM RESET INTERRUPT RECEIVED
```

```
9DBCfDEE 0711125500 T O errdemon ERROR LOGGING TURNED ON
```

```
499B30CC 0707114100 T H ent1 ETHERNET DOWN
```

```
499B30CC 0707113700 T H ent1 ETHERNET DOWN
```

```
C60BB505 0705101400 P S SYSPROC SW PROGRAM ABNORMALLY TERMINATED
```

```
35BFC499 0705101100 P H cd0 DISK OPERATION ERROR
```

```
0BA49C99 0705101100 T H scsi0 SCSI BUS ERROR
```

```
9DBCfDEE 0704153700 T O errdemon ERROR LOGGING TURNED ON
```

```
192AC071 0704153700 T O errdemon ERROR LOGGING TURNED OFF
```

```
9DBCfDEE 0704152600 T O errdemon ERROR LOGGING TURNED
```

每次启动时，错误日志功能就启动，在上例中，7月4日系统正常关机了两次，错误日志功能也正常关闭，192AC071说明了这个过程。在7月11日重启时，没有停止错误日志的记录，说明系统不是正常关机，12:52系统报了一次reset操作。下面是系统重启记录的详细情况。

```
# errpt -aj 1104AA28
```

```
-----
```

```
LABEL: SYS_RESET
```

```
IDENTIFIER: 1104AA28
```

```
Date/Time: Tue Jul 11 12:52:54
```

```
Sequence Number: 12
```

```
Machine Id: 000BC6DD4C00
```

```
Node Id: server3
```

```
Class: S
```

```
Type: TEMP
```

```
Resource Name: SYSPROC
```

Description

SYSTEM RESET INTERRUPT RECEIVED

Probable Causes

SYSTEM RESET INTERRUPT

Detail Data

KEY MODE SWITCH POSITION AT BOOT TIME

normal

KEY MODE SWITCH POSITION CURRENTLY

normal

AIX中国论坛