# init

init                          Unix                                   process          init
                    **/sbin/init**              init process         Unix
                                        init              init

- root                          init                        run level
- init process
- root                Ctrl-Alt-Del                     init
- daemon
              kill                                                    init

    myrun::ondemand:/home/wzhou/mydaemon
    /home/wzhou/mydaemon                                      kill              init process

-

            init process

```
[root@DEBUG root]# ps aux
USER       PID %CPU %MEM   VSZ   RSS TTY      STAT START   TIME COMMAND
root         1  0.2  0.1  1336   476 ?        S    19:52   0:04 init
root         2  0.0  0.0     0     0 ?        SW   19:52   0:00 [keventd]
root         3  0.0  0.0     0     0 ?        SW   19:52   0:00 [kapmd]
root         4  0.0  0.0     0     0 ?        SWN  19:52   0:00 [ksoftirqd_CPU0]
root         5  0.0  0.0     0     0 ?        SW   19:52   0:00 [kswapd]
root         6  0.0  0.0     0     0 ?        SW   19:52   0:00 [bdflush]
root         7  0.0  0.0     0     0 ?        SW   19:52   0:00 [kupdated]
root         8  0.0  0.0     0     0 ?        SW   19:53   0:00 [mdrecoveryd]
root        16  0.0  0.0     0     0 ?        SW   19:53   0:00 [kjournald]
root        72  0.0  0.0     0     0 ?        SW   19:53   0:00 [khubd]
root       250  0.0  0.2  1396   568 ?        S    19:53   0:00 syslogd -m 0
root       254  0.0  0.1  1336   428 ?        S    19:53   0:00 klogd -x
root       315  0.2  0.2  1500   576 ?        S    19:53   0:04 /usr/sbin/vmware-
root       340  0.0  0.5  3276  1464 ?        S    19:53   0:00 /usr/sbin/sshd
root       352  0.0  0.3  2264  1012 ?        S    19:53   0:00 login -- root
root       353  0.0  0.1  1316   404 tty2     S    19:53   0:00 /sbin/mingetty tt
root       354  0.0  0.1  1316   404 tty3     S    19:53   0:00 /sbin/mingetty tt
root       355  0.0  0.1  1316   404 tty4     S    19:53   0:00 /sbin/mingetty tt
root       356  0.0  0.1  1316   404 tty5     S    19:53   0:00 /sbin/mingetty tt
root       357  0.0  0.1  1316   404 tty6     S    19:53   0:00 /sbin/mingetty tt
root       360  0.0  0.5  4400  1416 tty1     S    19:53   0:00 -bash
root       451  0.0  0.2  2544   628 tty1     R    20:20   0:00 ps aux
[root@DEBUG root]# _
```

# init

init process                            **/etc/inittab**

/etc/inittab

```
[wzhou@dcmp10 ~]$ cat /etc/inittab
#
# inittab       This file describes how the INIT process should set up
#               the system in a certain run-level.
#
# Author:       Miquel van Smoorenburg, <miquels@drinkel.nl.mugnet.org>
#               Modified for RHS Linux by Marc Ewing and Donnie Barnes
#

# Default runlevel. The runlevels used by RHS are:
#   0 - halt (Do NOT set initdefault to this)
#   1 - Single user mode
#   2 - Multiuser, without NFS (The same as 3, if you do not have networking)
#   3 - Full multiuser mode
#   4 - unused
#   5 - X11
#   6 - reboot (Do NOT set initdefault to this)
#
id:5:initdefault:

# System initialization.
si::sysinit:/etc/rc.d/rc.sysinit

l0:0:wait:/etc/rc.d/rc 0
l1:1:wait:/etc/rc.d/rc 1
l2:2:wait:/etc/rc.d/rc 2
l3:3:wait:/etc/rc.d/rc 3
```

```
l4:4:wait:/etc/rc.d/rc 4
l5:5:wait:/etc/rc.d/rc 5
l6:6:wait:/etc/rc.d/rc 6

# Trap CTRL-ALT-DELETE
ca::ctrlaltdel:/sbin/shutdown -t3 -r now

# When our UPS tells us power has failed, assume we have a few minutes
# of power left.  Schedule a shutdown for 2 minutes from now.
# This does, of course, assume you have powerd installed and your
# UPS connected and working correctly.
pf::powerfail:/sbin/shutdown -f -h +2 "Power Failure; System Shutting Down"

# If power was restored before the shutdown kicked in, cancel it.
pr:12345:powerokwait:/sbin/shutdown -c "Power Restored; Shutdown Cancelled"


# Run gettys in standard runlevels
1:2345:respawn:/sbin/mingetty tty1
2:2345:respawn:/sbin/mingetty tty2
3:2345:respawn:/sbin/mingetty tty3
4:2345:respawn:/sbin/mingetty tty4
5:2345:respawn:/sbin/mingetty tty5
6:2345:respawn:/sbin/mingetty tty6

# Run xdm in runlevel 5
x:5:respawn:/etc/X11/prefdm -nodaemon
```

```
id:5:initdefault:
```

                          run level 5    X Window   Full multiuser mode

```
si::sysinit:/etc/rc.d/rc.sysinit
```

sysinit                                (run level)          /etc/rc.d/rc.sysinit
                                  init/main.c    start_kernel()

       /etc/rc.d/rc.sysinit

```
l0:0:wait:/etc/rc.d/rc 0                    run level  0        /etc/rc.d/rc          0
l1:1:wait:/etc/rc.d/rc 1                    run level  1        /etc/rc.d/rc          1
l2:2:wait:/etc/rc.d/rc 2                    run level  2        /etc/rc.d/rc          2
l3:3:wait:/etc/rc.d/rc 3                    run level  3        /etc/rc.d/rc          3
l4:4:wait:/etc/rc.d/rc 4                    run level  4        /etc/rc.d/rc          4
l5:5:wait:/etc/rc.d/rc 5                    run level  5        /etc/rc.d/rc          5
l6:6:wait:/etc/rc.d/rc 6                    run level  6        /etc/rc.d/rc          6
```

/etc/rc.d/rc                                                wait action    init process
        process

```
ca::ctrlaltdel:/sbin/shutdown -t3 -r now
```
                    run level      root        Ctrl+Alt+Del
/sbin/shutdown -t3 -r now
    init        Ctrl+Alt+Del                              shutdown        (now)
warning        3

```
# When our UPS tells us power has failed, assume we have a few minutes
# of power left.  Schedule a shutdown for 2 minutes from now.
# This does, of course, assume you have powerd installed and your
# UPS connected and working correctly.
pf::powerfail:/sbin/shutdown -f -h +2 "Power Failure; System Shutting Down"
```
                                          run level              "powerfail"

```
# If power was restored before the shutdown kicked in, cancel it.
pr:12345:powerokwait:/sbin/shutdown -c "Power Restored; Shutdown Cancelled"
```
         1, 2, 3, 4, 5                    "powerokwait" action          /sbin/shutdown -c "Power Restored; Shutdown
Cancelled"

```
# Run gettys in standard runlevels
1:2345:respawn:/sbin/mingetty tty1
2:2345:respawn:/sbin/mingetty tty2
3:2345:respawn:/sbin/mingetty tty3
```

```
4:2345:respawn:/sbin/mingetty tty4
5:2345:respawn:/sbin/mingetty tty5
6:2345:respawn:/sbin/mingetty tty6
```

6      init process   run level   2 3 4 5                          /sbin/mingetty                        tty1
   tty6

```
[root@DEBUG root]# ps aux
USER        PID %CPU %MEM   VSZ   RSS TTY       STAT START   TIME COMMAND
root          1  0.2  0.1  1336   476 ?         S     19:52   0:04 init
root          2  0.0  0.0     0     0 ?         SW    19:52   0:00 [keventd]
root          3  0.0  0.0     0     0 ?         SW    19:52   0:00 [kapmd]
root          4  0.0  0.0     0     0 ?         SWN   19:52   0:00 [ksoftirqd_CPU0]
root          5  0.0  0.0     0     0 ?         SW    19:52   0:00 [kswapd]
root          6  0.0  0.0     0     0 ?         SW    19:52   0:00 [bdflush]
root          7  0.0  0.0     0     0 ?         SW    19:52   0:00 [kupdated]
root          8  0.0  0.0     0     0 ?         SW    19:53   0:00 [mdrecoveryd]
root         16  0.0  0.0     0     0 ?         SW    19:53   0:00 [kjournald]
root         72  0.0  0.0     0     0 ?         SW    19:53   0:00 [khubd]
root        250  0.0  0.2  1396   568 ?         S     19:53   0:00 syslogd -m 0
root        254  0.0  0.1  1336   428 ?         S     19:53   0:00 klogd -x
root        315  0.2  0.2  1500   576 ?         S     19:53   0:04 /usr/sbin/vmware-
root        340  0.0  0.5  3276  1464 ?         S     19:53   0:00 /usr/sbin/sshd
root        352  0.0  0.3  2264  1012 ?         S     19:53   0:00 login -- root
root        353  0.0  0.1  1316   404 tty2      S     19:53   0:00 /sbin/mingetty tt
root        354  0.0  0.1  1316   404 tty3      S     19:53   0:00 /sbin/mingetty tt
root        355  0.0  0.1  1316   404 tty4      S     19:53   0:00 /sbin/mingetty tt
root        356  0.0  0.1  1316   404 tty5      S     19:53   0:00 /sbin/mingetty tt
root        357  0.0  0.1  1316   404 tty6      S     19:53   0:00 /sbin/mingetty tt
root        360  0.0  0.5  4400  1416 tty1      S     19:53   0:00 -bash
root        451  0.0  0.2  2544   628 tty1      R     20:20   0:00 ps aux
[root@DEBUG root]# _
```

                          6                 root          tty1,              "login -- root",        5
                mingetty

```
# Run xdm in runlevel 5
x:5:respawn:/etc/X11/prefdm -nodaemon
```

run level 5 /etc/X11/prefdm –nodaemon X Window, GUI

inittab init process
■ init process " initdefault" run level 5
■ init process " sysinit" action, /etc/rc.d/rc.sysinit
■ identifier " l5"
l5:5:wait:/etc/rc.d/rc 5
init process " wait" init process inittab action
" /etc/rc.d/rc 5"
■ 6 action
1:2345:respawn:/sbin/mingetty tty1
2:2345:respawn:/sbin/mingetty tty2
3:2345:respawn:/sbin/mingetty tty3
4:2345:respawn:/sbin/mingetty tty4
5:2345:respawn:/sbin/mingetty tty5
6:2345:respawn:/sbin/mingetty tty6
6 run level init process 2 3 4 5 " /sbin/mingetty tty5"
" respawn" /sbin/mingetty process crash ,init process
Linux Alt-F1,…,Alt-F6
tty1, exit, init process " respawn" exit
/sbin/mingetty process init process " /sbin/mingetty tty1"
tty1
■ /etc/X11/prefdm –nodaemon X Window
■ init process
ca::ctrlaltdel:/sbin/shutdown -t3 -r now
pf::powerfail:/sbin/shutdown -f -h +2 "Power Failure; System Shutting Down"
pr:12345:powerokwait:/sbin/shutdown -c "Power Restored; Shutdown Cancelled"

root Ctrl-Alt-Del init process " /sbin/shutdown -t3 -r now"
UPS " /sbin/shutdown -f -h +2 "Power Failure; System Shutting Down""
UPS " /sbin/shutdown -c "Power Restored; Shutdown Cancelled"" init process
root Ctrl-Alt-Del UPS Unix
signal init process signal

# init

init init man init man inittab init process

## init

```
NAME
       init, telinit - process control initialization

SYNOPSIS
       /sbin/init [ -a ] [ -s ] [ -b ] [ -z xxx ] [ 0123456Ss ]
       /sbin/telinit [ -t sec ] [ 0123456sSQqabcUu ]

DESCRIPTION
   Init
       Init  is the parent of all processes.  Its primary role is to create processes
       from a script stored in the file /etc/inittab (see  inittab(5)).   This  file
       usually  has  entries which cause init to spawn gettys on each line that users
       can log in.  It also controls autonomous processes required by any  particular
       system.

RUNLEVELS
       A  runlevel  is  a  software  configuration  of the system which allows only a
       selected group of processes to exist.  The processes spawned by init for  each
       of  these  runlevels are defined in the /etc/inittab file.  Init can be in one
       of eight runlevels: 0â€"6 and S or s.  The runlevel is changed by having a priv-
       ileged  user  run telinit, which sends appropriate signals to init, telling it
       which runlevel to change to.
```

Runlevels 0, 1, and 6 are reserved. Runlevel 0 is used to halt the system, runlevel 6 is used to reboot the system, and runlevel 1 is used to get the system down into single user mode. Runlevel S is not really meant to be used directly, but more for the scripts that are executed when entering runlevel 1. For more information on this, see the manpages for shutdown(8) and inittab(5).

Runlevels 7-9 are also valid, though not really documented. This is because "traditional" Unix variants don’t use them. In case you’re curious, runlevels S and s are in fact the same. Internally they are aliases for the same run-level.

BOOTING

After init is invoked as the last step of the kernel boot sequence, it looks for the file /etc/inittab to see if there is an entry of the type initdefault (see inittab(5)). The initdefault entry determines the initial runlevel of the system. If there is no such entry (or no /etc/inittab at all), a runlevel must be entered at the system console.

Runlevel S or s bring the system to single user mode and do not require an /etc/inittab file. In single user mode, a root shell is opened on /dev/con-sole.

When entering single user mode, init initializes the consoles stty settings to sane values. Clocal mode is set. Hardware speed and handshaking are not changed.

When entering a multi-user mode for the first time, init performs the boot and bootwait entries to allow file systems to be mounted before users can log in. Then all entries matching the runlevel are processed.

When starting a new process, init first checks whether the file /etc/initscript exists. If it does, it uses this script to start the process.

Each time a child terminates, init records the fact and the reason it died in /var/run/utmp and /var/log/wtmp, provided that these files exist.

```
CHANGING RUNLEVELS
        After it has spawned all of the processes specified, init waits for one of its
        descendant processes to die, a powerfail signal, or until it  is  signaled  by
        telinit  to  change the system's runlevel.  When one of the above three condi-
        tions occurs, it re-examines the /etc/inittab file.  New entries can be  added
        to  this  file  at  any  time.  However, init still waits for one of the above
        three conditions to occur.  To provide  for  an  instantaneous  response,  the
        telinit Q or q command can wake up init to re-examine the /etc/inittab file.

        If  init  is not in single user mode and receives a powerfail signal (SIGPWR),
        it reads the file /etc/powerstatus. It then starts a command based on the con-
        tents of this file:

        F(AIL) Power is failing, UPS is providing the power. Execute the powerwait and
               powerfail entries.

        O(K)   The power has been restored, execute the powerokwait entries.

        L(OW)  The power is failing and the UPS has a low battery. Execute the  power-
               failnow entries.

        If  /etc/powerstatus  doesn't exist or contains anything else then the letters
        F, O or L, init will behave as if it has read the letter F.

        Usage of SIGPWR and /etc/powerstatus is discouraged. Someone wanting to inter-
        act  with  init  should  use the /dev/initctl control channel - see the source
        code of the sysvinit package for more documentation about this.

        When init is requested to change the runlevel, it  sends  the  warning  signal
        SIGTERM  to  all  processes  that  are undefined in the new runlevel.  It then
        waits 5 seconds before forcibly terminating these processes  via  the  SIGKILL
        signal.   Note  that  init assumes that all these processes (and their descen-
        dants) remain in the same process group  which  init  originally  created  for
        them.   If  any  process  changes  its  process  group affiliation it will not
        receive these signals.  Such processes need to be terminated separately.

TELINIT
```

```
       /sbin/telinit is linked to /sbin/init.  It takes a one-character argument  and
       signals init to perform the appropriate action.  The following arguments serve
       as directives to telinit:


       0,1,2,3,4,5 or 6
             tell init to switch to the specified run level.


       a,b,c  tell init to process only those /etc/inittab file entries  having  run-
             level a,b or c.


       Q or q tell init to re-examine the /etc/inittab file.


       S or s tell init to switch to single user mode.


       U or u tell  init to re-execute itself (preserving the state). No re-examining
             of /etc/inittab file happens. Run level should be one of Ss12345,  oth-
             erwise request would be silently ignored.


       telinit  can  also tell init how long it should wait between sending processes
       the SIGTERM and SIGKILL signals.  The default is 5 seconds, but  this  can  be
       changed with the -t sec option.


       telinit can be invoked only by users with appropriate privileges.


       The  init binary checks if it is init or telinit by looking at its process id;
       the real init's process id is always 1.  From this it follows that instead  of
       calling telinit one can also just use init instead as a shortcut.

ENVIRONMENT
       Init sets the following environment variables for all its children:


       PATH   /usr/local/sbin:/sbin:/bin:/usr/sbin:/usr/bin


       INIT_VERSION
             As  the  name  says. Useful to determine if a script runs directly from
             init.
```

```
       RUNLEVEL
              The current system runlevel.

       PREVLEVEL
              The previous runlevel (useful after a runlevel switch).

       CONSOLE
              The system console. This is really inherited from the  kernel;  however
              if it is not set init will set it to /dev/console by default.

BOOTFLAGS
       It  is  possible  to pass a number of flags to init from the boot monitor (eg.
       LILO). Init accepts the following flags:

       -s, S, single
              Single user mode boot. In this mode  /etc/inittab  is  examined  and  the
              bootup  rc  scripts  are usually run before the single user mode shell is
              started.

       1-5  Runlevel to boot into.

       -b, emergency
              Boot directly into a single user shell without running any other  startup
              scripts.

       -a, auto
              The  LILO  boot  loader  adds  the  word "auto" to the command line if it
              booted the kernel with the default command line (without  user  interven-
              tion).  If this is found init sets the "AUTOBOOT" environment variable to
              "yes". Note that you cannot use this  for  any  security  measures  -  of
              course  the user could specify "auto" or -a on the command line manually.

       -z xxx
              The argument to -z is ignored. You can use this  to  expand  the  command
              line  a bit, so that it takes some more space on the stack. Init can then
              manipulate the command line so that ps(1) shows the current runlevel.
```

INTERFACE
       Init listens on a fifo in /dev, /dev/initctl, for messages.  Telinit uses this
       to  communicate  with  init. The interface is not very well documented or fin-
       ished. Those interested should study the initreq.h file in the src/  subdirec-
       tory of the init source code tar archive.

SIGNALS
       Init reacts to several signals:

       SIGHUP
           Has the same effect as telinit q.

       SIGUSR1
           On  receipt  of  this signals, init closes and re-opens its control fifo,
           /dev/initctl. Useful for bootscripts when /dev is remounted.

       SIGINT
           Normally the kernel sends  this  signal  to  init  when  CTRL-ALT-DEL  is
           pressed. It activates the ctrlaltdel action.

       SIGWINCH
           The  kernel  sends  this  signal  when the KeyboardSignal key is hit.  It
           activates the kbrequest action.

CONFORMING TO
       Init is compatible with the System V init. It works closely together with  the
       scripts  in the directories /etc/init.d and /etc/rc{runlevel}.d.  If your sys-
       tem uses this convention, there should be  a  README  file  in  the  directory
       /etc/init.d explaining how these scripts work.
FILES
       /etc/inittab
       /etc/initscript
       /dev/console
       /var/run/utmp
       /var/log/wtmp
       /dev/initctl

```
WARNINGS
       Init  assumes  that  processes and descendants of processes remain in the same
       process group which was originally created for them.  If the processes  change
       their  group, init canâ€™t kill them and you may end up with two processes read-
       ing from one terminal line.

DIAGNOSTICS
       If init finds that it is continuously respawning an entry more than  10  times
       in  2  minutes,  it  will assume that there is an error in the command string,
       generate an error message on the system console, and refuse  to  respawn  this
       entry  until  either 5 minutes has elapsed or it receives a signal.  This pre-
       vents it from eating up system resources when someone  makes  a  typographical
       error in the /etc/inittab file or the program for the entry is removed.

AUTHOR
       Miquel  van  Smoorenburg  (miquels@cistron.nl), initial manual page by Michael
       Haardt (u31b3hs@pool.informatik.rwth-aachen.de).

SEE ALSO
       getty(1), login(1),  sh(1),  runlevel(8),  shutdown(8),  kill(1),  inittab(5),
       initscript(5), utmp(5)

                                   18 April 2003                                INIT(8)
```

## /etc/inittab

```
INITTAB(5)                    Linux System Administrator's Manual                    INITTAB(5)

NAME
       inittab - format of the inittab file used by the sysv-compatible init process

DESCRIPTION
       The  inittab  file  describes which processes are started at bootup and during normal operation
       (e.g. /etc/init.d/boot, /etc/init.d/rc, gettys...).  Init(8) distinguishes multiple  runlevels,
```

each of which can have its own set of processes that are started.  Valid runlevels are 0-6 plus
A, B, and C for ondemand entries.  An entry in the inittab file has the following format:

        id:runlevels:action:process

Lines beginning with '#' are ignored.

id      is a unique sequence of 1-4 characters which identifies an entry in  inittab  (for  ver-
        sions of sysvinit compiled with the old libc5 (< 5.2.18) or a.out libraries the limit is
        2 characters).

        Note: traditionally, for getty and other login processes, the value of the id  field  is
        kept  the  same  as  the  suffix of the corresponding tty, e.g. 1 for tty1. Some ancient
        login accounting programs might expect this, though I can't think of any.

runlevels
        lists the runlevels for which the specified action should be taken.

action describes which action should be taken.

process
        specifies the process to be executed.  If the process field starts with a '+' character,
        init  will  not do utmp and wtmp accounting for that process.  This is needed for gettys
        that insist on doing their own utmp/wtmp housekeeping.  This is also a historic bug.

The runlevels field may contain multiple characters for different runlevels.  For example,  123
specifies that the process should be started in runlevels 1, 2, and 3.  The runlevels for onde-
mand entries may contain an A, B, or C.  The runlevels field of  sysinit,  boot,  and  bootwait
entries are ignored.

When  the  system runlevel is changed, any running processes that are not specified for the new
runlevel are killed, first with SIGTERM, then with SIGKILL.

Valid actions for the action field are:

respawn
        The process will be restarted whenever it terminates (e.g. getty).

```
     wait   The process will be started once when the specified runlevel is entered  and  init  will
            wait for its termination.

     once   The process will be executed once when the specified runlevel is entered.

     boot   The process will be executed during system boot.  The runlevels field is ignored.

     bootwait
            The  process  will  be executed during system boot, while init waits for its termination
            (e.g. /etc/rc).  The runlevels field is ignored.

     off    This does nothing.

     ondemand¹
            A process marked with an ondemand runlevel will be executed whenever the specified onde-
            mand runlevel is called.  However, no runlevel change will occur (ondemand runlevels are
            'a', 'b', and 'c').

     initdefault
            An initdefault entry specifies the runlevel which should be entered after  system  boot.
            If  none  exists,  init  will  ask  for  a runlevel on the console. The process field is
            ignored.

     sysinit
            The process will be executed during system boot. It will be executed before any boot  or
            bootwait entries.  The runlevels field is ignored.

     powerwait
            The  process  will  be executed when the power goes down. Init is usually informed about
            this by a process talking to a UPS connected to the computer.  Init will  wait  for  the
            process to finish before continuing.

     powerfail
            As for powerwait, except that init does not wait for the process's completion.
```

---

¹ Ondemand  respawn                          run level

```
       powerokwait
              This  process  will  be  executed as soon as init is informormed that the power has been
              restored.

       powerfailnow
              This process will be executed when init is told that the battery of the external UPS  is
              almost empty and the power is failing (provided that the external UPS and the monitoring
              process are able to detect this condition).

       ctrlaltdel
              The process will be executed when init receives the  SIGINT  signal.   This  means  that
              someone  on  the  system console has pressed the CTRL-ALT-DEL key combination. Typically
              one wants to execute some sort of shutdown either to get into single-user  level  or  to
              reboot the machine.

       kbrequest
              The  process will be executed when init receives a signal from the keyboard handler that
              a special key combination was pressed on the console keyboard.

              The documentation for this function is not complete yet; more documentation can be found
              in  the  kbd-x.xx packages (most recent was kbd-0.94 at the time of this writing). Basi-
              cally you want to map some keyboard combination  to  the  "KeyboardSignal"  action.  For
              example, to map Alt-Uparrow for this purpose use the following in your keymaps file:

              alt keycode 103 = KeyboardSignal

EXAMPLES
       This is an example of a inittab which resembles the old Linux inittab:

              # inittab for linux
              id:1:initdefault:
              rc::bootwait:/etc/rc
              1:1:respawn:/etc/getty 9600 tty1
              2:1:respawn:/etc/getty 9600 tty2
              3:1:respawn:/etc/getty 9600 tty3
              4:1:respawn:/etc/getty 9600 tty4
```

This inittab file executes /etc/rc during boot and starts gettys on tty1-tty4.

A more elaborate inittab with different runlevels (see the comments inside):

```
# Level to run in
id:2:initdefault:

# Boot-time system configuration/initialization script.
si::sysinit:/etc/init.d/rcS

# What to do in single-user mode.
~:S:wait:/sbin/sulogin

# /etc/init.d executes the S and K scripts upon change
# of runlevel.
#
# Runlevel 0 is halt.
# Runlevel 1 is single-user.
# Runlevels 2-5 are multi-user.
# Runlevel 6 is reboot.

l0:0:wait:/etc/init.d/rc 0
l1:1:wait:/etc/init.d/rc 1
l2:2:wait:/etc/init.d/rc 2
l3:3:wait:/etc/init.d/rc 3
l4:4:wait:/etc/init.d/rc 4
l5:5:wait:/etc/init.d/rc 5
l6:6:wait:/etc/init.d/rc 6

# What to do at the "3 finger salute".
ca::ctrlaltdel:/sbin/shutdown -t1 -h now

# Runlevel 2,3: getty on virtual consoles
# Runlevel   3: getty on terminal (ttyS0) and modem (ttyS1)
1:23:respawn:/sbin/getty tty1 VC linux
2:23:respawn:/sbin/getty tty2 VC linux
```

```
              3:23:respawn:/sbin/getty tty3 VC linux
              4:23:respawn:/sbin/getty tty4 VC linux
              S0:3:respawn:/sbin/getty -L 9600 ttyS0 vt320
              S1:3:respawn:/sbin/mgetty -x0 -D ttyS1

FILES
       /etc/inittab

AUTHOR
       Init  was written by Miquel van Smoorenburg (miquels@cistron.nl).  This manual page was written
       by Sebastian Lederer (lederer@francium.informatik.uni-bonn.de) and modified by  Michael  Haardt
       (u31b3hs@pool.informatik.rwth-aachen.de).

SEE ALSO
       init(8), telinit(8)
```

# init

## init process

init process    Linux                                    Linux

start_kernel()

                                         CPU              C                    C
        stack                            start_kernel()    C       main()                        main()
                                              linker                                    main()
        main()    argc  argv                         main()

  start_kernel()   Linux                                            init

```
   483      asmlinkage void __init start_kernel(void)              Linux
   484      {
   485          char * command_line;
   486          extern struct kernel_param __start___param[], __stop___param[];
   487
   488          smp_setup_processor_id();
   489
   490          /*
   491           * Need to run as early as possible, to initialize the
   492           * lockdep hash:
   493           */
   494          unwind_init();
```

```
495             lockdep_init();
496
497             local_irq_disable();
498             early_boot_irqs_off();
499             early_init_irq_lock_class();
500
501     /*
502      * Interrupts are still disabled. Do necessary setups, then
503      * enable them
504      */
505             lock_kernel();
506             boot_cpu_init();
507             page_address_init();
508             printk(KERN_NOTICE);
509             printk(linux_banner);



611             cpuset_init();
612             taskstats_init_early();
613             delayacct_init();
614
615             check_bugs();
616
617             acpi_early_init(); /* before LAPIC and SMP init */
618
619             /* Do the rest non-__init'ed, we're now alive */
620             rest_init();                    Linux
621     }
```

```
416     static void noinline rest_init(void)
417             __releases(kernel_lock)
418     {
419             kernel_thread(init, NULL, CLONE_FS | CLONE_SIGHAND);
420             numa_default_policy();                        Linux          Windows NT
421             unlock_kernel();                                      Linux
```

```
422                                                              init              init process
423              /*
424               * The boot idle thread must execute schedule()
425               * at least one to get things moving:
426               */
427              preempt_enable_no_resched();
428              schedule();
429              preempt_disable();
430
431              /* Call into cpu_idle with preempt disabled */
432              cpu_idle();
433      }
```

```
716      static int init(void * unused)                          initprocess
717      {
718              lock_kernel();
719              /*
720               * init can run on any cpu.
721               */
722              set_cpus_allowed(current, CPU_MASK_ALL);
723              /*
724               * Tell the world that we're going to be the grim
725               * reaper of innocent orphaned children.
726               *
727               * We don't want people to have to make incorrect
728               * assumptions about where in the task array this
729               * can be found.
730               */
731              init_pid_ns.child_reaper = current;
732
733              cad_pid = task_pid(current);
734
735              smp_prepare_cpus(max_cpus);
736
737              do_pre_smp_initcalls();
738
```

```
739            smp_init();
740            sched_init_smp();
741
742            cpuset_init_smp();
743
744            do_basic_setup();
745
746            /*
747             * check if there is an early userspace init.  If yes, let it do all
748             * the work
749             */
750
751            if (!ramdisk_execute_command)
752                    ramdisk_execute_command = "/init";
753
754            if (sys_access((const char __user *) ramdisk_execute_command, 0) != 0) {
755                    ramdisk_execute_command = NULL;
756                    prepare_namespace();
757            }
758
759            /*
760             * Ok, we have completed the initial bootup, and
761             * we're essentially up and running. Get rid of the
762             * initmem segments and start the user-mode stuff..
763             */
764            free_initmem();
765            unlock_kernel();
766            mark_rodata_ro();
767            system_state = SYSTEM_RUNNING;
768            numa_default_policy();
769
770            if (sys_open((const char __user *) "/dev/console", O_RDWR, 0) < 0)
771                    printk(KERN_WARNING "Warning: unable to open an initial console.\n");
772
773            (void) sys_dup(0);
774            (void) sys_dup(0);
775
```

```
776            if (ramdisk_execute_command) {
777                    run_init_process(ramdisk_execute_command);
778                    printk(KERN_WARNING "Failed to execute %s\n",
779                            ramdisk_execute_command);
780            }
781
782            /*
783             * We try each of these until one succeeds.
784             *
785             * The Bourne shell can be used instead of init if we are
786             * trying to recover a really broken machine.
787             */
788            if (execute_command) {
789                    run_init_process(execute_command);
790                    printk(KERN_WARNING "Failed to execute %s.  Attempting "
791                            "defaults...\n", execute_command);
792            }
793        run_init_process("/sbin/init");      run_init_process()
794        run_init_process("/etc/init");          sys_execve()
795        run_init_process("/bin/init");            init process
796        run_init_process("/bin/sh");
797
798            panic("No init found.  Try passing init= option to kernel.");
799    }
```

run_init_process        execve()        init                    " /sbin/init"              " /etc/init" ,
" /bin/init" ,        " /bin/sh"              init                                    4                      ,
                        init        init=/home/wzhou/init

```
710    static void run_init_process(char *init_filename)
711    {
712            argv_init[0] = init_filename;
713            kernel_execve(init_filename, argv_init, envp_init);
714    }
```

```
254        /*
255         * Do a system call from kernel instead of calling sys_execve so we
256         * end up with proper pt_regs.
257         */       sys_execve
258        int kernel_execve(const char *filename, char *const argv[], char *const envp[])
259        {
260                long __res;
261                asm volatile ("push %%ebx ; movl %2,%%ebx ; int $0x80 ; pop %%ebx"
262                : "=a" (__res)
263                : "0" (__NR_execve),"ri" (filename),"c" (argv), "d" (envp) : "memory");
264                return __res;
265        }
```

## init

init                                    init process

init
1.                                                        init
     /etc/inittab
2.               root              init                                    run level                          3 Console         Full
     multiuser mode        root
     # init 1                             Single user mode,                      Windows
             init                 run level                   pipe(    )              request                daemon
     init
3.               init        daemon                      /etc/inittab                                       pipe(    )    2
          run level    request(    )

Linux

Linux      kernel     debugger                                       Linux
" "

   init process                 " init 1"    " init 2"   " init 3"                                   init
init

  init     run level

## init 1

init 1                                         /sbin/init

```
2594  /*
2595   * Main entry for init and telinit.
2596   */
2597  int main(int argc, char **argv)              argc  1  argv[0] = " /sbin/init"
2598  {
2599        char            *p;
2600        int             f;
2601        int             isinit;
2602
2603        /* Get my own name */
2604        if ((p = strrchr(argv[0], '/')) != NULL)      argv[0] = /sbin/init
2605            p++;                                        p     init
2606        else
2607            p = argv[0];
2608        umask(022);
2609
2610        /* Quick check */
2611        if (geteuid() != 0) {                   root        init                  root
```

```
2612                fprintf(stderr, "%s: must be superuser.\n", p);
2613                exit(1);
2614            }
2615
2616        /*
2617         *    Is this telinit or init ?
2618         */
2619        isinit = (getpid() == 1);                        init process  PID  1        isinit = true

2620        for (f = 1; f < argc; f++) {          init 1
2621            if (!strcmp(argv[f], "-i") || !strcmp(argv[f], "--init"))
2622                isinit = 1;
2623                break;
2624        }
2625        if (!isinit) exit(telinit(p, argc, argv));      init 1
2626
2627        /*
2628         *    Check for re-exec
2629         */
2630        if (check_pipe(STATE_PIPE)) {                    init 1   init 3        pipe
2631                                                        init                      check_pipe()
2632            receive_state(STATE_PIPE);                    0  init 0          if          L2646
2633
2634            myname = istrdup(argv[0]);
2635            argv0 = argv[0];
2636            maxproclen = 0;
2637            for (f = 0; f < argc; f++)
2638                maxproclen += strlen(argv[f]) + 1;
2639            reload = 1;
2640            setproctitle("init [%c]",runlevel);
2641
2642            init_main();
2643        }
2644
2645        /* Check command line arguments */              init 1
2646        maxproclen = strlen(argv[0]) + 1;                L2666
```

```
2647          for(f = 1; f < argc; f++) {
2648                  if (!strcmp(argv[f], "single") || !strcmp(argv[f], "-s"))
2649                          dfl_level = 'S';
2650                  else if (!strcmp(argv[f], "-a") || !strcmp(argv[f], "auto"))
2651                          putenv("AUTOBOOT=YES");
2652                  else if (!strcmp(argv[f], "-b") || !strcmp(argv[f],"emergency"))
2653                          emerg_shell = 1;
2654                  else if (!strcmp(argv[f], "-z")) {
2655                          /* Ignore -z xxx */
2656                          if (argv[f + 1]) f++;
2657                  } else if (strchr("0123456789sS", argv[f][0])
2658                          && strlen(argv[f]) == 1)
2659                          dfl_level = argv[f][0];
2660                  /* "init u" in the very beginning makes no sense */
2661                  if (dfl_level == 's') dfl_level = 'S';
2662                  maxproclen += strlen(argv[f]) + 1;
2663          }
2664
2665          /* Start booting. */
2666          argv0 = argv[0];                              init 1   argv0 = /sbin/init
2667          argv[1] = NULL;
2668          setproctitle("init boot");
2669          init_main(dfl_level);                   init 1    init_main(0)  dfl_level          O
2670
2671          /*NOTREACHED*/
2672          return 0;
2673  }
```

OK init 1          init_main()

```
2340  /*
2341   *      The main loop
2342   */
2343  int init_main()
2344  {
2345    CHILD                  *ch;
2346    struct sigaction       sa;
2347    sigset_t          sgt;
```

```
2348    pid_t                    rc;
2349    int              f, st;
2350
2351    if (!reload) {                    init 1       reload             0             if
2352
2353  #if INITDEBUG                    debug init        debug init 1                /sbin/init
2354        /*                                                    debugger
2355         * Fork so we can debug the init process.
2356         */
2357        if ((f = fork()) > 0) {
2358            static const char killmsg[] = "PRNT: init killed.\r\n";
2359            pid_t rc;
2360
2361            while((rc = wait(&st)) != f)
2362                if (rc < 0 && errno == ECHILD)
2363                    break;
2364            write(1, killmsg, sizeof(killmsg) - 1);
2365            while(1) pause();
2366        }
2367  #endif
2368
2369  #ifdef __linux__                    init     ,FreeBSD             Macro        Linux
2370        /*
2371         *    Tell the kernel to send us SIGINT when CTRL-ALT-DEL
2372         *    is pressed, and that we want to handle keyboard signals.
2373         */
2374        init_reboot(BMAGIC_SOFT);                    reboot(BMAGIC_SOFT)        CTRL-ALT-DEL      init
                                                  process  SIGINT signal    man 2 reboot

2375        if ((f = open(VT_MASTER, O_RDWR | O_NOCTTY)) >= 0) {
2376            (void) ioctl(f, KDSIGACCEPT, SIGWINCH);
2377            close(f);
2378        } else
2379            (void) ioctl(0, KDSIGACCEPT, SIGWINCH);
2380  #endif
2381
```

```
2382        /*
2383         *    Ignore all signals.
2384         */
2385        for(f = 1; f <= NSIG; f++)
2386            SETSIG(sa, f, SIG_IGN, SA_RESTART);
2387    }
2388
```

signal handler     ignore        signals
signals   handler

signal

SIGALRM
SIGHUP                        telnet            Linux          top           telnet                    top

SIGINT                  Ctrl-C                        ,        CTRL-ALT-DEL        signal
SIGCHLD
SIGPWR            UPS                        init
SIGWINCH   WINdow CHange
SIGUSR1
SIGSTOP
SIGTSTP                        Ctrl-Z
SIGCONT     SIGSTOP        continue
SIGSEGV

  init process     SIGALRM
                   SIGHUP
                   SIGINT
                   SIGPWR
                   SIGWINCH
                   SIGUSR1        signal_handler    signal_handler            signal
got_signals

  init process     SIGCHLD        chld_handler        init process                    chld_handler

init process    SIGSTOP  SIGTSTP        stop_handler
init process    SIGCONT       cont_handler
init process    SIGSEGV       segv_handler       init process                    init process
              process

```
2389   SETSIG(sa, SIGALRM,  signal_handler, 0);                        init process    signal handler
2390   SETSIG(sa, SIGHUP,   signal_handler, 0);
2391   SETSIG(sa, SIGINT,   signal_handler, 0);
2392   SETSIG(sa, SIGCHLD,  chld_handler, SA_RESTART);
2393   SETSIG(sa, SIGPWR,   signal_handler, 0);
2394   SETSIG(sa, SIGWINCH, signal_handler, 0);
2395   SETSIG(sa, SIGUSR1,  signal_handler, 0);
2396   SETSIG(sa, SIGSTOP,  stop_handler, SA_RESTART);
2397   SETSIG(sa, SIGTSTP,  stop_handler, SA_RESTART);
2398   SETSIG(sa, SIGCONT,  cont_handler, SA_RESTART);
2399   SETSIG(sa, SIGSEGV,  (void (*)(int))segv_handler, SA_RESTART);
```

              signal handler
    SIGALRM  SIGHUP  SIGINT  SIGPWR  SIGWINCH  SIGUSR1      signal_handler

```
543   /*
544    *    We got a signal (HUP PWR WINCH ALRM INT)
545    */
546   void signal_handler(int sig)
547   {
548         ADDSET(got_signals, sig);            HUP PWR WINCH ALRM INT  signal
549   }                                       init_main()    process_signals()
```

got_signals

```
106   sig_atomic_t got_signals;     /* Set if we received a signal. */
```

    ADDSET                  signal

```
#define ADDSET(set, val)   ((set) |= (1 << (val)))
```

          got_signals        signal            process_signals()

```
2238  void process_signals()                  init process    SIGALRM SIGHUP SIGINT SIGPWR SIGWINCH
2239  {                                       SIGUSR1 signal
```

```
2240    CHILD            *ch;
2241    int      pwrstat;
2242    int      oldlevel;
2243    int      fd;
2244    char            c;
2245

2246    if (ISMEMBER(got_signals, SIGPWR)) {          SIGPWR signal    UPS       fail
2247        INITDBG(L_VB, "got SIGPWR");
2248        /* See _what_ kind of SIGPWR this is. */
2249        pwrstat = 0;
2250        if ((fd = open(PWRSTAT, O_RDONLY)) >= 0) {          /etc/powerstatus
2251            c = 0;                                                      " F"   " O"   " L"
2252            read(fd, &c, 1);                              SIGPWR signal        F    fail  O    OK        low

2253            pwrstat = c;                                 fail              pwrstat
2254            close(fd);
2255            unlink(PWRSTAT);                       /etc/powerstatus
2256        }
2257        do_power_fail(pwrstat);            powerfail       family      action
                                       action                inittab    " powerfail"

        pf::powerfail:/sbin/shutdown -f -h +2 "Power Failure; System Shutting Down"
            init process       /sbin/shutdown -f -h +2 "Power Failure; System Shutting Down"

2258        DELSET(got_signals, SIGPWR);    got_signals          SIGPWR
2259    }
2260
2261    if (ISMEMBER(got_signals, SIGINT)) {          SIGINT  signal
2262        INITDBG(L_VB, "got SIGINT");
2263        /* Tell ctrlaltdel entry to start up */
2264        for(ch = family; ch; ch = ch->next)
2265            if (ch->action == CTRLALTDEL)
2266                ch->flags &= ~XECUTED;                Ctrl-Alt-Del handler
2267        DELSET(got_signals, SIGINT);
2268    }
```

```
2269
2270    if (ISMEMBER(got_signals, SIGWINCH)) {
2271        INITDBG(L_VB, "got SIGWINCH");
2272        /* Tell kbrequest entry to start up */
2273        for(ch = family; ch; ch = ch->next)
2274            if (ch->action == KBREQUEST)
2275                ch->flags &= ~XECUTED;
2276        DELSET(got_signals, SIGWINCH);
2277    }
2278
2279    if (ISMEMBER(got_signals, SIGALRM)) {                    SIGALRM signal
2280        INITDBG(L_VB, "got SIGALRM");
2281        /* The timer went off: check it out */
2282        DELSET(got_signals, SIGALRM);
2283    }
2284
2285    if (ISMEMBER(got_signals, SIGCHLD)) {              init process
2286        INITDBG(L_VB, "got SIGCHLD");
2287        /* First set flag to 0 */
2288        DELSET(got_signals, SIGCHLD);                          signal
2289
2290        /* See which child this was */
2291        for(ch = family; ch; ch = ch->next)          family
2292            if (ch->flags & ZOMBIE) {
2293                INITDBG(L_VB, "Child died, PID= %d", ch->pid);
2294                ch->flags &= ~(RUNNING|ZOMBIE|WAITING);
2295                if (ch->process[0] != '+')
2296                    write_utmp_wtmp("", ch->id, ch->pid, DEAD_PROCESS, NULL);
2297            }
2298
2299    }
2300
2301    if (ISMEMBER(got_signals, SIGHUP)) {    signal SIGHUP                                init
2302        INITDBG(L_VB, "got SIGHUP");        process      inittab
2303 #if CHANGE_WAIT
2304        /* Are we waiting for a child? */
```

```
2305         for(ch = family; ch; ch = ch->next)
2306                 if (ch->flags & WAITING) break;
2307         if (ch == NULL)
2308 #endif
2309         {
2310                 /* We need to go into a new runlevel */
2311                 oldlevel = runlevel;
2312 #ifdef INITLVL
2313                 runlevel = read_level(0);
2314 #endif
2315                 if (runlevel == 'U') {
2316                         runlevel = oldlevel;
2317                         re_exec();
2318                 } else {
2319                         if (oldlevel != 'S' && runlevel == 'S') console_stty();
2320                         if (runlevel == '6' || runlevel == '0' ||
2321                             runlevel == '1') console_stty();
2322                         read_inittab();
2323                         fail_cancel();
2324                         setproctitle("init [%c]", runlevel);
2325                         DELSET(got_signals, SIGHUP);
2326                 }
2327         }
2328     }
2329     if (ISMEMBER(got_signals, SIGUSR1)) {                    signal          pipe     /dev/initctl
2330         /*
2331          *    SIGUSR1 means close and reopen /dev/initctl
2332          */
2333         INITDBG(L_VB, "got SIGUSR1");
2334         close(pipe_fd);
2335         pipe_fd = -1;
2336         DELSET(got_signals, SIGUSR1);
2337     }
2338 }
```

init_main()

```
2400
2401    console_init();                console
2402
2403    if (!reload) {                          init 1         0                  if
2404
2405        /* Close whatever files are open, and reset the console. */
2406        close(0);                           0  1  2      setsid()    Linux  Daemon           init
2407        close(1);              daemon process
2408        close(2);
2409        console_stty();
2410        setsid();
2411
2412        /*
2413         *    Set default PATH variable.
2414         */
2415        putenv(PATH_DFL);            init 1                  PATH="PATH=/bin:/usr/bin:/sbin:/usr/sbin"
2416
2417        /*
2418         *    Initialize /var/run/utmp (only works if /var is on
2419         *    root and mounted rw)
2420         */
2421        (void) close(open(UTMP_FILE, O_WRONLY|O_CREAT|O_TRUNC, 0644));       /var/run/utmp
2422
2423        /*
2424         *    Say hello to the world
2425         */
2426        initlog(L_CO, bootmsg, "booting");
2427
2428        /*
2429         *    See if we have to start an emergency shell.
2430         */
2431        if (emerg_shell) {                         shell    init 1                    if
2432            SETSIG(sa, SIGCHLD, SIG_DFL, SA_RESTART);
2433            if (spawn(&ch_emerg, &f) > 0) {
2434                while((rc = wait(&st)) != f)
2435                    if (rc < 0 && errno == ECHILD)
```

```
2436                               break;
2437               }
2438               SETSIG(sa, SIGCHLD,  chld_handler, SA_RESTART);
2439        }
2440
2441        /*
2442         *    Start normal boot procedure.
2443         */
2444        runlevel = '#';                              init 1
2445        read_inittab();             /etc/inittab
2446
2447    } else {  re-exec      init process           inittab          init 1              L2455
2448        /*
2449         *    Restart: unblock signals and let the show go on
2450         */
2451        initlog(L_CO, bootmsg, "reloading");
2452        sigfillset(&sgt);
2453        sigprocmask(SIG_UNBLOCK, &sgt, NULL);
2454    }
2455    start_if_needed();              family            /etc/inittab       action
2456                                          inittab
2457    while(1) {           init process       init process                         init 3
2458                       daemon process      init
2459        /* See if we need to make the boot transitions. */
2460        boot_transitions();
2461        INITDBG(L_VB, "init_main: waiting..");
2462
2463        /* Check if there are processes to be waited on. */
2464        for(ch = family; ch; ch = ch->next)
2465         if ((ch->flags & RUNNING) && ch->action != BOOT) break;



2466
2467  #if CHANGE_WAIT
```

```
2468          /* Wait until we get hit by some signal. */
2469          while (ch != NULL && got_signals == 0) {     daemon process
2470           if (ISMEMBER(got_signals, SIGHUP)) {                    SIGHUP       init process     inittab
2471                  /* See if there are processes to be waited on. */
2472               for(ch = family; ch; ch = ch->next)                  ″ wait″          ″ boot″
2473                    if (ch->flags & WAITING) break;
2474           }
2475           if (ch != NULL) check_init_fifo();
2476          }
2477  #else /* CHANGE_WAIT */
2478       if (ch != NULL && got_signals == 0) check_init_fifo();
2479  #endif /* CHANGE_WAIT */
```

check_init_fifo()        init 2   init 3             ″ /dev/initctl″ pipe

```
2480
2481       /* Check the 'failing' flags */
2482       fail_check();
2483
2484       /* Process any signals. */
2485       process_signals();                        signal       init process     signal        got_signals
2486                                        init process                          init process     signal
                                                 signal               L2457
2487       /* See what we need to start up (again) */
2488       start_if_needed();                 family        node
2489     }                                    (action)
2490     /*NOTREACHED*/
2491  }
```

init process                    inittab                          Inittab
init process

```
/* Information about a process in the in-core inittab */
typedef struct _child_ {
  int flags;                 /* Status of this entry */
  int exstat;                /* Exit status of process */
```

```
   int pid;                      /* Pid of this process */
   time_t tm;                    /* When respawned last */
   int count;                    /* Times respawned in the last 2 minutes */
   char id[8];                   /* Inittab id (must be unique) */
   char rlevel[12];              /* run levels */
   int action;                   /* what to do (see list below) */
   char process[128];            /* The command line */
   struct _child_ *new;          /* New entry (after inittab re-read) */
   struct _child_ *next;         /* For the linked list */
} CHILD;
```

Inittab

3:2345:respawn:/sbin/mingetty tty3

    init process         CHILD                           field

| id[8] | rlevel[12] | action | process[128] |
|-------|------------|--------|--------------|
| 3 | 2345 | Respawn | /sbin/mingetty tty3 |

flags        process

exstat     process                              exit()

pid      process   process identifier

tm   respawn      process               init process respawn

count   respawn   ondemanded     process               2           (spawn)


init process            process    next               family

family

next

CHILD Node → CHILD Node → CHILD Node → [ ] → CHILD Node N

/etc/inittab

```
id:5:initdefault:                    CHILD Node

si::sysinit:/etc/rc.d/rc.sysinit

                                     CHILD Node

l0:0:wait:/etc/rc.d/rc 0     CHILD Node
l1:1:wait:/etc/rc.d/rc 1
l2:2:wait:/etc/rc.d/rc 2
l3:3:wait:/etc/rc.d/rc 3
l4:4:wait:/etc/rc.d/rc 4
l5:5:wait:/etc/rc.d/rc 5
l6:6:wait:/etc/rc.d/rc 6


ca::ctrlaltdel:/sbin/shutdown -t3 -r now

3:2345:respawn:/sbin/mingetty tty3  CHILD Node N
```

read_inittab
inittab
family
  node ,

```
1108  void read_inittab(void)
1109  {
1110    FILE            *fp;                /* The INITTAB file */
1111    CHILD           *ch, *old, *i;     /* Pointers to CHILD structure */
1112    CHILD           *head = NULL;      /* Head of linked list */
1113  #ifdef INITLVL
1114    struct stat     st;                /* To stat INITLVL */
1115  #endif
1116    sigset_t  nmask, omask;            /* For blocking SIGCHLD. */
1117    char            buf[256];          /* Line buffer */
1118    char            err[64];           /* Error message. */
1119    char            *id, *rlevel,
1120            *action, *process;         /* Fields of a line */
1121    char            *p;
1122    int        lineNo = 0;             /* Line number in INITTAB file */
1123    int        actionNo;               /* Decoded action field */
1124    int        f;                      /* Counter */
1125    int        round;                  /* round 0 for SIGTERM, 1 for SIGKILL */
1126    int        foundOne = 0;           /* No killing no sleep */
1127    int        talk;                   /* Talk to the user */
1128    int        done = 0;               /* Ready yet? */
1129
1130  #if DEBUG
1131    if (newFamily != NULL) {
1132        INITDBG(L_VB, "PANIC newFamily != NULL");
1133        exit(1);
1134    }
1135    INITDBG(L_VB, "Reading inittab");
1136  #endif
1137
1138    /*
1139     *  Open INITTAB and real line by line.
1140     */
1141    if ((fp = fopen(INITTAB, "r")) == NULL)                    /etc/inittab
1142        initlog(L_VB, "No inittab file found");
1143
```

```
1144    while(!done) {                                    inittab        newFamily          newFamily
1145        /*                          family
1146         *    Add single user shell entry at the end.
1147         */
1148        if (fp == NULL || fgets(buf, sizeof(buf), fp) == NULL) {
1149            done = 1;
1150            /*
1151             *    See if we have a single user entry.
1152             */
1153            for(old = newFamily; old; old = old->next)
1154                if (strpbrk(old->rlevel, "S")) break;
1155            if (old == NULL)
1156                snprintf(buf, sizeof(buf), "~~:S:wait:%s\n", SULOGIN);
1157            else
1158                continue;
1159        }
1160        lineNo++;
1161        /*
1162         *    Skip comments and empty lines
1163         */
1164        for(p = buf; *p == ' ' || *p == '\t'; p++)
1165            ;
1166        if (*p == '#' || *p == '\n') continue;              " #"
1167
1168        /*
1169         *    Decode the fields
1170         */
        id:runlevels:action:process    4
1171        id =     strsep(&p, ":");                        " :"                strsep
1172        rlevel =  strsep(&p, ":");
1173        action =  strsep(&p, ":");
1174        process = strsep(&p, "\n");
1175
                        init manual                                     127
1176        /*
1177         *    Check if syntax is OK. Be very verbose here, to
```

```
1178              *     avoid newbie postings on comp.os.linux.setup :)
1179              */
1180             err[0] = 0;
1181             if (!id || !*id) strcpy(err, "missing id field");
1182             if (!rlevel)     strcpy(err, "missing runlevel field");
1183             if (!process)    strcpy(err, "missing process field");
1184             if (!action || !*action)
1185                     strcpy(err, "missing action field");
1186             if (id && strlen(id) > sizeof(utproto.ut_id))
1187                 sprintf(err, "id field too long (max %d characters)",
1188                         (int)sizeof(utproto.ut_id));
1189             if (rlevel && strlen(rlevel) > 11)
1190                 strcpy(err, "rlevel field too long (max 11 characters)");
1191             if (process && strlen(process) > 127)
1192                 strcpy(err, "process field too long");
1193             if (action && strlen(action) > 32)
1194                 strcpy(err, "action field too long");
1195             if (err[0] != 0) {
1196                 initlog(L_VB, "%s[%d]: %s", INITTAB, lineNo, err);
1197                 INITDBG(L_VB, "%s:%s:%s:%s", id, rlevel, action, process);
1198                 continue;
1199             }
1200
1201             /*
1202              *    Decode the "action" field
1203              */
        init     action          actions[]                                      identifier
1204             actionNo = -1;
1205             for(f = 0; actions[f].name; f++)
1206                 if (strcasecmp(action, actions[f].name) == 0) {
1207                     actionNo = actions[f].act;
1208                     break;
1209                 }
1210             if (actionNo == -1) {                          action(    actions[]        )
1211                 initlog(L_VB, "%s[%d]: %s: unknown action field",
1212                         INITTAB, lineNo, action);
1213                 continue;
```

```
1214          }
1215
1216          /*
1217           *    See if the id field is unique
1218           */

                    identifier                                                    CHILD
                                                          id                              /etc/inittab
    id

1219          for(old = newFamily; old; old = old->next) {
1220                  if(strcmp(old->id, id) == 0 && strcmp(id, "~~")) {
1221                          initlog(L_VB, "%s[%d]: duplicate ID field \"%s\"",
1222                                  INITTAB, lineNo, id);
1223                          break;
1224                  }
1225          }
1226          if (old) continue;
1227
1228          /*
1229           *    Allocate a CHILD structure
1230           */
1231          ch = imalloc(sizeof(CHILD));                        CHILD node
1232
1233          /*
1234           *    And fill it in.                               CHILD node
1235           */
1236          ch->action = actionNo;          action
1237          strncpy(ch->id, id, sizeof(utproto.ut_id) + 1); /* Hack for different libs. */
1238          strncpy(ch->process, process, sizeof(ch->process) - 1);
1239          if (rlevel[0]) {                          run level
1240                  for(f = 0; f < sizeof(rlevel) - 1 && rlevel[f]; f++) {
1241                          ch->rlevel[f] = rlevel[f];
1242                          if (ch->rlevel[f] == 's') ch->rlevel[f] = 'S';
1243                  }
1244                  strncpy(ch->rlevel, rlevel, sizeof(ch->rlevel) - 1);
```

```
1245            } else {                    run level          run level              process
1246                    strcpy(ch->rlevel, "0123456789");
1247                    if (ISPOWER(ch->action))
1248                            strcpy(ch->rlevel, "S0123456789");
1249            }
                action
1250            /*
1251             *    We have the fake runlevel '#' for SYSINIT  and
1252             *    '*' for BOOT and BOOTWAIT.
1253             */

                        SYSINIT action   '#'       BOOT action   '*'                  run level   0   9    'S'
        " #"   " *"            run level           SYSINIT                          BOOT    action

1254            if (ch->action == SYSINIT) strcpy(ch->rlevel, "#");
1255            if (ch->action == BOOT || ch->action == BOOTWAIT)
1256                    strcpy(ch->rlevel, "*");
1257
1258            /*
1259             *    Now add it to the linked list. Special for powerfail.
1260             */

    /etc/inittab                         newFamily                       family
   init      run level      family                         init           /etc/inittab

1261            if (ISPOWER(ch->action)) {        action          POWERWAIT  POWERFAIL  POWEROKWAIT
1262                                             POWERFAILNOW  CTRLALTDEL                      Ctrl+Alt+Del

1263                    /*
1264                     *    Disable by default
1265                     */
1266                    ch->flags |= XECUTED;                     startup()                      action   flag
1267                                             XECUTED                                      ISPOWER()   action
1268                    /*
                                                 Ctrl+Alt+Del                       /etc/inittab   CTRLALTDEL
```

```
                                  Action       process              disable  (        XECUTED
                                  )            Ctrl+Alt+Del          enable
                                         action       family


1269            *     Tricky: insert at the front of the list..
1270            */
1271          old = NULL;
1272          for(i = newFamily; i; i = i->next) {
1273                  if (!ISPOWER(i->action)) break;
1274                  old = i;
1275          }
1276          /*
1277           *    Now add after entry "old"
1278           */
1279          if (old) {
1280                  ch->next = i;
1281                  old->next = ch;
1282                  if (i == NULL) head = ch;
1283          } else {
1284                  ch->next = newFamily;
1285                  newFamily = ch;
1286                  if (ch->next == NULL) head = ch;
1287          }
1288      } else {        action          KBREQUEST              init manual    SIGWINCH  signal
1289          /*
1290           *    Just add at end of the list
1291           */
1292          if (ch->action == KBREQUEST) ch->flags |= XECUTED;
1293          ch->next = NULL;
1294          if (head)
1295                  head->next = ch;
1296          else
1297                  newFamily = ch;
1298          head = ch;
1299      }
1300
```

```
1301          /*
1302           *    Walk through the old list comparing id fields
1303           */
1304          for(old = family; old; old = old->next)
1305                  if (strcmp(old->id, ch->id) == 0) {
1306                          old->new = ch;
1307                          break;
1308                  }

1309      }
1310      /*
1311       *  We're done.
1312       */
1313      if (fp) fclose(fp);         /etc/inittab          init 1
1314                                       daemon       init   init 3
1315      /*                    kernel     init(init 1)      daemon          init  init 3
1316       *  Loop through the list of children, and see if they need to
1317       *  be killed.
1318       */
1319
1320      INITDBG(L_VB, "Checking for children to kill");
1321      for(round = 0; round < 2; round++) {
1322        talk = 1;
1323        for(ch = family; ch; ch = ch->next) {                    init  init 1          family
1324          ch->flags &= ~KILLME;                            round = 0  talk = 1  foundOne = 0
1325                                                   L1393



1326          /*
1327           *    Is this line deleted?
1328           */
1329          if (ch->new == NULL) ch->flags |= KILLME;
```

```
1330
1331        /*
1332         *    If the entry has changed, kill it anyway. Note that
1333         *    we do not check ch->process, only the "action" field.
1334         *    This way, you can turn an entry "off" immediately, but
1335         *    changes in the command line will only become effective
1336         *    after the running version has exited.
1337         */
1338        if (ch->new && ch->action != ch->new->action) ch->flags |= KILLME;
1339
1340        /*
1341         *    Only BOOT processes may live in all levels
1342         */
1343        if (ch->action != BOOT &&
1344            strchr(ch->rlevel, runlevel) == NULL) {
1345            /*
1346             *    Ondemand procedures live always,
1347             *    except in single user
1348             */
1349            if (runlevel == 'S' || !(ch->flags & DEMAND))
1350                    ch->flags |= KILLME;
1351        }
1352
1353        /*
1354         *    Now, if this process may live note so in the new list
1355         */
1356        if ((ch->flags & KILLME) == 0) {
1357            ch->new->flags  = ch->flags;
1358            ch->new->pid    = ch->pid;
1359            ch->new->exstat = ch->exstat;
1360            continue;
1361        }
1362
1363
1364        /*
1365         *    Is this process still around?
1366         */
```

```
1367          if ((ch->flags & RUNNING) == 0) {
1368                  ch->flags &= ~KILLME;
1369                  continue;
1370          }
1371          INITDBG(L_VB, "Killing \"%s\"", ch->process);
1372          switch(round) {
1373                  case 0: /* Send TERM signal */
1374                          if (talk)
1375                                  initlog(L_CO,
1376                                          "Sending processes the TERM signal");
1377                          kill(-(ch->pid), SIGTERM);
1378                          foundOne = 1;
1379                          break;
1380                  case 1: /* Send KILL signal and collect status */
1381                          if (talk)
1382                                  initlog(L_CO,
1383                                          "Sending processes the KILL signal");
1384                          kill(-(ch->pid), SIGKILL);
1385                          break;
1386          }
1387          talk = 0;
1388
1389     }
1390     /*
1391      *      See if we have to wait 5 seconds
1392      */
1393     if (foundOne && round == 0) {                            init  init 1    round = 0,  foundOne = 0
1394         /*                                            if           L1419
1395          *   Yup, but check every second if we still have children.
1396          */
1397         for(f = 0; f < sltime; f++) {
1398                 for(ch = family; ch; ch = ch->next) {
1399                         if (!(ch->flags & KILLME)) continue;
1400                         if ((ch->flags & RUNNING) && !(ch->flags & ZOMBIE))
1401                                 break;
1402                 }
1403                 if (ch == NULL) {
```

```
1404                    /*
1405                     *     No running children, skip SIGKILL
1406                     */
1407                    round = 1;
1408                    foundOne = 0; /* Skip the sleep below. */
1409                    break;
1410                }
1411            do_sleep(1);
1412        }
1413    }
1414    }
1415
1416    /*
1417     *  Now give all processes the chance to die and collect exit statuses.
1418     */
1419    if (foundOne) do_sleep(1);                    init 1        foundOne = 0
1420    for(ch = family; ch; ch = ch->next)          family                              L1437
1421        if (ch->flags & KILLME) {
1422            if (!(ch->flags & ZOMBIE))
1423                initlog(L_CO, "Pid %d [id %s] seems to hang", ch->pid,
1424                        ch->id);
1425            else {
1426                INITDBG(L_VB, "Updating utmp for pid %d [id %s]",
1427                        ch->pid, ch->id);
1428                ch->flags &= ~RUNNING;
1429                if (ch->process[0] != '+')
1430                  write_utmp_wtmp("", ch->id, ch->pid, DEAD_PROCESS, NULL);
1431            }
1432        }
1433
1434    /*
1435     *  Both rounds done; clean up the list.
1436     */
1437    sigemptyset(&nmask);
1438    sigaddset(&nmask, SIGCHLD);
1439    sigprocmask(SIG_BLOCK, &nmask, &omask);
1440    for(ch = family; ch; ch = old) {              init 1        family                L1444
```

```
1441         old = ch->next;
1442         free(ch);
1443     }
1444     family = newFamily;                    newFamily                            /etc/inittab
                                                family

1445     for(ch = family; ch; ch = ch->next) ch->new = NULL;
1446     newFamily = NULL;
1447     sigprocmask(SIG_SETMASK, &omask, NULL);
1448
1449 #ifdef INITLVL
1450     /*
1451      *  Dispose of INITLVL file.
1452      */    /etc/initrunlvl          symbol link
1453     if (lstat(INITLVL, &st) >= 0 && S_ISLNK(st.st_mode)) {        /etc/initrunlvl        symbol link
1454         /*
1455          *    INITLVL is a symbolic link, so just truncate the file.
1456          */
1457         close(open(INITLVL, O_WRONLY|O_TRUNC));
1458     } else {
1459         /*
1460          *    Delete INITLVL file.
1461          */
1462         unlink(INITLVL);
1463     }
1464 #endif
1465 #ifdef INITLVL2
1466     /*
1467      *  Dispose of INITLVL2 file.
1468      */    /var/log/initrunlvl          symbol link
1469     if (lstat(INITLVL2, &st) >= 0 && S_ISLNK(st.st_mode)) {
1470         /*
1471          *    INITLVL2 is a symbolic link, so just truncate the file.
1472          */
1473         close(open(INITLVL2, O_WRONLY|O_TRUNC));
1474     } else {
1475         /*
```

```
1476          *     Delete INITLVL2 file.
1477          */
1478         unlink(INITLVL2);
1479     }
1480  #endif
1481  }
```

                        init 1      read_inittab()       L1313                                          daemon
init 3                  init 3                                  inittab                         CHILD node
family                  init process                                      start_if_needed()

```
1483  /*
1484   *     Walk through the family list and start up children.
1485   *     The entries that do not belong here at all are removed
1486   *     from the list.
1487   */
```

                              family                                      CHILD      process[128]

```
1488  void start_if_needed(void)
1489  {
1490         CHILD *ch;          /* Pointer to child */
1491         int delete;         /* Delete this entry from list? */
1492
1493         INITDBG(L_VB, "Checking for children to start");
```

         read_inittab()      /etc/inittab                              family          node

```
1494
1495         for(ch = family; ch; ch = ch->next) {
1496
1497  #if DEBUG
1498              if (ch->rlevel[0] == 'C') {
1499                   INITDBG(L_VB, "%s: flags %d", ch->process, ch->flags);
```

```
1500              }
1501  #endif
1502
1503              /* Are we waiting for this process? Then quit here. */
1504              if (ch->flags & WAITING) break;      WAITING        init process
1505                                                                          process


1506              /* Already running? OK, don't touch it */
1507              if (ch->flags & RUNNING) continue;    node         process     running
1508
1509              /* See if we have to start it up */
1510              delete = 1;
1511              if (strchr(ch->rlevel, runlevel) ||      runlevel              run level,        node
                                                                  node       DEMAND(        )
                                                          " #*Ss"       sysinit  boot
                                                          action  (" #*Ss"     run level        )

1512                  ((ch->flags & DEMAND) && !strchr("#*Ss", runlevel))) {
1513                      startup(ch);
1514                      delete = 0;                                  node     clear
1515              }
1516
1517              if (delete) {                              node              level
1518                      /* FIXME: is this OK? */
1519                      ch->flags &= ~(RUNNING|WAITING);
1520                      if (!ISPOWER(ch->action) && ch->action != KBREQUEST)
1521                          ch->flags &= ~XECUTED;
1522                      ch->pid = 0;
1523              } else
1524                      /* Do we have to wait for this process? */
1525                      if (ch->flags & WAITING) break;
1526          }
1527      /* Done. */
1528  }
```

```
1063  /*
1064   *    Start a child running!
1065   */


         CHILD                                 l2:2:wait:/etc/rc.d/rc 2                         " /etc/rc.d/rc 2"
      " /etc/rc.d/rc"                    " 2"

1066  void startup(CHILD *ch)
1067  {
1068          /*
1069           *    See if it's disabled
1070           */
1071          if (ch->flags & FAILING) return;
1072
1073          switch(ch->action) {
1074
1075                  case SYSINIT:
1076                  case BOOTWAIT:
1077                  case WAIT:
1078                  case POWERWAIT:
1079                  case POWERFAILNOW:
1080                  case POWEROKWAIT:
1081                  case CTRLALTDEL:
1082                          if (!(ch->flags & XECUTED)) ch->flags |= WAITING;
                                                              action                    XECUTED
                                                              WAITING              break,          L1091
                                                    spawn()                              WAITING



1083                  case KBREQUEST:
```

```
1084              case BOOT:
1085              case POWERFAIL:
1086              case ONCE:
1087                  if (ch->flags & XECUTED) break;        XECUTED
1088              case ONDEMAND:
1089              case RESPAWN:
1090                  ch->flags |= RUNNING;
1091                  if (spawn(ch, &(ch->pid)) < 0) break;              ch->pid        pid spawn()
1092                  /*
1093                   *    Do NOT log if process field starts with '+'
1094                   *    FIXME: that's for compatibility with *very*
1095                   *    old getties - probably it can be taken out.
1096                   */
1097                  if (ch->process[0] != '+')
1098                      write_utmp_wtmp("", ch->id, ch->pid,
1099                          INIT_PROCESS, "");
1100                  break;
1101          }
1102  }
```

spawn

launch    CHILD              process[128]

```
823   /*
824    *    Fork and execute.
825    *
826    *    This function is too long and indents too deep.
827    *
828    */        ch                    *res              pid
```

    /etc/inittab
si::sysinit:/etc/rc.d/rc.sysinit


  CHILD *ch
ch->id = " si"

```
ch->rlevel = 0123456S
ch->action = SYSINIT
ch->process[128] = /etc/rc.d/rc.sysinit

829   int spawn(CHILD *ch, int *res)
830   {
831     char *args[16];        /* Argv array */
832     char buf[136];         /* Line buffer */
833     int f, st, rc;         /* Scratch variables */
834     char *ptr;             /* Ditto */
835     time_t t;              /* System time */
836     int oldAlarm;                /* Previous alarm value */
837     char *proc = ch->process;    /* Command line */
838     pid_t pid, pgrp;             /* child, console process group. */
839     sigset_t nmask, omask;       /* For blocking SIGCHLD */
840     struct sigaction sa;
841
842     *res = -1;
843     buf[sizeof(buf) - 1] = 0;
844
845     /* Skip '+' if it's there */
846     if (proc[0] == '+') proc++;
847
848     ch->flags |= XECUTED;                       ch
849
```

init process      action  "RESPAWN"  "ONDEMAND"
                                          init process                    2          10
          5

```
850     if (ch->action == RESPAWN || ch->action == ONDEMAND) {
851         /* Is the date stamp from less than 2 minutes ago? */
852         time(&t);                    ch->tm      process   (spawn)
853         if (ch->tm + TESTTIME > t) {  TESTTIME   2                      2
```

```
854                ch->count++;
855            } else {
856                ch->count = 0;                    2                 2
857                ch->tm = t;
858            }
859
860        /* Do we try to respawn too fast? */
861        if (ch->count >= MAXSPAWN) {                      2              10              2              10
862                                          FAILING
863          initlog(L_VB,
864              "Id \"%s\" respawning too fast: disabled for %d minutes",
865              ch->id, SLEEPTIME / 60);
866          ch->flags &= ~RUNNING;
867          ch->flags |= FAILING;              FAILING
868
869          /* Remember the time we stopped */
870          ch->tm = t;
871
872          /* Try again in 5 minutes */
873          oldAlarm = alarm(0);            alarm
874          if (oldAlarm > SLEEPTIME || oldAlarm <= 0) oldAlarm = SLEEPTIME;          disable 5
875          alarm(oldAlarm);              alarm         5
876          return(-1);                      -1
877        }
878    }
879
880    /* See if there is an "initscript" (except in single user mode). */
881    if (access(INITSCRIPT, R_OK) == 0 && runlevel != 'S') {      /etc/initscript
882        /* Build command line using "initscript" */          /etc/initscript          run level
883        args[1] = SHELL;                                Single Mode(       )
884        args[2] = INITSCRIPT;
885        args[3] = ch->id;                    L1045     execvp()
886        args[4] = ch->rlevel;
887        args[5] = "unknown";                  execvp(" /bin/sh" ,args + 1),  args[1]=" /bin/sh"
888        for(f = 0; actions[f].name; f++) {                      args[2] = " /etc/initscript"
```

```
889                 if (ch->action == actions[f].act) {                          args[3] = " si"
890                         args[5] = actions[f].name;                           args[4] = " 0123456S"
891                         break;                                             args[5]="  sysinit"
892                 }                                                          args[6]=" /etc/rc.d/rc.sysinit"
893         }                                                                  args[7]= NULL
894         args[6] = proc;
895         args[7] = NULL;
896     } else if (strpbrk(proc, "~`!$^&*()=|\\{}[];\"'<>?")) {                    " ~`!$^&*()=|\\{}[];\"'<>?"
897     /* See if we need to fire off a shell for this command */
898         /* Give command line to shell */
899         args[1] = SHELL;                        /bin/sh –c exec /etc/rc.d/rc.sysinit
900         args[2] = "-c";
901         strcpy(buf, "exec ");
902         strncat(buf, proc, sizeof(buf) - strlen(buf) - 1);
903         args[3] = buf;
904         args[4] = NULL;
905     } else {
906         /* Split up command line arguments */
907         buf[0] = 0;
908         strncat(buf, proc, sizeof(buf) - 1);
909         ptr = buf;
910         for(f = 1; f < 15; f++) {
911                 /* Skip white space */
912                 while(*ptr == ' ' || *ptr == '\t') ptr++;
913                 args[f] = ptr;
914
915                 /* May be trailing space.. */
916                 if (*ptr == 0) break;
917
918                 /* Skip this `word' */
919                 while(*ptr && *ptr != ' ' && *ptr != '\t' && *ptr != '#')
920                         ptr++;
921
922                 /* If end-of-line, break */
923                 if (*ptr == '#' || *ptr == 0) {
924                         f++;
```

```
925                     *ptr = 0;
926                     break;
927                 }
928                 /* End word with \0 and continue */
929                 *ptr++ = 0;
930         }
931         args[f] = NULL;
932     }
933     args[0] = args[1];
934     while(1) {
935         /*
936          *    Block sigchild while forking.
937          */
938         sigemptyset(&nmask);
939         sigaddset(&nmask, SIGCHLD);
940         sigprocmask(SIG_BLOCK, &nmask, &omask);

            daemon

941
942         if ((pid = fork()) == 0) {     init process
943
944             close(0);               file handle 0,1,2   STDIN  STDOUT  STDERR
945             close(1);
946             close(2);
947             if (pipe_fd >= 0) close(pipe_fd);
948
949             sigprocmask(SIG_SETMASK, &omask, NULL);
950
951             /*
952              *    In sysinit, boot, bootwait or single user mode:
953              *    for any wait-type subprocess we _force_ the console
954              *    to be its controlling tty.
955              */
956             if (strchr("*#sS", runlevel) && ch->flags & WAITING) {
957                 /*
958                  *    We fork once extra. This is so that we can
```

```
959                         *     wait and change the process group and session
960                         *     of the console after exit of the leader.
961                         */
962                     setsid();
963                     if ((f = console_open(O_RDWR|O_NOCTTY)) >= 0) {
964                             /* Take over controlling tty by force */
965                             (void)ioctl(f, TIOCSCTTY, 1);
966                             dup(f);
967                             dup(f);
968                     }
969                     if ((pid = fork()) < 0) {      fork                        init process
970                             initlog(L_VB, "cannot fork");                             L942   fork()
971                             exit(1);
972                     }
973                     if (pid > 0) {                                waitpid
974                             /*
975                              *     Ignore keyboard signals etc.
976                              *     Then wait for child to exit.
977                              */
978                             SETSIG(sa, SIGINT, SIG_IGN, SA_RESTART);
979                             SETSIG(sa, SIGTSTP, SIG_IGN, SA_RESTART);
980                             SETSIG(sa, SIGQUIT, SIG_IGN, SA_RESTART);
981                             SETSIG(sa, SIGCHLD, SIG_DFL, SA_RESTART);
982
983                             while ((rc = waitpid(pid, &st, 0)) != pid)
984                                     if (rc < 0 && errno == ECHILD)
985                                             break;
986
987                             /*
988                              *     Small optimization. See if stealing
989                              *     controlling tty back is needed.
990                              */
991                             pgrp = tcgetpgrp(f);
992                             if (pgrp != getpid())
993                                     exit(0);
994
995                             /*
```

```
996                          *    Steal controlling tty away. We do
997                          *    this with a temporary process.
998                          */
999                         if ((pid = fork()) < 0) {
1000                                initlog(L_VB, "cannot fork");
1001                                exit(1);
1002                        }
1003                        if (pid == 0) {
1004                                setsid();
1005                                (void)ioctl(f, TIOCSCTTY, 1);
1006                                exit(0);
1007                        }
1008                        while((rc = waitpid(pid, &st, 0)) != pid)
1009                                if (rc < 0 && errno == ECHILD)
1010                                        break;
1011                        exit(0);
1012                }

1014                /* Set ioctl settings to default ones */
1015                console_stty();

1017        } else {                        process   context
1018                setsid();
1019                if ((f = console_open(O_RDWR|O_NOCTTY)) < 0) {
1020                        initlog(L_VB, "open(%s): %s", console_dev,
1021                                strerror(errno));
1022                        f = open("/dev/null", O_RDWR);
1023                }
1024                dup(f);            L944  L945  L946        0,1,2   handle
1025                dup(f);      handle          0  1  2 handle      system console      1,2,3
                               Handle        console_open()  /dev/console       f   0    L1024
                                  process   1 handle  /dev/console L1025       process   2 handle
                               /dev/console       STDIN(      )  STDOUT           STDERR
```

```
1026                }
1027
1028                /* Reset all the signals, set up environment */
1029                for(f = 1; f < NSIG; f++) SETSIG(sa, f, SIG_DFL, SA_RESTART);
                                                                    process   signal

1030                environ = init_buildenv(1);
1031
1032                /*
1033                 *    Execute prog. In case of ENOEXEC try again
1034                 *    as a shell script.
1035                 */
1036                execvp(args[1], args + 1);
1037                if (errno == ENOEXEC) {
1038                        args[1] = SHELL;
1039                        args[2] = "-c";
1040                        strcpy(buf, "exec ");
1041                        strncat(buf, proc, sizeof(buf) - strlen(buf) - 1);
1042                        args[3] = buf;
1043                        args[4] = NULL;
1044                        execvp(args[1], args + 1);
1045                }
1046                initlog(L_VB, "cannot execute \"%s\"", args[1]);
1047                exit(1);
1048        }
1049        *res = pid;                              pid       init process
1050        sigprocmask(SIG_SETMASK, &omask, NULL);
1051
1052        INITDBG(L_VB, "Started id %s (pid %d)", ch->id, pid);
1053
1054        if (pid == -1) {
1055                initlog(L_VB, "cannot fork, retry..");
1056                do_sleep(5);
```

```
1057            continue;
1058        }
1059        return(pid);
1060    }
1061  }
```

check_init_fifo

```
1991  /*
1992   *     Read from the init FIFO. Processes like telnetd and rlogind can
1993   *     ask us to create login processes on their behalf.
1994   *
1995   *     FIXME:      this needs to be finished. NOT that it is buggy, but we need
1996   *             to add the telnetd/rlogind stuff so people can start using it.
1997   *             Maybe move to using an AF_UNIX socket so we can use
1998   *             the 2.2 kernel credential stuff to see who we're talking to.
1999   *
2000   */
2001  void check_init_fifo(void)
2002  {
2003    struct init_request   request;
2004    struct timeval  tv;
2005    struct stat           st, st2;
2006    fd_set          fds;
2007    int             n;
2008    int             quit = 0;
2009
2010    /*
2011     *  First, try to create /dev/initctl if not present.
2012     */
2013    if (stat(INIT_FIFO, &st2) < 0 && errno == ENOENT)                        (name pipe)/dev/initctl
2014        (void)mkfifo(INIT_FIFO, 0600);                                       root
2015
2016    /*
2017     *  If /dev/initctl is open, stat the file to see if it
2018     *  is still the _same_ inode.
```

```
2019      */
2020     if (pipe_fd >= 0) {              pipe_fd      /dev/initctl   file handle      >= 0        open
2021         fstat(pipe_fd, &st);
2022         if (stat(INIT_FIFO, &st2) < 0 ||
2023             st.st_dev != st2.st_dev ||                                      ino
2024             st.st_ino != st2.st_ino) {
2025             close(pipe_fd);                                                        open   pipe
2026             pipe_fd = -1;                                      L2033   if
2027         }
2028     }
2029
2030     /*
2031      *  Now finally try to open /dev/initctl
2032      */
2033     if (pipe_fd < 0) {                                          pipe_fd  /dev/initctl   file handle -1
2034         if ((pipe_fd = open(INIT_FIFO, O_RDWR|O_NONBLOCK)) >= 0) {
2035             fstat(pipe_fd, &st);
2036             if (!S_ISFIFO(st.st_mode)) {
2037                 initlog(L_VB, "%s is not a fifo", INIT_FIFO);
2038                 close(pipe_fd);
2039                 pipe_fd = -1;
2040             }
2041         }
2042         if (pipe_fd >= 0) {
2043             /*
2044              *   Don't use fd's 0, 1 or 2.
2045              */
2046             (void) dup2(pipe_fd, PIPE_FD);       /dev/initctl   file handle    PIPE_FD 10
2047             close(pipe_fd);
2048             pipe_fd = PIPE_FD;                      /dev/initctl    file handle   10(PIPE_FD)
2049
2050             /*
2051              *   Return to caller - we'll be back later.
2052              */
2053         }
2054     }
```

2055

/dev/initctl

```
2056    /* Wait for data to appear, _if_ the pipe was opened. */
2057    if (pipe_fd >= 0) while(!quit) {
2058
2059        /* Do select, return on EINTR. */
2060        FD_ZERO(&fds);
2061        FD_SET(pipe_fd, &fds);
2062        tv.tv_sec = 5;                          select              5  (timeout  5)
2063        tv.tv_usec = 0;
2064        n = select(pipe_fd + 1, &fds, NULL, NULL, &tv);      select      /dev/initctl
2065        if (n <= 0) {
2066            if (n == 0 || errno == EINTR) return;        init 3    init 2      request
2067            continue;
2068        }
2069
```

select                /dev/initctl                              init X        X run level
request           select

```
2070        /* Read the data, return on EINTR. */
2071        n = read(pipe_fd, &request, sizeof(request));
```

/dev/initctl    request    request

```
struct init_request {
    int    magic;            /* Magic number                */
    int    cmd;              /* What kind of request        */
    int    runlevel;         /* Runlevel to change to       */
    int    sleeptime;        /* Time between TERM and KILL   */
    union {
        struct init_request_bsd bsd;
        char            data[368];
    } i;
```

```
    };


2072        if (n == 0) {
2073                /*
2074                 *     End of file. This can't happen under Linux (because
2075                 *     the pipe is opened O_RDWR - see select() in the
2076                 *     kernel) but you never know...
2077                 */
2078                close(pipe_fd);
2079                pipe_fd = -1;
2080                return;
2081        }
2082        if (n <= 0) {
2083                if (errno == EINTR) return;              signal      select
2084                initlog(L_VB, "error reading initrequest");
2085                continue;
2086        }
2087


                                request      /dev/initctl


2088        /*
2089         *     This is a convenient point to also try to
2090         *     find the console device or check if it changed.
2091         */
2092        console_init();
2093
2094        /*
2095         *     Process request.
2096         */
2097        if (request.magic != INIT_MAGIC || n != sizeof(request)) {           request
2098                initlog(L_VB, "got bogus initrequest");
2099                continue;
2100        }
2101        switch(request.cmd) {
```

```
2102            case INIT_CMD_RUNLVL:                          run level
2103                    sltime = request.sleeptime;
2104                    fifo_new_level(request.runlevel);               run level           inittab
                                                                        run level
2105                    quit = 1;
2106                    break;
2107            case INIT_CMD_POWERFAIL:
2108                    sltime = request.sleeptime;
2109                    do_power_fail('F');
2110                    quit = 1;
2111                    break;
2112            case INIT_CMD_POWERFAILNOW:
2113                    sltime = request.sleeptime;
2114                    do_power_fail('L');
2115                    quit = 1;
2116                    break;
2117            case INIT_CMD_POWEROK:
2118                    sltime = request.sleeptime;
2119                    do_power_fail('O');
2120                    quit = 1;
2121                    break;
```

```
2122            case INIT_CMD_SETENV:
2123                    initcmd_setenv(request.i.data, sizeof(request.i.data));
2124                    break;
2125            case INIT_CMD_CHANGECONS:
2126                    if (user_console) {
2127                            free(user_console);
```

```
2128                        user_console = NULL;
2129                    }
2130                    if (!request.i.bsd.reserved[0])
2131                        user_console = NULL;
2132                    else
2133                        user_console = strdup(request.i.bsd.reserved);
2134                console_init();
2135                quit = 1;
2136                break;
2137            default:
2138                initlog(L_VB, "got unimplemented initrequest.");
2139                break;
2140        }
2141    }
2142
2143    /*
2144     *  We come here if the pipe couldn't be opened.
2145     */
2146    if (pipe_fd < 0) pause();
2147
2148  }
```

**init 2**


init 2        root            init                    init 3(daemon    )                        init
process  SIGTERM  SIGKILL              (sltime)  init

     main()

```
2597  int main(int argc, char **argv)
2598  {
2599      char             *p;
2600      int              f;
2601      int              isinit;
2602
2603      /* Get my own name */
```

```
2604        if ((p = strrchr(argv[0], '/')) != NULL)
2605            p++;
2606        else
2607            p = argv[0];
2608        umask(022);                        argv[0] = /sbin/init    p    init
2609
2610        /* Quick check */
2611        if (geteuid() != 0) {                        root        init
2612            fprintf(stderr, "%s: must be superuser.\n", p);
2613            exit(1);
2614        }
2615
2616        /*
2617         *    Is this telinit or init ?
2618         */
2619        isinit = (getpid() == 1);          init process(init 2)  pid        1 1    init 2
                                               isinit = 0
2620        for (f = 1; f < argc; f++) {  init 2       init 1                  -i  -init

2621            if (!strcmp(argv[f], "-i") || !strcmp(argv[f], "--init"))
2622                isinit = 1;
2623                break;
2624        }
2625        if (!isinit) exit(telinit(p, argc, argv));        isinit = 0    init 2        telinit()
2626
2627        /*
2628         *    Check for re-exec
2629         */
```

telinit()    init 2                    programe    init process            " init"

```
2502  int telinit(char *progname, int argc, char **argv)
```

```
2503   {
2504   #ifdef TELINIT_USES_INITLVL
2505           FILE                *fp;
2506   #endif
2507           struct init_request    request;
```

init 2   init 3         " /dev/initctl"         pipe      init 3            pipe            init 2   init 3
        init 2                                          init 3          init_request   "        "

```
struct init_request {
     int    magic;              /* Magic number                 */
     int    cmd;                /* What kind of request         */
     int    runlevel;           /* Runlevel to change to        */
     int    sleeptime;          /* Time between TERM and KILL    */
     union {
          struct init_request_bsd bsd;
          char              data[368];
     } i;
};
```

```
2508           struct sigaction  sa;
2509           int               f, fd, l;
2510           char              *env = NULL;
2511
2512           memset(&request, 0, sizeof(request));
2513           request.magic     = INIT_MAGIC;
2514
```

request

```
2515           while ((f = getopt(argc, argv, "t:e:")) != EOF) switch(f) {
2516                   case 't':          t       init    process  SIGTERM  SIGKILL           (sltime)
2517                       sltime = atoi(optarg);
2518                       break;
2519                   case 'e':
```

```
2520                    if (env == NULL)
2521                            env = request.i.data;
2522                    l = strlen(optarg);
2523                    if (env + l + 2 > request.i.data + sizeof(request.i.data)) {
2524                            fprintf(stderr, "%s: -e option data "
2525                                    "too large\n", progname);
2526                            exit(1);
2527                    }
2528                    memcpy(env, optarg, l);
2529                    env += l;
2530                    *env++ = 0;
2531                    break;
2532            default:
2533                    usage(progname);
2534                    break;
2535        }

2536
2537        if (env) *env++ = 0;

2538
2539        if (env) {
2540            if (argc != optind)
2541                    usage(progname);
2542            request.cmd = INIT_CMD_SETENV;
2543        } else {
2544            if (argc - optind != 1 || strlen(argv[optind]) != 1)
2545                    usage(progname);
2546            if (!strchr("0123456789SsQqAaBbCcUu", argv[optind][0]))
2547                    usage(progname);
2548            request.cmd = INIT_CMD_RUNLVL;                       request        check_init_fifo     (L2102)
2549            request.runlevel  = env ? 0 : argv[optind][0];        request.runlevel              run level
2550            request.sleeptime = sltime;
2551        }

2552
2553        /* Open the fifo and write a command. */
2554        /* Make sure we don't hang on opening /dev/initctl */
2555        SETSIG(sa, SIGALRM, signal_handler, 0);
2556        alarm(3);                       alarm    3            SIGALRM signal
```

```
2557         if ((fd = open(INIT_FIFO, O_WRONLY)) >= 0 &&                " /dev/initctl"
2558             write(fd, &request, sizeof(request)) == sizeof(request)) {          request
2559             close(fd);
2560             alarm(0);                                    alarm      line 2556     alarm(3)
                                         Line 2556   Line 2560                 /dev/initctl
2561             return 0;                        3            3                     SIGALRM signal        Line 2585
2562         }                                ISMEMBER          true,

2564 #ifdef TELINIT_USES_INITLVL
2565         if (request.cmd == INIT_CMD_RUNLVL) {
2566             /* Fallthrough to the old method. */
2567
2568             /* Now write the new runlevel. */
2569             if ((fp = fopen(INITLVL, "w")) == NULL) {         /etc/initrunlvl
2570                 fprintf(stderr, "%s: cannot create %s\n",
2571                     progname, INITLVL);
2572                 exit(1);
2573             }
2574             fprintf(fp, "%s %d", argv[optind], sltime);
2575             fclose(fp);
2576
2577             /* And tell init about the pending runlevel change. */
2578             if (kill(INITPID, SIGHUP) < 0) perror(progname);         1         SIGHUP signal
2579                                                                   init 2        /etc/inittab

2580             return 0;
2581         }
2582 #endif
2583
2584         fprintf(stderr, "%s: ", progname);
2585         if (ISMEMBER(got_signals, SIGALRM)) {
2586             fprintf(stderr, "timeout opening/writing control channel %s\n",
2587                 INIT_FIFO);
2588         } else {
2589             perror(INIT_FIFO);
2590         }
```

```
2591        return 1;
2592  }
```

init 2   init 3

init 3

init 2                          init 3

telinit()

/dev/initctl pipe

select
pipe

request

check_init
_fifo

init 2       pipe
request       init 3
        pipe
request
check_init_fifo()
        request

init 3          daemon
        init 2       request

**init 3**


init 3          daemon process              init 1              init 1                          init process
    Client-Server          server process                                          init process
                                    init process      init 1    init 3




   init 1

```
2340  /*
2341   *      The main loop
2342   */
2343  int init_main()
2344  {


                                        init 1


2456
2457    while(1) {                      init process
2458
2459        /* See if we need to make the boot transitions. */
2460        boot_transitions();
2461        INITDBG(L_VB, "init_main: waiting..");
2462
2463        /* Check if there are processes to be waited on. */
2464        for(ch = family; ch; ch = ch->next)                      init process
2465         if ((ch->flags & RUNNING) && ch->action != BOOT) break;
```

```
2466
2467  #if CHANGE_WAIT
2468        /* Wait until we get hit by some signal. */
2469        while (ch != NULL && got_signals == 0) {
2470          if (ISMEMBER(got_signals, SIGHUP)) {
2471                /* See if there are processes to be waited on. */
2472                for(ch = family; ch; ch = ch->next)
2473                    if (ch->flags & WAITING) break;
2474          }
2475          if (ch != NULL) check_init_fifo();
2476        }
2477  #else /* CHANGE_WAIT */
2478        if (ch != NULL && got_signals == 0) check_init_fifo();       check_init_fifo        " init 1"
2479  #endif /* CHANGE_WAIT */                                                    init 3          init 2
2480                                                                        request
2481        /* Check the 'failing' flags */
2482        fail_check();
2483
2484        /* Process any signals. */
2485        process_signals();          init                        signal    signal handler
2486                                                              " init 1"
2487        /* See what we need to start up (again) */
2488        start_if_needed();                    " init 1"
2489    }
2490    /*NOTREACHED*/
2491  }
```

init process                    signal
▪       select              **/dev/initctl**                              init process      request  init 3
        request                     check_init_fifo
▪     signal     init process    init   signal handler              signal              got_signals
        signal                        init 3        process_signals()

    init process       (signal)
```
  SETSIG(sa, SIGALRM,  signal_handler, 0);
```

```
   SETSIG(sa, SIGHUP,   signal_handler, 0);
   SETSIG(sa, SIGINT,   signal_handler, 0);
   SETSIG(sa, SIGPWR,   signal_handler, 0);
   SETSIG(sa, SIGWINCH, signal_handler, 0);
   SETSIG(sa, SIGUSR1,  signal_handler, 0);
```

init 1    HUP PWR WINCH ALRM INT signal              signal                      signal
     process_signals()

```
/*
 *    We got a signal (HUP PWR WINCH ALRM INT)
 */
void signal_handler(int sig)
{
      ADDSET(got_signals, sig);
}
```

ADDSET      macro

```
#define ADDSET(set, val)   ((set) |=  (1 << (val)))
```


   process_signals()

```
   if (ISMEMBER(got_signals, SIGPWR)) {         ISMEMBER      macro
         INITDBG(L_VB, "got SIGPWR");                  #define ISMEMBER(set, val) ((set) & (1 << (val)))
         /* See _what_ kind of SIGPWR this is. */         signal
         pwrstat = 0;
         if ((fd = open(PWRSTAT, O_RDONLY)) >= 0) {
               c = 0;
               read(fd, &c, 1);                          fail
               pwrstat = c;
               close(fd);
               unlink(PWRSTAT);
         }
         do_power_fail(pwrstat);
         DELSET(got_signals, SIGPWR);
   }

   if (ISMEMBER(got_signals, SIGINT)) {                Ctrl+Alt+Del        SIGINT signal   init
         INITDBG(L_VB, "got SIGINT");             process
```

```
        /* Tell ctrlaltdel entry to start up */
        for(ch = family; ch; ch = ch->next)
              if (ch->action == CTRLALTDEL)
                    ch->flags &= ~XECUTED;              disable
        DELSET(got_signals, SIGINT);
}

if (ISMEMBER(got_signals, SIGWINCH)) {
      INITDBG(L_VB, "got SIGWINCH");
      /* Tell kbrequest entry to start up */
      for(ch = family; ch; ch = ch->next)
            if (ch->action == KBREQUEST)
                  ch->flags &= ~XECUTED;
      DELSET(got_signals, SIGWINCH);
}

if (ISMEMBER(got_signals, SIGALRM)) {
      INITDBG(L_VB, "got SIGALRM");
      /* The timer went off: check it out */
      DELSET(got_signals, SIGALRM);
}

if (ISMEMBER(got_signals, SIGCHLD)) {
      INITDBG(L_VB, "got SIGCHLD");
      /* First set flag to 0 */
      DELSET(got_signals, SIGCHLD);

      /* See which child this was */
      for(ch = family; ch; ch = ch->next)
          if (ch->flags & ZOMBIE) {
              INITDBG(L_VB, "Child died, PID= %d", ch->pid);
              ch->flags &= ~(RUNNING|ZOMBIE|WAITING);
              if (ch->process[0] != '+')
                    write_utmp_wtmp("", ch->id, ch->pid, DEAD_PROCESS, NULL);
          }

}
```

```
   if (ISMEMBER(got_signals, SIGHUP)) {
        INITDBG(L_VB, "got SIGHUP");
#if CHANGE_WAIT
        /* Are we waiting for a child? */
        for(ch = family; ch; ch = ch->next)
            if (ch->flags & WAITING) break;
        if (ch == NULL)
#endif
        {
            /* We need to go into a new runlevel */
            oldlevel = runlevel;
#ifdef INITLVL
            runlevel = read_level(0);
#endif
            if (runlevel == 'U') {
                runlevel = oldlevel;
                re_exec();
            } else {
                if (oldlevel != 'S' && runlevel == 'S') console_stty();
                if (runlevel == '6' || runlevel == '0' ||
                    runlevel == '1') console_stty();
                read_inittab();
                fail_cancel();
                setproctitle("init [%c]", runlevel);
                DELSET(got_signals, SIGHUP);
            }
        }
    }
   if (ISMEMBER(got_signals, SIGUSR1)) {
        /*
         *    SIGUSR1 means close and reopen /dev/initctl
         */
        INITDBG(L_VB, "got SIGUSR1");
        close(pipe_fd);
        pipe_fd = -1;
        DELSET(got_signals, SIGUSR1);
```

```
        }
```

■　　　UPS　　　　fail signal

```
        if ((fd = open(PWRSTAT, O_RDONLY)) >= 0) {
                c = 0;
                read(fd, &c, 1);
                pwrstat = c;
                close(fd);
                unlink(PWRSTAT);
        }
        do_power_fail(pwrstat);
```

　　　/etc/powerstatus　　　　　"F"　"L"　"O"　　　　　　　Fail　　　Low　　　Ok(　　)

```
    pwrstat   powerfail                    /etc/inittab
    pwrstat            power
"O"        OK
"L"        Low
"F"         Fail

1757  /*
1758   *    Start up powerfail entries.
1759   */
1760  void do_power_fail(int pwrstat)
1761  {
1762        CHILD *ch;
1763
1764        /*
1765         *    Tell powerwait & powerfail entries to start up
1766         */
1767        for (ch = family; ch; ch = ch->next) {        family
1768              if (pwrstat == 'O') {                              Ok    signal
1769                    /*
1770                     *    The power is OK again.
1771                     */
```

```
1772                    if (ch->action == POWEROKWAIT)        XECUTED   disable   process
1773                        ch->flags &= ~XECUTED;                 disable        POWEROKWAIT   process

1774            } else if (pwrstat == 'L') {                          Low   signal
1775                    /*
1776                     *    Low battery, shut down now.
1777                     */
1778                    if (ch->action == POWERFAILNOW)
1779                        ch->flags &= ~XECUTED;                POWERFAILNOW   process
1780            } else {
1781                    /*
1782                     *    Power is failing, shutdown imminent              Fail   signal
1783                     */
1784                    if (ch->action == POWERFAIL || ch->action == POWERWAIT)
1785                        ch->flags &= ~XECUTED;                POWERFAIL   POWERWAIT   process
1786            }
1787        }
1788    }
```

```
# When our UPS tells us power has failed, assume we have a few minutes
# of power left.  Schedule a shutdown for 2 minutes from now.
# This does, of course, assume you have powerd installed and your
# UPS connected and working correctly.
pf::powerfail:/sbin/shutdown -f -h +2 "Power Failure; System Shutting Down"

# If power was restored before the shutdown kicked in, cancel it.
pr:12345:powerokwait:/sbin/shutdown -c "Power Restored; Shutdown Cancelled"
```

           /sbin/shutdown -f -h +2 "Power Failure; System Shutting Down"        POWERFAIL
/sbin/shutdown -c "Power Restored; Shutdown Cancelled"(     POWEROKWAIT)   init process
disable          XECUTED                              enable
       start_if_needed()

```
2484        /* Process any signals. */
2485        process_signals();                      signal
```

```
2486
2487        /* See what we need to start up (again) */
2488        start_if_needed();              family        node
2489    }                              (action)
```

■              Delete    Ctrl-C              process group   SIGINT signal

```
2263        /* Tell ctrlaltdel entry to start up */
2264        for(ch = family; ch; ch = ch->next)
2265            if (ch->action == CTRLALTDEL)
2266                ch->flags &= ~XECUTED;     Ctrl-Alt-Del handler
```

■      UNIX

          ioctl

```
2272        /* Tell kbrequest entry to start up */
2273        for(ch = family; ch; ch = ch->next)
2274            if (ch->action == KBREQUEST)
2275                ch->flags &= ~XECUTED;          KBREQUEST    process
```

■      timeout   alarm

```
2280        INITDBG(L_VB, "got SIGALRM");
2281        /* The timer went off: check it out */
2282        DELSET(got_signals, SIGALRM);
```

■                      utmp   wtmp

```
2290        /* See which child this was */
2291        for(ch = family; ch; ch = ch->next)           family
2292            if (ch->flags & ZOMBIE) {
2293              INITDBG(L_VB, "Child died, PID= %d", ch->pid);
2294              ch->flags &= ~(RUNNING|ZOMBIE|WAITING);
2295              if (ch->process[0] != '+')
2296                    write_utmp_wtmp("", ch->id, ch->pid, DEAD_PROCESS, NULL);
2297            }
```

■　　　　　　　/etc/inittab　　　　　　　　　　　　　　　　　init process
　　　　　SIGHUP signal　　init process　　inittab

```
2303  #if CHANGE_WAIT
2304      /* Are we waiting for a child? */
2305      for(ch = family; ch; ch = ch->next)                      process
2306          if (ch->flags & WAITING) break;                  inittab
2307      if (ch == NULL)                  process              ch      NULL      NULL
2308  #endif                                  inittab
2309      {
2310          /* We need to go into a new runlevel */
2311          oldlevel = runlevel;                  run level    oldlevel
2312  #ifdef INITLVL
2313          runlevel = read_level(0);          inittab      run level
2314  #endif
2315          if (runlevel == 'U') {          init          " U or u tell  init to re-execute itself
2316              runlevel = oldlevel;    (preserving the state). No re-examining of /etc/inittab
2317              re_exec();              file happens"      re_exec()              init
                                          Process
2318          } else {
2319              if (oldlevel != 'S' && runlevel == 'S') console_stty();
2320              if (runlevel == '6' || runlevel == '0' ||
2321                  runlevel == '1') console_stty();
2322              read_inittab();                      init 1      read_inittab()
2323              fail_cancel();                  init 3
2324              setproctitle("init [%c]", runlevel);
```

■　　init process　　SIGUSR1 signal(　　　　　)　　　　/dev/initctl pipe

```
    close(pipe_fd);          init          signal    init              /dev/initctl
    pipe_fd = -1;
```

```
1707  /*
1708   *     This procedure is called after every signal (SIGHUP, SIGALRM..)
1709   *
1710   *     Only clear the 'failing' flag if the process is sleeping
1711   *     longer than 5 minutes, or inittab was read again due
1712   *     to user interaction.
1713   */
1714  void fail_check(void)
1715  {
1716          CHILD *ch;                              /* Pointer to child structure */
1717          time_t      t;                          /* System time */
1718          time_t      next_alarm = 0;    /* When to set next alarm */
1719
1720          time(&t);
1721
1722          for(ch = family; ch; ch = ch->next) {           init process
1723
1724                  if (ch->flags & FAILING) {                      family          fail   node
1725                          /* Can we free this sucker? */
1726                          if (ch->tm + SLEEPTIME < t) {
1727                                  ch->flags &= ~FAILING;
1728                                  ch->count = 0;
1729                                  ch->tm = 0;
1730                          } else {
1731                                  /* No, we'll look again later */
1732                                  if (next_alarm == 0 ||
1733                                      ch->tm + SLEEPTIME > next_alarm)
1734                                          next_alarm = ch->tm + SLEEPTIME;
1735                          }
1736                  }
1737          }
```

```
1738        if (next_alarm) {
1739            next_alarm -= t;
1740            if (next_alarm < 1) next_alarm = 1;
1741            alarm(next_alarm);
1742        }
1743    }
```

re_exec

```
1827    /*
1828     *    Attempt to re-exec.
1829     */
1830    void re_exec(void)
1831    {
1832        CHILD       *ch;
1833        sigset_t    mask, oldset;
1834        pid_t       pid;
1835        char        **env;
1836        int         fd;
1837
1838        if (strchr("S12345",runlevel) == NULL)          telinit          " Run level should be one of
1839            return;                                  Ss12345, otherwise request would be silently ignored"
1840
1841        /*
1842         *    Reset the alarm, and block all signals.
1843         */
1844        alarm(0);                                  alarm signal
1845        sigfillset(&mask);
1846        sigprocmask(SIG_BLOCK, &mask, &oldset);
1847
1848        /*
1849         *    construct a pipe fd --> STATE_PIPE and write a signature
1850         */
1851        fd = make_pipe(STATE_PIPE);
1852
```

```
1853          /*
1854           * It's a backup day today, so I'm pissed off.  Being a BOFH, however,
1855           * does have it's advantages...
1856           */
1857          fail_cancel();
1858          close(pipe_fd);
1859          pipe_fd = -1;
1860          DELSET(got_signals, SIGCHLD);
1861          DELSET(got_signals, SIGHUP);
1862          DELSET(got_signals, SIGUSR1);
1863
1864          /*
1865           *    That should be cleaned.
1866           */
1867          for(ch = family; ch; ch = ch->next)           init process
1868              if (ch->flags & ZOMBIE) {
1869                  INITDBG(L_VB, "Child died, PID= %d", ch->pid);
1870                  ch->flags &= ~(RUNNING|ZOMBIE|WAITING);
1871                  if (ch->process[0] != '+')
1872                      write_utmp_wtmp("", ch->id, ch->pid, DEAD_PROCESS, NULL);
1873              }
1874
1875          if ((pid = fork()) == 0) {
1876              /*
1877               *    Child sends state information to the parent.
1878               */
1879              send_state(fd);               init        state pipe        init process     process
1880              exit(0);
1881          }
1882
1883          /*
1884           *    The existing init process execs a new init binary.
1885           */
1886          env = init_buildenv(0);
1887          execl(myname, myname, "--init", NULL, env);           init process   /sbin/init
1888                                                    init_buildenv()
```

```
1889          /*
1890           *    We shouldn't be here, something failed.
1891           *    Bitch, close the state pipe, unblock signals and return.
1892           */
1893          close(fd);
1894          close(STATE_PIPE);
1895          sigprocmask(SIG_SETMASK, &oldset, NULL);
1896          init_freeenv(env);
1897          initlog(L_CO, "Attempt to re-exec failed");
1898  }
```

read_inittab      init 3

```
1108  void read_inittab(void)
1109  {
1110    FILE            *fp;              /* The INITTAB file */
1111    CHILD           *ch, *old, *i;   /* Pointers to CHILD structure */
1112    CHILD           *head = NULL;    /* Head of linked list */
1113  #ifdef INITLVL
1114    struct stat     st;              /* To stat INITLVL */
1115  #endif
1116    sigset_t  nmask, omask;          /* For blocking SIGCHLD. */
1117    char            buf[256];        /* Line buffer */
1118    char            err[64];         /* Error message. */
1119    char            *id, *rlevel,
1120            *action, *process;       /* Fields of a line */
1121    char            *p;
1122    int     lineNo = 0;       /* Line number in INITTAB file */
1123    int     actionNo;         /* Decoded action field */
1124    int     f;                /* Counter */
1125    int     round;            /* round 0 for SIGTERM, 1 for SIGKILL */
1126    int     foundOne = 0;     /* No killing no sleep */
1127    int     talk;             /* Talk to the user */
1128    int     done = 0;         /* Ready yet? */
1129
1130  #if DEBUG
1131    if (newFamily != NULL) {
```

```
1132            INITDBG(L_VB, "PANIC newFamily != NULL");
1133            exit(1);
1134         }
1135      INITDBG(L_VB, "Reading inittab");
1136  #endif
1137
1138      /*
1139       *  Open INITTAB and real line by line.
1140       */
1141      if ((fp = fopen(INITTAB, "r")) == NULL)                    /etc/inittab
1142            initlog(L_VB, "No inittab file found");
1143
1144      while(!done) {                                    inittab        newFamily        newFamily
1145         /*                                family
1146          *    Add single user shell entry at the end.
1147          */
1148         if (fp == NULL || fgets(buf, sizeof(buf), fp) == NULL) {
1149               done = 1;
1150               /*
1151                *    See if we have a single user entry.
1152                */
1153               for(old = newFamily; old; old = old->next)
1154                     if (strpbrk(old->rlevel, "S")) break;
1155               if (old == NULL)
1156                     snprintf(buf, sizeof(buf), "~~:S:wait:%s\n", SULOGIN);
1157               else
1158                     continue;
1159         }
1160         lineNo++;
1161         /*
1162          *    Skip comments and empty lines
1163          */
1164         for(p = buf; *p == ' ' || *p == '\t'; p++)
1165                ;
1166         if (*p == '#' || *p == '\n') continue;             " #"
1167
```

```
1168          /*
1169           *    Decode the fields
1170           */
          id:runlevels:action:process    4
1171          id =      strsep(&p, ":");                           " :"                    strsep
1172          rlevel =  strsep(&p, ":");
1173          action =  strsep(&p, ":");
1174          process = strsep(&p, "\n");
1175

                         init manual                                           127
1176          /*
1177           *    Check if syntax is OK. Be very verbose here, to
1178           *    avoid newbie postings on comp.os.linux.setup :)
1179           */
1180          err[0] = 0;
1181          if (!id || !*id) strcpy(err, "missing id field");
1182          if (!rlevel)     strcpy(err, "missing runlevel field");
1183          if (!process)    strcpy(err, "missing process field");
1184          if (!action || !*action)
1185                  strcpy(err, "missing action field");
1186          if (id && strlen(id) > sizeof(utproto.ut_id))
1187              sprintf(err, "id field too long (max %d characters)",
1188                      (int)sizeof(utproto.ut_id));
1189          if (rlevel && strlen(rlevel) > 11)
1190              strcpy(err, "rlevel field too long (max 11 characters)");
1191          if (process && strlen(process) > 127)
1192              strcpy(err, "process field too long");
1193          if (action && strlen(action) > 32)
1194              strcpy(err, "action field too long");
1195          if (err[0] != 0) {
1196              initlog(L_VB, "%s[%d]: %s", INITTAB, lineNo, err);
1197              INITDBG(L_VB, "%s:%s:%s:%s", id, rlevel, action, process);
1198              continue;
1199          }
1200
1201          /*
```

```
1202              *      Decode the "action" field
1203              */
      init        action          actions[]                                               identifier
1204         actionNo = -1;
1205         for(f = 0; actions[f].name; f++)
1206                 if (strcasecmp(action, actions[f].name) == 0) {
1207                         actionNo = actions[f].act;
1208                         break;
1209                 }
1210         if (actionNo == -1) {                               action(    actions[]      )
1211                 initlog(L_VB, "%s[%d]: %s: unknown action field",
1212                         INITTAB, lineNo, action);
1213                 continue;
1214         }
1215
1216         /*
1217          *      See if the id field is unique
1218          */


                     identifier                                                          CHILD
                                                          id                             /etc/inittab
      id

1219         for(old = newFamily; old; old = old->next) {
1220                 if(strcmp(old->id, id) == 0 && strcmp(id, "~~")) {
1221                         initlog(L_VB, "%s[%d]: duplicate ID field \"%s\"",
1222                                 INITTAB, lineNo, id);
1223                         break;
1224                 }
1225         }
1226         if (old) continue;
1227
1228         /*
1229          *      Allocate a CHILD structure
1230          */
1231         ch = imalloc(sizeof(CHILD));                         CHILD node
```

```
1232
1233          /*
1234           *    And fill it in.                              CHILD node
1235           */
1236          ch->action = actionNo;          action
1237          strncpy(ch->id, id, sizeof(utproto.ut_id) + 1); /* Hack for different libs. */
1238          strncpy(ch->process, process, sizeof(ch->process) - 1);
1239          if (rlevel[0]) {                      run level
1240              for(f = 0; f < sizeof(rlevel) - 1 && rlevel[f]; f++) {
1241                  ch->rlevel[f] = rlevel[f];
1242                  if (ch->rlevel[f] == 's') ch->rlevel[f] = 'S';
1243              }
1244              strncpy(ch->rlevel, rlevel, sizeof(ch->rlevel) - 1);
1245          } else {              run level          run level           process
1246              strcpy(ch->rlevel, "0123456789");
1247              if (ISPOWER(ch->action))
1248                  strcpy(ch->rlevel, "S0123456789");
1249          }
              action
1250          /*
1251           *    We have the fake runlevel '#' for SYSINIT  and
1252           *    '*' for BOOT and BOOTWAIT.
1253           */

                  SYSINIT action  '#'      BOOT action  '*'                  run level  0  9   'S'
      " #"   " *"          run level          SYSINIT                            BOOT    action

1254          if (ch->action == SYSINIT) strcpy(ch->rlevel, "#");
1255          if (ch->action == BOOT || ch->action == BOOTWAIT)
1256              strcpy(ch->rlevel, "*");
1257
1258          /*
1259           *    Now add it to the linked list. Special for powerfail.
1260           */
```

```
     /etc/inittab                        newFamily                      family
   init       run level      family                         init            /etc/inittab

1261          if (ISPOWER(ch->action)) {         action           POWERWAIT POWERFAIL POWEROKWAIT
1262                                           POWERFAILNOW CTRLALTDEL                    Ctrl+Alt+Del

1263             /*
1264              *    Disable by default
1265              */
1266             ch->flags |= XECUTED;                        startup()                      action   flag
1267                                             XECUTED                          ISPOWER()   action
1268             /*
                                                   Ctrl+Alt+Del                    /etc/inittab  CTRLALTDEL
                                         Action       process                      disable (       XECUTED
                                          )                Ctrl+Alt+Del          enable
                                             action       family


1269              *    Tricky: insert at the front of the list..
1270              */
1271             old = NULL;
1272             for(i = newFamily; i; i = i->next) {
1273                     if (!ISPOWER(i->action)) break;
1274                     old = i;
1275             }
1276             /*
1277              *    Now add after entry "old"
1278              */
1279             if (old) {
1280                     ch->next = i;
1281                     old->next = ch;
1282                     if (i == NULL) head = ch;
1283             } else {
1284                     ch->next = newFamily;
1285                     newFamily = ch;
```

```
1286                    if (ch->next == NULL) head = ch;
1287              }
1288        } else {        action         KBREQUEST                    init manual     SIGWINCH  signal
1289            /*
1290             *      Just add at end of the list
1291             */
1292            if (ch->action == KBREQUEST) ch->flags |= XECUTED;
1293            ch->next = NULL;
1294            if (head)
1295                    head->next = ch;
1296            else
1297                    newFamily = ch;
1298            head = ch;
1299        }

1301        /*
1302         *      Walk through the old list comparing id fields
1303         */
1304        for(old = family; old; old = old->next)
1305            if (strcmp(old->id, ch->id) == 0) {
1306                    old->new = ch;
1307                    break;
1308            }

1309    }
1310    /*
1311     *  We're done.
1312     */
1313    if (fp) fclose(fp);       /etc/inittab          init 1
1314                                        daemon       init   init 3
                          kernel     init(init 1)     daemon          init  init 3


1315    /*
1316     *  Loop through the list of children, and see if they need to
1317     *  be killed.
```

```
1318      */
1319
1320     INITDBG(L_VB, "Checking for children to kill");
1321     for(round = 0; round < 2; round++) {              round 0 for SIGTERM, 1 for SIGKILL
1322       talk = 1;                                  process  SIGTERM signal      SIGKILL signal
1323       for(ch = family; ch; ch = ch->next) {                      init  init 1          family
1324         ch->flags &= ~KILLME;                                round = 0  talk = 1  foundOne = 0
1325                                                     L1393
```

init 3    family
L1321   L1414        init 3           init 1
                   init      process(  family    )          kill       inittab       process

1           process
      BOOT
      S --- Single User Mode
      DEMAND
2       process                      SIGTERM signal      sltime    sltime      5
    SIGKILL signal

```
1326          /*
1327           *    Is this line deleted?
1328           */
1329          if (ch->new == NULL) ch->flags |= KILLME;
1330
1331          /*
1332           *    If the entry has changed, kill it anyway. Note that
1333           *    we do not check ch->process, only the "action" field.
1334           *    This way, you can turn an entry "off" immediately, but
1335           *    changes in the command line will only become effective
1336           *    after the running version has exited.
1337           */
1338          if (ch->new && ch->action != ch->new->action) ch->flags |= KILLME;
1339
1340          /*
```

```
1341              *    Only BOOT processes may live in all levels
1342              */
1343          if (ch->action != BOOT &&
1344              strchr(ch->rlevel, runlevel) == NULL) {
1345                  /*
1346                   *    Ondemand procedures live always,
1347                   *    except in single user
1348                   */
1349                  if (runlevel == 'S' || !(ch->flags & DEMAND))
1350                          ch->flags |= KILLME;
1351          }
1352
1353          /*
1354           *    Now, if this process may live note so in the new list
1355           */
1356          if ((ch->flags & KILLME) == 0) {
1357                  ch->new->flags  = ch->flags;
1358                  ch->new->pid    = ch->pid;
1359                  ch->new->exstat = ch->exstat;
1360                  continue;
1361          }
1362
1363
1364          /*
1365           *    Is this process still around?
1366           */
1367          if ((ch->flags & RUNNING) == 0) {
1368                  ch->flags &= ~KILLME;
1369                  continue;
1370          }
1371          INITDBG(L_VB, "Killing \"%s\"", ch->process);
1372          switch(round) {                                    SIGTERM signal
1373                  case 0: /* Send TERM signal */
1374                          if (talk)
1375                                  initlog(L_CO,
1376                                          "Sending processes the TERM signal");
1377                          kill(-(ch->pid), SIGTERM);        SIGTERM signal        id  ch->pid
```

```
1378                         foundOne = 1;                                        process
1379                         break;
1380                 case 1: /* Send KILL signal and collect status */
1381                         if (talk)                                    SIGKILL signal
1382                                 initlog(L_CO,
1383                                         "Sending processes the KILL signal");
1384                         kill(-(ch->pid), SIGKILL);
1385                         break;
1386             }
1387         talk = 0;
1388
1389         }
1390         /*
1391          *      See if we have to wait 5 seconds
1392          */
1393         if (foundOne && round == 0) {                        init  init 1   round = 0,  foundOne = 0
1394             /*                                          if            L1419
1395              *    Yup, but check every second if we still have children.
1396              */                                        init 3              L1397   L1411
1397           for(f = 0; f < sltime; f++) {
1398                 for(ch = family; ch; ch = ch->next) {
1399                         if (!(ch->flags & KILLME)) continue;
1400                         if ((ch->flags & RUNNING) && !(ch->flags & ZOMBIE))
1401                                 break;
1402                 }
1403                 if (ch == NULL) {
1404                         /*
1405                          *    No running children, skip SIGKILL
1406                          */
1407                         round = 1;
1408                         foundOne = 0; /* Skip the sleep below. */
1409                         break;
1410                 }
1411                 do_sleep(1);                            1 * sltime
1412           }
1413         }
```

```
1414        }
1415
1416        /*
1417         *  Now give all processes the chance to die and collect exit statuses.
1418         */
1419        if (foundOne) do_sleep(1);                   init 1        foundOne = 0
1420        for(ch = family; ch; ch = ch->next)             family                              L1437
1421            if (ch->flags & KILLME) {                init 3          L1420    L1430
1422                if (!(ch->flags & ZOMBIE))
1423                    initlog(L_CO, "Pid %d [id %s] seems to hang", ch->pid,
1424                            ch->id);
1425                else {
1426                    INITDBG(L_VB, "Updating utmp for pid %d [id %s]",
1427                            ch->pid, ch->id);
1428                    ch->flags &= ~RUNNING;
1429                    if (ch->process[0] != '+')
1430                        write_utmp_wtmp("", ch->id, ch->pid, DEAD_PROCESS, NULL);
1431                }
1432            }
1433
1434        /*
1435         *  Both rounds done; clean up the list.
1436         */
1437        sigemptyset(&nmask);
1438        sigaddset(&nmask, SIGCHLD);
1439        sigprocmask(SIG_BLOCK, &nmask, &omask);
1440        for(ch = family; ch; ch = old) {                init 1        family                                L1444
1441            old = ch->next;             init 3                      family                    inittab
1442                                                    newFamily
1443            free(ch);
1444        }
1444        family = newFamily;                        newFamily                              /etc/inittab
                                                                family

1445        for(ch = family; ch; ch = ch->next) ch->new = NULL;
1446        newFamily = NULL;
```

```
1447      sigprocmask(SIG_SETMASK, &omask, NULL);
1448
1449  #ifdef INITLVL
1450      /*
1451       *  Dispose of INITLVL file.
1452       */      /etc/initrunlvl          symbol link
1453      if (lstat(INITLVL, &st) >= 0 && S_ISLNK(st.st_mode)) {        /etc/initrunlvl        symbol link
1454          /*
1455           *    INITLVL is a symbolic link, so just truncate the file.
1456           */
1457          close(open(INITLVL, O_WRONLY|O_TRUNC));
1458      } else {
1459          /*
1460           *    Delete INITLVL file.
1461           */
1462          unlink(INITLVL);
1463      }
1464  #endif
1465  #ifdef INITLVL2
1466      /*
1467       *  Dispose of INITLVL2 file.
1468       */      /var/log/initrunlvl          symbol link
1469      if (lstat(INITLVL2, &st) >= 0 && S_ISLNK(st.st_mode)) {
1470          /*
1471           *    INITLVL2 is a symbolic link, so just truncate the file.
1472           */
1473          close(open(INITLVL2, O_WRONLY|O_TRUNC));
1474      } else {
1475          /*
1476           *    Delete INITLVL2 file.
1477           */
1478          unlink(INITLVL2);
1479      }
1480  #endif
1481  }
```

Linux                                                    init
console              GUI          X Window          [1] Linux                      init process

---

[1]          Linux                           start_kernel

*Walter Zhou*

[mailto:z-l-dragon@hotmail.com](mailto:z-l-dragon@hotmail.com)

1. sysvinit-2.86
2. VMware + Redhat 8.0

## inittab  action

- **respawn**     process                                    init                              inittab
                init
- **wait**     process                                    inittab
- **once**     process
                                init
- **boot**                    init

- **bootwait**                                                    init

- **powerfail**    init                SIGPWR
- **powerwait**    init                SIGPWR
- **off**                    init        SIGTERM                        SIGKILL                    5

- **ondemand**        respawn                                rstate        a  b  c
- **sysinit**                                                        init
                init

- **initdefault** init rstate
  rstate init 0123456 6
  inittab initdefault

init process inittab shutdown( )

```
# Trap CTRL-ALT-DELETE
ca::ctrlaltdel:/sbin/shutdown -t3 -r now                          CTRL-ALT-DELETE

# When our UPS tells us power has failed, assume we have a few minutes
# of power left.  Schedule a shutdown for 2 minutes from now.
# This does, of course, assume you have powerd installed and your
# UPS connected and working correctly.                         UPS        fail
pf::powerfail:/sbin/shutdown -f -h +2 "Power Failure; System Shutting Down"

# If power was restored before the shutdown kicked in, cancel it.       UPS
pr:12345:powerokwait:/sbin/shutdown -c "Power Restored; Shutdown Cancelled"
```

init process shutdown CTRL-ALT-DELETE " -t3 -r
now" (now) 3 (t3) (r)

init process shutdown init process shutdown
init process

Shutdown

```
/*
 *    Show usage message.
 */
void usage(void)
```

```
{
	fprintf(stderr,
	"Usage:\t  shutdown [-akrhHPfnc] [-t secs] time [warning message]\n"
	"\t\t  -a:       use /etc/shutdown.allow\n"
	"\t\t  -k:       don't really shutdown, only warn.\n"
	"\t\t  -r:       reboot after shutdown.\n"
	"\t\t  -h:       halt after shutdown.\n"
	"\t\t  -P:       halt action is to turn off power.\n"
	"\t\t  -H:       halt action is to just halt.\n"
	"\t\t  -f:       do a 'fast' reboot (skip fsck).\n"
	"\t\t  -F:       Force fsck on reboot.\n"
	"\t\t  -n:       do not go through \"init\" but go down real fast.\n"
	"\t\t  -c:       cancel a running shutdown.\n"
	"\t\t  -t secs: delay between warning and kill signal.\n"
	"\t\t  ** the \"time\" argument is mandatory! (try \"now\") **\n");
	exit(1);
}
```

-a              /etc/shutdown.allow                    shutdown
-k
-r
-h        halt
-P
-f                               fsck(          )
-F                                fsck(          )
-c              shutdown
-t secs:

     shutdown

| /etc/nologin | |
|---|---|

```
/*
 *    Create the /etc/nologin file.
 */
void donologin(int min)
```

```
{
     FILE *fp;
     time_t t;

     time(&t);
     t += 60 * min;

     if ((fp = fopen(NOLOGIN, "w")) != NULL) {          create /etc/nologin
          fprintf(fp, "\rThe system is going down on %s\r\n", ctime(&t));
          if (message[0]) fputs(message, fp);
          fclose(fp);
     }
}
```

| /fastboot | fsck |
| --- | --- |

```
     while((c = getopt(argc, argv, "HPacqkrhnfFyt:g:i:")) != EOF) {          shutdown
          switch(c) {
               case 'H':
                    halttype = "HALT";
                    break;
               case 'P':
                    halttype = "POWERDOWN";
                    break;
               case 'a': /* Access control. */
                    useacl = 1;
                    break;
               case 'c': /* Cancel an already running shutdown. */
                    cancel = 1;
                    break;
               case 'k': /* Don't really shutdown, only warn.*/
                    dontshut = 1;
                    break;
               case 'r': /* Automatic reboot */
                    down_level[0] = '6';
                    break;
```

```
                    case 'h': /* Halt after shutdown */
                            down_level[0] = '0';
                            break;
                    case 'f': /* Don't perform fsck after next boot */
                            fastboot = 1;
                            break;
```

```
        chdir("/");
        if (fastboot)  close(open(FASTBOOT,  O_CREAT | O_RDWR, 0644));        fastboot         /fastboot
```

| /forcefsck | fsck |
|---|---|

```
        if (forcefsck) close(open(FORCEFSCK, O_CREAT | O_RDWR, 0644));        fastfsck         /forcefsck
```

| /etc/shutdown.allow | shutdown                                                                shutdown |
|---|---|
|                     | root                shutdown |

```
        /* Process the options. */
        while((c = getopt(argc, argv, "HPacqkrhnfFyt:g:i:")) != EOF) {
                switch(c) {
                        case 'H':
                                halttype = "HALT";
                                break;
                        case 'P':
                                halttype = "POWERDOWN";
                                break;
                        case 'a': /* Access control. */     -a option        /etc/shutdown.allow
                                useacl = 1;
                                break;
```

```c
        /* Do we need to use the shutdown.allow file ? */
        if (useacl && (fp = fopen(SDALLOW, "r")) != NULL) {              /etc/shutdown.allow

                /* Read /etc/shutdown.allow. */
                i = 0;
                while(fgets(buf, 128, fp)) {
                        if (buf[0] == '#' || buf[0] == '\n') continue;    #
                        if (i > 31) continue;                            32
                        for(sp = buf; *sp; sp++) if (*sp == '\n') *sp = 0;
                        downusers[i++] = strdup(buf);                    downusers
                }
                if (i < 32) downusers[i] = 0;
                fclose(fp);

                /* Now walk through /var/run/utmp to find logged in users. */
                while(!user_ok && (ut = getutent()) != NULL) {

                        /* See if this is a user process on a VC. */
                        if (ut->ut_type != USER_PROCESS) continue;
                        sprintf(term, "/dev/%.*s", UT_LINESIZE, ut->ut_line);
                        if (stat(term, &st) < 0) continue;
#ifdef major /* glibc */
                        if (major(st.st_rdev) != 4 ||
                            minor(st.st_rdev) > 63) continue;
#else
                        if ((st.st_rdev & 0xFFC0) != 0x0400) continue;
#endif
                        /* Root is always OK. */
                        if (strcmp(ut->ut_user, "root") == 0) {    root
                            user_ok++;
                            break;
                        }

                        /* See if this is an allowed user. */
                        for(i = 0; i < 32 && downusers[i]; i++)
                                if (!strncmp(downusers[i], ut->ut_user,      shutdown              downusers
```

```
                        UT_NAMESIZE)) {
                            user_ok++;                              shutdown
                            break;
                    }
            }
            endutent();

            /* See if user was allowed. */
            if (!user_ok) {
                if ((fp = fopen(CONSOLE, "w")) != NULL) {
                    fprintf(fp, "\rshutdown: no authorized users "
                                "logged in.\r\n");
                    fclose(fp);
                }
                exit(1);
            }
    }
```

**Shutdown**

```
/*
 * shutdown.c      Shut the system down.
 *
 * Usage:    shutdown [-krhfnc] time [warning message]
 *              -k: don't really shutdown, only warn.
 *              -r: reboot after shutdown.
 *              -h: halt after shutdown.
 *              -f: do a 'fast' reboot (skip fsck).
 *              -F: Force fsck on reboot.
 *              -n: do not go through init but do it ourselves.
 *              -c: cancel an already running shutdown.
 *              -t secs: delay between SIGTERM and SIGKILL for init.
 *
 * Author:   Miquel van Smoorenburg, miquels@cistron.nl
```

```
 *
 * Version: @(#)shutdown  2.86-1  31-Jul-2004  miquels@cistron.nl
 *
 *          This file is part of the sysvinit suite,
 *          Copyright 1991-2004 Miquel van Smoorenburg.
 *
 *          This program is free software; you can redistribute it and/or
 *          modify it under the terms of the GNU General Public License
 *          as published by the Free Software Foundation; either version
 *          2 of the License, or (at your option) any later version.
 */
#include <sys/types.h>
#include <sys/stat.h>
#include <sys/wait.h>
#include <time.h>
#include <string.h>
#include <unistd.h>
#include <errno.h>
#include <stdlib.h>
#include <stdio.h>
#include <signal.h>
#include <fcntl.h>
#include <stdarg.h>
#include <utmp.h>
#include <syslog.h>
#include "paths.h"
#include "reboot.h"
#include "initreq.h"

char *Version = "@(#) shutdown 2.86-1 31-Jul-2004 miquels@cistron.nl";

#define MESSAGELEN      256

int dontshut = 0; /* Don't shutdown, only warn  */
char down_level[2];     /* What runlevel to go to.    */
int dosync = 1;         /* Sync before reboot or halt */
int fastboot = 0; /* Do a 'fast' reboot         */
```

```c
int forcefsck = 0;          /* Force fsck on reboot       */
char message[MESSAGELEN];    /* Warning message       */
char *sltime = 0; /* Sleep time                 */
char newstate[64];          /* What are we gonna do       */
int doself = 0;             /* Don't use init       */
int got_alrm = 0;

char *clean_env[] = {
       "HOME=/",
       "PATH=/bin:/usr/bin:/sbin:/usr/sbin",
       "TERM=dumb",
       NULL,
};

/* From "wall.c" */
extern void wall(char *, int, int);

/* From "utmp.c" */
extern void write_wtmp(char *user, char *id, int pid, int type, char *line);

/*
 *     Sleep without being interrupted.
 */
void hardsleep(int secs)
{
       struct timespec ts, rem;

       ts.tv_sec = secs;
       ts.tv_nsec = 0;

       while(nanosleep(&ts, &rem) < 0 && errno == EINTR)
              ts = rem;
}

/*
 *     Break off an already running shutdown.
 */
```

```c
void stopit(int sig)
{
	unlink(NOLOGIN);
	unlink(FASTBOOT);
	unlink(FORCEFSCK);
	unlink(SDPID);
	printf("\r\nShutdown cancelled.\r\n");
	exit(0);
}


/*
 *	Show usage message.
 */
void usage(void)
{
	fprintf(stderr,
	"Usage:\t  shutdown [-akrhHPfnc] [-t secs] time [warning message]\n"
	"\t\t  -a:       use /etc/shutdown.allow\n"
	"\t\t  -k:       don't really shutdown, only warn.\n"
	"\t\t  -r:       reboot after shutdown.\n"
	"\t\t  -h:       halt after shutdown.\n"
	"\t\t  -P:       halt action is to turn off power.\n"
	"\t\t  -H:       halt action is to just halt.\n"
	"\t\t  -f:       do a 'fast' reboot (skip fsck).\n"
	"\t\t  -F:       Force fsck on reboot.\n"
	"\t\t  -n:       do not go through \"init\" but go down real fast.\n"
	"\t\t  -c:       cancel a running shutdown.\n"
	"\t\t  -t secs: delay between warning and kill signal.\n"
	"\t\t  ** the \"time\" argument is mandatory! (try \"now\") **\n");
	exit(1);
}



void alrm_handler(int sig)
{
	got_alrm = sig;
}
```

```c
/*
 *    Set environment variables in the init process.
 */
int init_setenv(char *name, char *value)
{
	struct init_request	request;
	struct sigaction  sa;
	int			fd;
	int			nl, vl;

	memset(&request, 0, sizeof(request));
	request.magic = INIT_MAGIC;
	request.cmd = INIT_CMD_SETENV;
	nl = strlen(name);
	vl = value ? strlen(value) : 0;

	if (nl + vl + 3 >= sizeof(request.i.data))
		return -1;

	memcpy(request.i.data, name, nl);
	if (value) {
		request.i.data[nl] = '=';
		memcpy(request.i.data + nl + 1, value, vl);
	}

	  /*
	 *    Open the fifo and write the command.
	   *  Make sure we don't hang on opening /dev/initctl
	 */
	memset(&sa, 0, sizeof(sa));
	sa.sa_handler = alrm_handler;
	sigaction(SIGALRM, &sa, NULL);
	got_alrm = 0;
	  alarm(3);
	  if ((fd = open(INIT_FIFO, O_WRONLY)) >= 0 &&
```

```c
            write(fd, &request, sizeof(request)) == sizeof(request)) {
                    close(fd);
                    alarm(0);
                    return 0;
        }

        fprintf(stderr, "shutdown: ");
        if (got_alrm) {
                fprintf(stderr, "timeout opening/writing control channel %s\n",
                        INIT_FIFO);
        } else {
                perror(INIT_FIFO);
        }
        return -1;
}


/*
 *      Tell everyone the system is going down in 'mins' minutes.
 */
void warn(int mins)
{
        char buf[MESSAGELEN + sizeof(newstate)];
        int len;

        buf[0] = 0;
        strncat(buf, message, sizeof(buf) - 1);
        len = strlen(buf);

        if (mins == 0)
                snprintf(buf + len, sizeof(buf) - len,
                        "\rThe system is going down %s NOW!\r\n",
                        newstate);
        else
                snprintf(buf + len, sizeof(buf) - len,
                        "\rThe system is going DOWN %s in %d minute%s!\r\n",
                            newstate, mins, mins == 1 ? "" : "s");
```

```c
        wall(buf, 1, 0);
}

/*
 *      Create the /etc/nologin file.
 */
void donologin(int min)
{
        FILE *fp;
        time_t t;

        time(&t);
        t += 60 * min;

        if ((fp = fopen(NOLOGIN, "w")) != NULL) {
                fprintf(fp, "\rThe system is going down on %s\r\n", ctime(&t));
                if (message[0]) fputs(message, fp);
                fclose(fp);
        }
}

/*
 *      Spawn an external program.
 */
int spawn(int noerr, char *prog, ...)
{
        va_list       ap;
        pid_t pid, rc;
        int   i;
        char  *argv[8];

        i = 0;
        while ((pid = fork()) < 0 && i < 10) {
                perror("fork");
                sleep(5);
                i++;
        }
```

```c
        if (pid < 0) return -1;

        if (pid > 0) {
                while((rc = wait(&i)) != pid)
                        if (rc < 0 && errno == ECHILD)
                                break;
                return (rc == pid) ? WEXITSTATUS(i) : -1;
        }

        if (noerr) fclose(stderr);

        argv[0] = prog;
        va_start(ap, prog);
        for (i = 1; i < 7 && (argv[i] = va_arg(ap, char *)) != NULL; i++)
                ;
        argv[i] = NULL;
        va_end(ap);

        chdir("/");
        environ = clean_env;

        execvp(argv[0], argv);
        perror(argv[0]);
        exit(1);

        /*NOTREACHED*/
        return 0;
}

/*
 *      Kill all processes, call /etc/init.d/halt (if present)
 */
void fastdown()
{
        int do_halt = (down_level[0] == '0');
        int i;
```

```c
#if 0
    char cmd[128];
    char *script;

    /*
     *     Currently, the halt script is either init.d/halt OR rc.d/rc.0,
     *     likewise for the reboot script. Test for the presence
     *     of either.
     */
    if (do_halt) {
        if (access(HALTSCRIPT1, X_OK) == 0)
            script = HALTSCRIPT1;
        else
            script = HALTSCRIPT2;
    } else {
        if (access(REBOOTSCRIPT1, X_OK) == 0)
            script = REBOOTSCRIPT1;
        else
            script = REBOOTSCRIPT2;
    }
#endif

    /* First close all files. */
    for(i = 0; i < 3; i++)
        if (!isatty(i)) {
            close(i);
            open("/dev/null", O_RDWR);
        }
    for(i = 3; i < 20; i++) close(i);
    close(255);

    /* First idle init. */
    if (kill(1, SIGTSTP) < 0) {
        fprintf(stderr, "shutdown: can't idle init.\r\n");
        exit(1);
    }
```

```c
        /* Kill all processes. */
        fprintf(stderr, "shutdown: sending all processes the TERM signal...\r\n");
        kill(-1, SIGTERM);
        sleep(sltime ? atoi(sltime) : 3);
        fprintf(stderr, "shutdown: sending all processes the KILL signal.\r\n");
        (void) kill(-1, SIGKILL);

#if 0
        /* See if we can run /etc/init.d/halt */
        if (access(script, X_OK) == 0) {
                spawn(1, cmd, "fast", NULL);
                fprintf(stderr, "shutdown: %s returned - falling back "
                                "on default routines\r\n", script);
        }
#endif

        /* script failed or not present: do it ourself. */
        sleep(1); /* Give init the chance to collect zombies. */

        /* Record the fact that we're going down */
        write_wtmp("shutdown", "~~", 0, RUN_LVL, "~~");

        /* This is for those who have quota installed. */
        spawn(1, "accton", NULL);
        spawn(1, "quotaoff", "-a", NULL);

        sync();
        fprintf(stderr, "shutdown: turning off swap\r\n");
        spawn(0, "swapoff", "-a", NULL);
        fprintf(stderr, "shutdown: unmounting all file systems\r\n");
        spawn(0, "umount", "-a", NULL);

        /* We're done, halt or reboot now. */
        if (do_halt) {
                fprintf(stderr, "The system is halted. Press CTRL-ALT-DEL "
                                "or turn off power\r\n");
                init_reboot(BMAGIC_HALT);
```

```
            exit(0);
      }

      fprintf(stderr, "Please stand by while rebooting the system.\r\n");
      init_reboot(BMAGIC_REBOOT);
      exit(0);
}

/*
 *     Go to runlevel 0, 1 or 6.
 */
void shutdown(char *halttype)
{
      char   *args[8];
      int    argp = 0;
      int    do_halt = (down_level[0] == '0');

      /* Warn for the last time */
      warn(0);
      if (dontshut) {
            hardsleep(1);
            stopit(0);
      }
      openlog("shutdown", LOG_PID, LOG_USER);
      if (do_halt)
            syslog(LOG_NOTICE, "shutting down for system halt");
      else
            syslog(LOG_NOTICE, "shutting down for system reboot");
      closelog();

      /* See if we have to do it ourself. */
      if (doself) fastdown();

      /* Create the arguments for init. */
      args[argp++] = INIT;
      if (sltime) {
            args[argp++] = "-t";
```

```
            args[argp++] = sltime;
        }
        args[argp++] = down_level;
        args[argp]   = (char *)NULL;

        unlink(SDPID);
        unlink(NOLOGIN);

        /* Now execute init to change runlevel. */
        sync();
        init_setenv("INIT_HALT", halttype);
        execv(INIT, args);

        /* Oops - failed. */
        fprintf(stderr, "\rshutdown: cannot execute %s\r\n", INIT);
        unlink(FASTBOOT);
        unlink(FORCEFSCK);
        init_setenv("INIT_HALT", NULL);
        openlog("shutdown", LOG_PID, LOG_USER);
        syslog(LOG_NOTICE, "shutdown failed");
        closelog();
        exit(1);
}

/*
 *      returns if a warning is to be sent for wt
 */
static int needwarning(int wt)
{
        int ret;

        if (wt < 10)
                ret = 1;
        else if (wt < 60)
                ret = (wt % 15 == 0);
        else if (wt < 180)
                ret = (wt % 30 == 0);
```

```
        else
                ret = (wt % 60 == 0);

        return ret;
}

/*
 *      Main program.
 *      Process the options and do the final countdown.
 */
int main(int argc, char **argv)
{
        FILE                *fp;
        extern int          getopt();
        extern int          optind;
        struct sigaction    sa;
        struct tm           *lt;
        struct stat         st;
        struct utmp         *ut;
        time_t                      t;
        uid_t               realuid;
        char                *halttype;
        char                *downusers[32];
        char                buf[128];
        char                term[UT_LINESIZE + 6];
        char                *sp;
        char                *when = NULL;
        int                 c, i, wt;
        int                 hours, mins;
        int                 didnolog = 0;
        int                 cancel = 0;
        int                 useacl = 0;
        int                 pid = 0;
        int                 user_ok = 0;

        /* We can be installed setuid root (executable for a special group) */
        realuid = getuid();
```

```c
        setuid(geteuid());

        if (getuid() != 0) {
                fprintf(stderr, "shutdown: you must be root to do that!\n");
                exit(1);
        }
        strcpy(down_level, "1");
        halttype = NULL;

        /* Process the options. */
        while((c = getopt(argc, argv, "HPacqkrhnfFyt:g:i:")) != EOF) {
                switch(c) {
                        case 'H':
                                halttype = "HALT";
                                break;
                        case 'P':
                                halttype = "POWERDOWN";
                                break;
                        case 'a': /* Access control. */
                                useacl = 1;
                                break;
                        case 'c': /* Cancel an already running shutdown. */
                                cancel = 1;
                                break;
                        case 'k': /* Don't really shutdown, only warn.*/
                                dontshut = 1;
                                break;
                        case 'r': /* Automatic reboot */
                                down_level[0] = '6';
                                break;
                        case 'h': /* Halt after shutdown */
                                down_level[0] = '0';
                                break;
                        case 'f': /* Don't perform fsck after next boot */
                                fastboot = 1;
                                break;
                        case 'F': /* Force fsck after next boot */
```

```c
				forcefsck = 1;
				break;
			case 'n': /* Don't switch runlevels. */
				doself = 1;
				break;
			case 't': /* Delay between TERM and KILL */
				sltime = optarg;
				break;
			case 'y': /* Ignored for sysV compatibility */
				break;
			case 'g': /* sysv style to specify time. */
				when = optarg;
				break;
			case 'i': /* Level to go to. */
				if (!strchr("0156aAbBcCsS", optarg[0])) {
					fprintf(stderr,
					"shutdown: `%s': bad runlevel\n",
					optarg);
					exit(1);
				}
				down_level[0] = optarg[0];
				break;
			default:
				usage();
				break;
		}
	}

	/* Do we need to use the shutdown.allow file ? */
	if (useacl && (fp = fopen(SDALLOW, "r")) != NULL) {

		/* Read /etc/shutdown.allow. */
		i = 0;
		while(fgets(buf, 128, fp)) {
			if (buf[0] == '#' || buf[0] == '\n') continue;
			if (i > 31) continue;
			for(sp = buf; *sp; sp++) if (*sp == '\n') *sp = 0;
```

```
                downnusers[i++] = strdup(buf);
        }
        if (i < 32) downnusers[i] = 0;
        fclose(fp);

        /* Now walk through /var/run/utmp to find logged in users. */
        while(!user_ok && (ut = getutent()) != NULL) {

                /* See if this is a user process on a VC. */
                if (ut->ut_type != USER_PROCESS) continue;
                sprintf(term, "/dev/%.*s", UT_LINESIZE, ut->ut_line);
                if (stat(term, &st) < 0) continue;
#ifdef major /* glibc */
                if (major(st.st_rdev) != 4 ||
                    minor(st.st_rdev) > 63) continue;
#else
                if ((st.st_rdev & 0xFFC0) != 0x0400) continue;
#endif
                /* Root is always OK. */
                if (strcmp(ut->ut_user, "root") == 0) {
                        user_ok++;
                        break;
                }

                /* See if this is an allowed user. */
                for(i = 0; i < 32 && downnusers[i]; i++)
                        if (!strncmp(downnusers[i], ut->ut_user,
                            UT_NAMESIZE)) {
                                user_ok++;
                                break;
                        }
        }
        endutent();

        /* See if user was allowed. */
        if (!user_ok) {
                if ((fp = fopen(CONSOLE, "w")) != NULL) {
```

```
                                fprintf(fp, "\rshutdown: no authorized users "
                                                        "logged in.\r\n");
                                fclose(fp);
                        }
                        exit(1);
                }
        }

        /* Read pid of running shutdown from a file */
        if ((fp = fopen(SDPID, "r")) != NULL) {
                fscanf(fp, "%d", &pid);
                fclose(fp);
        }

        /* Read remaining words, skip time if needed. */
        message[0] = 0;
        for(c = optind + (!cancel && !when); c < argc; c++) {
                if (strlen(message) + strlen(argv[c]) + 4 > MESSAGELEN)
                        break;
                strcat(message, argv[c]);
                strcat(message, " ");
        }
        if (message[0]) strcat(message, "\r\n");

        /* See if we want to run or cancel. */
        if (cancel) {
                if (pid <= 0) {
                        fprintf(stderr, "shutdown: cannot find pid "
                                        "of running shutdown.\n");
                        exit(1);
                }
                init_setenv("INIT_HALT", NULL);
                if (kill(pid, SIGINT) < 0) {
                        fprintf(stderr, "shutdown: not running.\n");
                        exit(1);
                }
                if (message[0]) wall(message, 1, 0);
```

```c
            exit(0);
    }

    /* Check syntax. */
    if (when == NULL) {
            if (optind == argc) usage();
            when = argv[optind++];
    }

    /* See if we are already running. */
    if (pid > 0 && kill(pid, 0) == 0) {
            fprintf(stderr, "\rshutdown: already running.\r\n");
            exit(1);
    }

    /* Extra check. */
    if (doself && down_level[0] != '0' && down_level[0] != '6') {
            fprintf(stderr,
            "shutdown: can use \"-n\" for halt or reboot only.\r\n");
            exit(1);
    }

    /* Tell users what we're gonna do. */
    switch(down_level[0]) {
            case '0':
                    strcpy(newstate, "for system halt");
                    break;
            case '6':
                    strcpy(newstate, "for reboot");
                    break;
            case '1':
                    strcpy(newstate, "to maintenance mode");
                    break;
            default:
                    sprintf(newstate, "to runlevel %s", down_level);
                    break;
    }
```

```
        /* Create a new PID file. */
        unlink(SDPID);
        umask(022);
        if ((fp = fopen(SDPID, "w")) != NULL) {
                fprintf(fp, "%d\n", getpid());
                fclose(fp);
        } else if (errno != EROFS)
                fprintf(stderr, "shutdown: warning: cannot open %s\n", SDPID);

        /*
         *      Catch some common signals.
         */
        signal(SIGQUIT, SIG_IGN);
        signal(SIGCHLD, SIG_IGN);
        signal(SIGHUP,  SIG_IGN);
        signal(SIGTSTP, SIG_IGN);
        signal(SIGTTIN, SIG_IGN);
        signal(SIGTTOU, SIG_IGN);

        memset(&sa, 0, sizeof(sa));
        sa.sa_handler = stopit;
        sigaction(SIGINT, &sa, NULL);

        /* Go to the root directory */
        chdir("/");
        if (fastboot)  close(open(FASTBOOT,  O_CREAT | O_RDWR, 0644));
        if (forcefsck) close(open(FORCEFSCK, O_CREAT | O_RDWR, 0644));

        /* Alias now and take care of old '+mins' notation. */
        if (!strcmp(when, "now")) strcpy(when, "0");
        if (when[0] == '+') when++;

        /* Decode shutdown time. */
        for (sp = when; *sp; sp++) {
            if (*sp != ':' && (*sp < '0' || *sp > '9'))
                    usage();
```

```
        }
        if (strchr(when, ':') == NULL) {
                /* Time in minutes. */
                wt = atoi(when);
                if (wt == 0 && when[0] != '0') usage();
        } else {
                /* Time in hh:mm format. */
                if (sscanf(when, "%d:%2d", &hours, &mins) != 2) usage();
                if (hours > 23 || mins > 59) usage();
                time(&t);
                lt = localtime(&t);
                wt = (60*hours + mins) - (60*lt->tm_hour + lt->tm_min);
                if (wt < 0) wt += 1440;
        }
        /* Shutdown NOW if time == 0 */
        if (wt == 0) shutdown(halttype);

        /* Give warnings on regular intervals and finally shutdown. */
        if (wt < 15 && !needwarning(wt)) warn(wt);
        while(wt) {
                if (wt <= 5 && !didnolog) {
                        donologin(wt);
                        didnolog++;
                }
                if (needwarning(wt)) warn(wt);
                hardsleep(60);
                wt--;
        }
        shutdown(halttype);

        return 0; /* Never happens */
}
```