Linux Capabilities and Namespaces

Capabilities

Michael Kerrisk, man7.org © 2025

August 2025

mtk@man7.org

3.4 Setting and viewing file capabilities3.5 Exercises3.6 Capabilities-dumb and capabilities-aware applications	utlin	e Rev: #	caf166f4161b
 3.2 Process and file capabilities 3.3 Permitted and effective capabilities 3.4 Setting and viewing file capabilities 3.5 Exercises 3.6 Capabilities-dumb and capabilities-aware applications 	3 (Capabilities	3-1
 3.3 Permitted and effective capabilities 3.4 Setting and viewing file capabilities 3.5 Exercises 3.6 Capabilities-dumb and capabilities-aware applications 	3.1	Overview	3-3
 3.4 Setting and viewing file capabilities 3.5 Exercises 3.6 Capabilities-dumb and capabilities-aware applications 	3.2	Process and file capabilities	3-8
3.5 Exercises3.6 Capabilities-dumb and capabilities-aware applications	3.3	Permitted and effective capabilities	3-13
3.6 Capabilities-dumb and capabilities-aware applications	3.4	Setting and viewing file capabilities	3-16
	3.5	Exercises	3-22
	3.6	Capabilities-dumb and capabilities-aware applications	3-28
3.7 Text-form capabilities	3.7	Text-form capabilities	3-32
3.8 Exercises	3.8	Exercises	3-35

Outline

3	Capabilities	3-1
3.1	Overview	3-3
3.2	Process and file capabilities	3-8
3.3	Permitted and effective capabilities	3-13
3.4	Setting and viewing file capabilities	3-16
3.5	Exercises	3-22
3.6	Capabilities-dumb and capabilities-aware applications	3-28
3.7	Text-form capabilities	3-32
3.8	Exercises	3-35

Rationale for capabilities

- Traditional UNIX privilege model divides users into two groups:
 - Normal users, subject to privilege checking based on UID and **GIDs**
 - Effective UID 0 (superuser) bypasses many of those checks
- Coarse granularity is a problem:
 - E.g., to give a process power to change system time, we must also give it power to bypass file permission checks
 - → No limit on possible damage if program is compromised



[TLPI §39.1]

Rationale for capabilities

- Capabilities divide power of superuser into small pieces
 - 41 capabilities, as at Linux 6.16
 - Traditional superuser == process that has full set of capabilities
- Goal: replace set-UID-root programs with programs that have capabilities
 - Compromise in set-UID-*root* binary ⇒ very dangerous
 - Compromise in binary with file capabilities \Rightarrow less dangerous
- Capabilities are not specified by POSIX
 - A 1990s standardization effort was ultimately abandoned
 - Some other implementations have something similar
 - E.g., Solaris, FreeBSD



Linux Capabilities and Namespaces

©2025 M. Kerrisk

Capabilities

3-5 §3.1

A selection of Linux capabilities

Capability	Permits process to
CAP_CHOWN	Make arbitrary changes to file UIDs and GIDs
CAP_DAC_OVERRIDE	Bypass file RWX permission checks
CAP_DAC_READ_SEARCH	Bypass file R and directory X permission checks
CAP_IPC_LOCK	Lock memory
CAP_FOWNER	<pre>chmod(), utime(), set ACLs on arbitrary files</pre>
CAP_KILL	Send signals to arbitrary processes
CAP_NET_ADMIN	Various network-related operations
CAP_SETFCAP	Set file capabilities
CAP_SETGID	Make arbitrary changes to process's (own) GIDs
CAP_SETPCAP	Make changes to process's (own) capabilities
CAP_SETUID	Make arbitrary changes to process's (own) UIDs
CAP_SYS_ADMIN	Perform a wide range of system admin tasks
CAP_SYS_BOOT	Reboot the system
CAP_SYS_NICE	Change process priority and scheduling policy
CAP_SYS_MODULE	Load and unload kernel modules
CAP_SYS_RESOURCE	Raise process resource limits, override some limits
CAP_SYS_TIME	Modify the system clock

More details: capabilities(7) manual page and TLPI §39.2

man7.org

Linux Capabilities and Namespaces ©2025 M. Kerrisk Capabilities 3-6 §3.1

Supporting capabilities

- To support implementation of capabilities, the kernel must:
 - Check process capabilities for each privileged operation
 - Cf. traditional check: is process's effective UID 0?
 - Provide system calls allowing a process to modify its capabilities
 - So process can raise (add) and lower (remove) capabilities
 - (Capabilities analog of set*id() calls)
 - Support attaching capabilities to executable files
 - When file is executed, process gains attached capabilities
 - (Capabilities analog of set-UID-*root* program)
- Implemented as follows:
 - Support for first two pieces available since Linux 2.2 (1999)
 - Support for file capabilities added in Linux 2.6.24 (2008)
 - (Delay due to design concerns rather than technical reasons)



[TLPI §39.4]

Linux Capabilities and Namespaces

©2025 M. Kerrisk

Capabilities

3-7 §3.1

3 Capabilities	3-1
3.1 Overview	3-3
3.2 Process and file capabilities	3-8
3.3 Permitted and effective capabilities	3-13
3.4 Setting and viewing file capabilities	3-16
3.5 Exercises	3-22
3.6 Capabilities-dumb and capabilities-aware applications	3-28
3.7 Text-form capabilities	3-32
3.8 Exercises	3-35

Process and file capabilities

- Processes and (binary) files can each have capabilities
- Process capabilities define power of process to do privileged operations
 - Traditional superuser == process that has **all** capabilities
- File capabilities are a mechanism to give a process capabilities when it execs the file
 - Stored in security.capability extended attribute
 - (File metadata; getfattr -m <file>)



[TLPI §39.3]

Linux Capabilities and Namespaces

©2025 M. Kerrisk

Capabilities

3-9 §3.2

Process and file capability sets

- Capability set: bit mask representing a group of capabilities
- Each **process**[†] has 3[‡] capability sets:
 - Permitted
 - Effective
 - Inheritable

†In truth, capabilities are a per-thread attribute ‡In truth, there are more capability sets

- An **executable file** may have 3 associated capability sets:
 - Permitted
 - Effective
 - Inheritable
- Inheritable capabilities are little used; can mostly ignore

man7.org

Linux Capabilities and Namespaces ©2025 M. Kerrisk Capabilities 3-10 §3.2

Viewing process capabilities

• /proc/PID/status fields (hexadecimal bit masks):

```
$ cat /proc/4091/status
CapInh: 0000000000000000
CapPrm: 000000000200020
CapEff: 0000000000000000
```

- See <sys/capability.h> for capability bit numbers
 - Here: CAP KILL (bit 5), CAP SYS ADMIN (bit 21)
- getpcaps(1) (part of libcap package):

```
$ getpcaps 4091
Capabilities for `4091': = cap_kill,cap_sys_admin+p
```

- More readable notation, but a little tricky to interpret
- Here, single '=' means all sets are empty
- capsh(1) can be used to decode hex masks:

```
$ capsh --decode=200020
0x0000000000200020=cap_kill,cap_sys_admin
```

man7.org

Linux Capabilities and Namespaces

©2025 M. Kerrisk

Capabilities

3-11 §3.2

Modifying process capabilities

- A process can modify its capability sets by:
 - Raising a capability (adding it to set)
 - Synonyms: add, enable
 - Lowering a capability (removing it from set)
 - Synonyms: drop, clear, remove, disable
 - Mostly, we'll defer discussion of the APIs until later
- There are various rules about changes a process can make to its capability sets



Outline

3	Capabilities	3-1
3.1	Overview	3-3
3.2	Process and file capabilities	3-8
3.3	Permitted and effective capabilities	3-13
3.4	Setting and viewing file capabilities	3-16
3.5	Exercises	3-22
3.6	Capabilities-dumb and capabilities-aware applications	3-28
3.7	Text-form capabilities	3-32
3.8	Exercises	3-35

Process permitted and effective capabilities

- Permitted: capabilities that process may employ
 - "Upper bound" on effective capability set
 - Once dropped from permitted set, a capability can't be reacquired
 - (But see discussion of execve() later)
 - Can't drop while capability is also in effective set
- Effective: capabilities that are currently in effect for process
 - I.e., capabilities that are examined when checking if a process can perform a privileged operation
 - Capabilities can be dropped from effective set and reacquired
 - Operate with least privilege....
 - Reacquisition possible only if capability is in permitted set



[TLPI §39.3.3]

File permitted and effective capabilities

- Permitted: a set of capabilities that may be added to process's permitted set during exec()
- - If set, all capabilities in process's new permitted set are also enabled in effective set
 - Useful for so-called *capabilities-dumb* applications
 - If not set, process's new effective set is empty
- File capabilities allow implementation of capabilities analog of set-UID-root program
 - Notable difference: setting effective bit off allows a program to start in unprivileged state
 - Set-UID/set-GID programs always start in privileged state



[TLPI §39.3.4]

Linux Capabilities and Namespaces

©2025 M. Kerrisk

Capabilities

3-15 §3.3

3 Capabilities	3-1
3.1 Overview	3-3
3.2 Process and file capabilities	3-8
3.3 Permitted and effective capabilities	3-13
3.4 Setting and viewing file capabilities	3-16
3.5 Exercises	3-22
3.6 Capabilities-dumb and capabilities-aware applications	3-28
3.7 Text-form capabilities	3-32
3.8 Exercises	3-35

Setting and viewing file capabilities from the shell

- setcap(8) sets capabilities on files
 - Requires privilege (CAP_SETFCAP "set file capabilities")
 - E.g., to set CAP_SYS_TIME as a permitted and effective capability on an executable file:

```
$ cp /bin/date mydate
$ sudo setcap "cap_sys_time=pe" mydate
```

getcap(8) displays capabilities associated with a file

```
$ getcap mydate
mydate = cap_sys_time+ep
```

filecap(8) searches for files that have capabilities:

man7.org

filecap is part of the libcap-ng-utils package

[TLPI §39.3.6]

Linux Capabilities and Namespaces

©2025 M. Kerrisk

Capabilities

3-17 §3.4

cap/demo_file_caps.c

- Display process capabilities
- Report result of opening file named in argv[1] (if present)



cap/demo_file_caps.c

- All steps in demos are done from unprivileged user ID 1000
- Binary has no capabilities ⇒ process gains no capabilities
 - "=" in the output means "all capability sets empty"
- open() of /etc/shadow fails
 - Because /etc/shadow is readable only by privileged process
 - Process needs CAP_DAC_READ_SEARCH capability



man7.org

Linux Capabilities and Namespaces

©2025 M. Kerrisk

Capabilities

3-19 §3.4

cap/demo_file_caps.c

```
$ sudo setcap cap_dac_read_search=p demo_file_caps
$ ./demo_file_caps /etc/shadow
Capabilities: = cap_dac_read_search+p
Open failed: Permission denied
```

- Binary confers permitted capability to process, but capability is not effective
- Process gains capability in permitted set
- open() of /etc/shadow fails
 - Because CAP_DAC_READ_SEARCH is not in effective set



cap/demo_file_caps.c

```
$ sudo setcap cap_dac_read_search=pe demo_file_caps
$ ./demo_file_caps /etc/shadow
Capabilities: = cap_dac_read_search+ep
Successfully opened /etc/shadow
```

- Binary confers permitted capability and has effective bit on
- Process gains capability in permitted and effective sets
- open() of /etc/shadow succeeds



Linux Capabilities and Namespaces

©2025 M. Kerrisk

Capabilities

3-21 §3.4

3 Capabilities	3-1
3.1 Overview	3-3
3.2 Process and file capabilities	3-8
3.3 Permitted and effective capabilities	3-13
3.4 Setting and viewing file capabilities	3-16
3.5 Exercises	3-22
3.6 Capabilities-dumb and capabilities-aware applications	3-28
3.7 Text-form capabilities	3-32
3.8 Exercises	3-35

Notes for online practical sessions

- Small groups in breakout rooms
 - Write a note into Slack if you have a preferred group
- We will go faster, if groups collaborate on solving the exercise(s)
 - You can share a screen in your room
- I will circulate regularly between rooms to answer questions
- Zoom has an "Ask for help" button...
- Keep an eye on the #general Slack channel
 - Perhaps with further info about exercise;
 - Or a note that the exercise merges into a break
- When your room has finished, write a message in the Slack channel: "***** Room X has finished *****"
 - Then I have an idea of how many people have finished



Linux Capabilities and Namespaces

©2025 M. Kerrisk

Capabilities

3-23 §3.5

Shared screen etiquette

- It may help your colleagues if you use a larger than normal font!
 - In many environments (e.g., *xterm*, *VS Code*), we can adjust the font size with Control+Shift+"+" and Control+"-"
 - Or (e.g., emacs) hold down Control key and use mouse wheel
- Long shell prompts make reading your shell session difficult
 - Use PS1='\$ ' or PS1='# '
- Low contrast color themes are difficult to read; change this if you can
- Turn on line numbering in your editor
 - In vim use: :set number
 - In emacs use: M-x display-line-numbers-mode <RETURN>
 M-x means Left-Alt+x
- For collaborative editing, relative line-numbering is evil....
 - In vim use: :set nornu
 - In *emacs*, the following should suffice:

M-: means Left-Alt+Shift+:

M-: (display-line-numbers-mode) <RETURN>
M-: (setq display-line-numbers 'absolute) <RETURN>

man7.org

©2025 M. Kerrisk

Capabilities

3-24 §3.5

Using tmate in in-person practical sessions

In order to share an X-term session with others, do the following:

• Enter the command *tmate* in an X-term, and you'll see the following:

```
$ tmate
...
Connecting to ssh.tmate.io...
Note: clear your terminal before sharing readonly access
web session read only: ...
ssh session read only: ssh SOmErAnDOm5Tr1Ng@lon1.tmate.io
web session: ...
ssh session: ssh SOmEoTheRrAnDOm5Tr1Ng@lon1.tmate.io
```

- Share last "ssh" string with colleague(s) via Slack or another channel
 - Or: "ssh session read only" string gives others read-only access
- Your colleagues should paste that string into an X-term...
- Now, you are sharing an X-term session in which anyone can type
 - Any "mate" can cut the connection to the session with the 3-character sequence <ENTER> \sim .
- To see above message again: tmate show-messages

Linux Capabilities and Namespaces

©2025 M. Kerrisk

Capabilities

3-25 §3.5

Exercises

Ompile and run the cap/demo_file_caps program, without adding any capabilities to the file, and verify that when you run the binary, the process has no capabilities:

```
$ cc -o demo_file_caps demo_file_caps.c -lcap
$ ./demo_file_caps
```

- The string "=" means all capability sets empty.
- Now make the binary set-UID-root:

```
$ sudo chown root demo_file_caps # Change owner to root
$ sudo chmod u+s demo_file_caps # Turn on set-UID bit
$ ls -l demo_file_caps # Verify
-rwsr-xr-x. 1 root mtk 8624 Oct 1 13:19 demo_file_caps
```

- 3 Run the binary and verify that the process gains all capabilities. (The string "=ep" means "all capabilities in the permitted + effective sets".)
 - If the process does not gain all capabilities, check whether the filesystem is mounted with the nosuid option (findmnt -T <dir>). If it is, either remount the filesystem without that option or do the exercise on a filesystem that is not mounted with nosuid (typically, /tmp should work).
- Take the existing set-UID-root binary, add a permitted capability to it, and set the effective capability bit:

\$ sudo setcap cap_dac_read_search=pe demo_file_caps

man7.org

Linux Capabilities and Namespaces ©2025 M. Kerrisk

Capabilities

Exercises

- When you now run the binary, what capabilities does the process have?
 - \$./demo_file_caps
- 6 Suppose you assign empty capability sets to the binary. When you execute the binary, what capabilities does the process then have?

```
$ sudo setcap = demo_file_caps
$ getcap demo_file_caps
$ ./demo_file_caps
```

Use the following command to remove capabilities from the binary and verify that when executed, the binary once more grants all capabilities to the process:

```
$ sudo setcap -r demo_file_caps
$ getcap demo_file_caps
$ ./demo_file_caps
```

Use the following command to find the binaries on your system that have capabilities attached:

\$ sudo filecap -a 2> /dev/null

Write the name of your distribution, and paste the list of binaries into the Slack man7.org channel.

Linux Capabilities and Namespaces

©2025 M. Kerrisk

Capabilities

3-27 §3.5

3 Capabilities	3-1
3.1 Overview	3-3
3.2 Process and file capabilities	3-8
3.3 Permitted and effective capabilities	3-13
3.4 Setting and viewing file capabilities	3-16
3.5 Exercises	3-22
3.6 Capabilities-dumb and capabilities-aware applications	3-28
3.7 Text-form capabilities	3-32
3.8 Exercises	3-35

Capabilities-dumb and capabilities-aware applications

• Capabilities-dumb application:

- (Typically) an existing set-UID-root binary whose code we can't change
 - Thus, binary does not know to use capabilities APIs (Binary simply uses traditional set*uid() APIs)
- But want to make legacy binary less dangerous than set-UID-root
- Converse is **capabilities-aware** application
 - Program that was written/modified to use capabilities APIs
 - Set binary up with file effective capability bit off
 - Program "knows" it must use capabilities APIs to enable effective capabilities



Linux Capabilities and Namespaces

©2025 M. Kerrisk

Capabilities

3-29 §3.6

Adding capabilities to a capabilities-dumb application

To convert existing set-UID-root binary to use file capabilities:

- Setup:
 - Binary remains set-UID-root
 - Enable a subset of file permitted capabilities + set effective bit on
 - l.e., capabilities-dumb == binary with effective bit on
 - (Note: code of binary isn't changed)
- Operation:
 - When binary is executed, process gets (just the) specified subset of capabilities in its permitted and effective sets
 - IOW: file-capabilities override effect of set-UID-root bit, which would normally confer all capabilities to process
 - Process UID changes between zero and nonzero automatically raise/lower process's capabilities
 - (Covered in more detail later)



Linux Capabilities and Namespaces

man7.org

How do I work out what capabilities a program needs?

Some possibilities to discover what capabilities are needed by an arbitrary program:

- System call manual pages (section 2) are a good start
 - Look for capability requirements documented in DESCRIPTION or ERRORS
- Run the program (without capabilities) under *strace(1)*:
 - System call failures due to lack of capabilities normally return EPERM in *errno*
 - <u>M</u> But not all EPERM errors are due to lack of capabilities
 - If program displays an error message that seems to relate to capabilities, look in trace output for nearby EPERM errors
 - You may want to use the -v option so that strace doesn't abbreviate strings
- In extreme cases, you may need to read kernel source



Linux Capabilities and Namespaces

©2025 M. Kerrisk

Capabilities

3-31 §3.6

2 Canabilities	2 1
3 Capabilities	3-1
3.1 Overview	3-3
3.2 Process and file capabilities	3-8
3.3 Permitted and effective capabilities	3-13
3.4 Setting and viewing file capabilities	3-16
3.5 Exercises	3-22
3.6 Capabilities-dumb and capabilities-aware applications	3-28
3.7 Text-form capabilities	3-32
3.8 Exercises	3-35

Textual representation of capabilities

- Both setcap(8) and getcap(8) work with textual representations of capabilities
 - Syntax described in cap_from_text(3) manual page
- String read left to right, containing space-separated clauses
 - (The capability sets are initially considered to be empty)
- Clause: caps-list operator flags [operator flags] ...
 - caps-list is comma-separated list of capability names, or all
 - *operator* is +, -, or =
 - flags is zero or more of p (permitted), e (effective), or
 i (inheritable)
 - Clause can contain multiple [operator flags] parts:
 - E.g., "cap_sys_time+p-i" (is same as "cap_sys_time+p cap_sys_time-i")



man7.org

Linux Capabilities and Namespaces

©2025 M. Kerrisk

Capabilities

3-33 §3.7

Textual representation of capabilities

Operators:

- + operator: raise capabilities in sets specified by *flags*
- operator: lower capabilities in sets specified by flags
- = operator:
 - Raise capabilities in sets specified by flags; lower those capabilities in remaining sets
 - So, "CAP_KILL=p" is same as "CAP_KILL+p-ie"
 - caps-list can be omitted; defaults to all
 - flags can be omitted ⇒ clear capabilities from all sets
 ⇒ Thus: "=" means clear all capabilities in all sets
- What does "=p cap_kill,cap_sys_admin+e" mean?
 - All capabilities in permitted set, plus CAP_KILL and CAP_SYS_ADMIN in effective set



man7.org

Outline

3	Capabilities	3-1
3.1	Overview	3-3
3.2	Process and file capabilities	3-8
3.3	Permitted and effective capabilities	3-13
3.4	Setting and viewing file capabilities	3-16
3.5	Exercises	3-22
3.6	Capabilities-dumb and capabilities-aware applications	3-28
3.7	Text-form capabilities	3-32
3.8	Exercises	3-35

Exercises

- What capability bits are enabled by each of the following text-form capability specifications?
 - "=p"

 - "cap_setuid=p cap_sys_time+pie"
 - "=p cap_kill-p"
 - "cap_kill=p = cap_sys_admin+pe"
 - "cap_chown=i cap_kill=pe cap_setfcap,cap_chown=p"
- The program cap/cap_text.c takes a single command-line argument, which is a text-form capability string. It converts that string to an in-memory representation and then iterates through the set of all capabilities, printing out the state of each capability within the permitted, effective, and inheritable sets. It thus provides a method of verifying your interpretation of text-form capability strings. Try supplying each of the above strings as an argument to the program (remember to enclose the entire string in quotes!) and check the results against your answers to the previous exercise.



Exercises

- The *pscap* command (part of *libcap-ng*) displays a list of the processes on the system that have permitted, effective, or inheritable capabilities. In addition to showing the PPID, PID, UID, command, and capabilities for each of the displayed processes, output lines may be annotated with one of the following characters:
 - +: the process has a nonempty capability bounding set
 - 0: the process has a nonempty ambient capability set (later)
 - *: the process is in a child user namespace (later)

Use the *pscap* command to display the processes that have capabilities on your system. (By default, PID 1 (init) is excluded from the list; use the -a option to include PID 1, if you wish.)



Linux Capabilities and Namespaces

 $@2025\ M.\ Kerrisk$

Capabilities

3-37 §3.8

This page intentionally blank