Open Source Summit Europe

# The Linux capabilities model

Michael Kerrisk, man7.org © 2019

mtk@man7.org
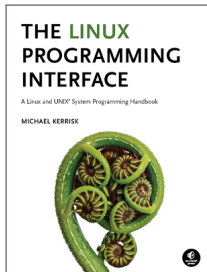
29 October 2019, Lyon, France

# Outline

# Who am I?

- Maintainer of Linux *man-pages* project since 2004
  - ≈1050 pages, mainly for system calls & C library functions
    - https://www.kernel.org/doc/man-pages/
    - (I wrote a lot of those pages...)
- Author of a book on the Linux programming interface
  - http://man7.org/tlpi/
- **Trainer**/writer/engineer
  http://man7.org/training/
- Email: mtk@man7.org
  Twitter: @mkerrisk

THE LINUX
PROGRAMMING
INTERFACE
A Linux and UNIX® System Programming Handbook

MICHAEL KERRISK

*man7.org*

# Outline

## 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

# Rationale for capabilities

- Capabilities divide power of superuser into small pieces
  - 38 capabilities, as at Linux 5.4
  - Traditional superuser == process that has full set of capabilities
- Goal: replace set-UID-*root* programs with programs that have capabilities
  - Set-UID-*root* program compromised ⇒ very dangerous
  - Compromise in binary with file capabilities ⇒ less dangerous

# A selection of Linux capabilities

| Capability | Permits process to |
|---|---|
| CAP_CHOWoN | 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_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)* man page

*man7.org*

# Outline

# Process and file capabilities

- Processes and (executable) 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

*man7.org*

# 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*

# Viewing process capabilities

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

```
$ cat /proc/4091/status
...
CapInh: 0000000000000000
CapPrm: 0000000000200020
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 inheritable + effective sets are empty

# 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
- There are various rules about changes a process can make to its capability sets
  - (APIs are *libcap* library, *capset(2)*, *capget(2)*, *prctl(2)*; we won't look at these)

*man7.org*

# Outline

# 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 *exec* 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
    - Reacquisition possible only if capability is in permitted set

*man7.org*

# File permitted and effective capabilities

- *Permitted* : a set of capabilities that may be added to process's permitted set during *exec()*
- *Effective* : a **single bit** that determines state of process's new effective set after *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

# Outline

# Setting and viewing file capabilities from the shell

- *setcap(8)* sets capabilities on files
    - Only available to privileged users (`CAP_SETFCAP`)
    - 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
    ```

    (This is the capabilities equivalent of a set-UID program)

- *getcap(8)* displays capabilities associated with a file

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

# cap/demo_file_caps.c

```c
int main(int argc, char *argv[]) {
  cap_t caps;
  int fd;
  char *str;

  caps = cap_get_proc();    /* Fetch process capabilities */
  str = cap_to_text(caps, NULL);
  printf("Capabilities: %s\n", str);
  ...
  if (argc > 1) {
    fd = open(argv[1], O_RDONLY);
    if (fd >= 0)
      printf("Successfully opened %s\n", argv[1]);
    else
      printf("Open failed: %s\n", strerror(errno));
  }
  exit(EXIT_SUCCESS);
}
```

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

## cap/demo_file_caps.c

```
$ id -u
1000
$ cc -o demo_file_caps demo_file_caps.c -lcap
$ ./demo_file_caps /etc/shadow
Capabilities: =
Open failed: Permission denied
$ ls -l /etc/shadow
----------. 1 root root 1974 Mar 15 08:09 /etc/shadow
```

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

*man7.org*

# 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

# Outline

# Transformation of process capabilities during *exec*

- During *execve()*, process's capabilities are transformed:

```
P'( perm ) = F ( perm ) & P ( bset )

P'( eff ) = F ( eff ) ? P'( perm ) : 0
```

  - *P()* / *P'()*: process capability set before/after *exec*
  - *F()*: file capability set (**of file that is being execed**)
- New permitted set for process comes from file permitted set ANDed with *capability bounding set* (discussed soon)
  - ⚠ Note that *P(perm)* has no effect on *P'(perm)*
- New effective set is either 0 or same as new permitted set
- ⚠ Transformation **rules above are a simplification**
  - (More details later)

*man7.org*

# Transformation of process capabilities during *exec*

- Commonly, process bounding set contains all capabilities
- Therefore transformation rule for process permitted set:

```
P'(perm) = F(perm) & P(bset)
```

commonly simplifies to:

```
P'(perm) = F(perm)
```

*man7.org*

# Outline

# The capability bounding set

- Per-process attribute (actually: per-thread)
- A "safety catch" to limit capabilities that can be gained during *exec*
    - Limits capabilities that can be granted by file permitted set
    - Limits capabilities that can be added to process inheritable set (later)
- Use case: remove some capabilities from bounding set to ensure process never regains them on *execve()*
    - E.g., *systemd* reduces bounding set before executing some daemons
        - Guarantees that daemon can **never** get certain capabilities

*man7.org*

# The capability bounding set

- Inherited by child of *fork()*, preserved across *execve()*
  - *init* starts with capability bounding set containing **all capabilities**
- To view: `/proc/PID/status` CapBnd field
- Can (irreversibly) drop capabilities from bounding set
  - *prctl()* `PR_CAPBSET_DROP`
  - Requires `CAP_SETPCAP` effective capability
  - Doesn't change permitted, effective, and inheritable sets

*man7.org*

# Outline

# Inheritable and ambient capabilities

- Processes[†] and files can each have a set of **inheritable** capabilities, but:
  - Inheritable capabilities turned out not to be fit for purpose
  - They are little used
  - You can pretty much ignore them
- Process[†] **ambient** capabilities were added in Linux 4.3:
  - Added to solve the problem that inheritable capabilities didn't solve

    [†]In truth, capabilities are a per-thread attribute

*man7.org*

## Ambient capabilities

- Problem scenario (not solved by inheritable capabilities):
    - We have a parent process that has capabilities
    - Parent wants to create a child process that executes an **unprivileged** helper program
    - Helper **should** have same capabilities as parent process
    - But **child loses capabilities** on *exec* because of transformation rule: $P'(perm) = F(perm) \And P(bset)$
- Ambient capabilities provide a way for child to preserve some its capabilities across *exec*:
    - Child copies some of its permitted capabilities into its ambient set
    - During *exec* of **unprivileged** binary, ambient capabilities are added to process's new permitted and effective sets

# Outline

# Capabilities and *execve()*

- During *execve()*, process capabilities transform as follows:

```
P'(amb) = (privileged-binary) ? 0 : P(amb)

P'(perm) = (P(inh) & F(inh)) | (F(perm) & P(bset))
                  | P'(amb)

P'(eff) = F(eff) ? P'(perm) : P'(amb)

P'(inh) = P(inh)

P'(bset) = P(bset)
```

- *P()* / *P'()*: process capability set before/after *exec*
- *F()*: file capability set
- *privileged-binary* == binary that is set-UID or set-GID or has file capabilities attached

man7.org

# Capabilities and *execve()* – simplified

```
P'(amb) = (privileged-binary) ? 0 : P(amb)

P'(perm) = F(perm) | P'(amb)

P'(eff) = F(eff) ? P'(perm) : P'(amb)
```

Simplification, based on:

- Inheritable capabilities are normally unused
- Process bounding set is (usually) all bits on

# Outline

# Capabilities: the promise

- Can be used to make a program more secure
  - Reduce power of program $\Rightarrow$ attacks become more difficult
- But not a panacea

*man7.org*

# Capabilities: the problems

- It's (too) complicated!
- Less familiar to sysadmins
- More work to program
    - New, more complex set of APIs for changing privilege states
- Some capabilities can be leveraged to full power of *root* in some circumstances
    - See "False Boundaries and Arbitrary Code Execution" http://forums.grsecurity.net/viewtopic.php?f=7&t=2522

*man7.org*

## Capabilities: the problems

- Some capabilities are too broad
    - Capability required to do single operation may also allow many other operations
        - Kernel developer dilemma: for new privileged operation ⇒ add new capability or re-use an existing capability?
    - Most prominent example: `CAP_SYS_ADMIN`
        - Accounts for ~~nearly 40% (Linux 3.2, 2012)~~ over 45% (Linux 5.2) of all capability checks in kernel! ☺
        - See https://lwn.net/Articles/486306/; Michael Kerrisk, "CAP_SYS_ADMIN: the new root", March 2012

# Thanks!

Michael Kerrisk, Trainer and Consultant
http://man7.org/training/

mtk@man7.org    @mkerrisk

Slides at http://man7.org/conf/
Source code at http://man7.org/tlpi/code/