

## SETCAP

Утилита для управления возможностями файлов в Linux.

### Синтаксис

```
setcap [-q] [-n <rootuid>] [-v] {capabilities|-r} filename [ ... capabilitiesN fileN ]
```

### Описание

Без опции **-v** (verify - проверка) **setcap** устанавливает указанные возможности для каждого файла, заданного параметром filename. Необязательный аргумент **-n <rootuid>** можно использовать для установки возможности, используемого лишь в пространстве имен пользователя, владеющего идентификатором root. Опция **-v** служит для проверки связывания указанных возможностей с файлом. При указании опций **-v** и **-n** проверяется также аргумент **-n <rootuid>**.

Возможности указываются с помощью **cap\_from\_text**, как описано ниже.

Функция **cap\_from\_text()** выделяет и инициализирует состояние возможности в рабочем хранилище. Функция устанавливает содержимое вновь созданного состояния возможности в соответствии с понятной человеку строкой символов в стиле языка C (nul-завершение), указанной buf\_p и возвращает указатель на созданное состояние. Текстовое представление возможности состоит из одного или множества разделенных пробелами «слов», каждое из которых указывает те или иные операции из набора возможностей, которые включаются или отключаются после выполнения команды (по умолчанию все отключено).

Каждое слово представляет собой список разделенных запятыми имен возможностей (или all), за которым следует список действий. Список действий включает последовательность пар «оператор-флаги». Операторами могут служить =, + и -, а флагами - e, i и p. Регистр флагов имеет значение и они указывают группу (набор), к которой относятся возможности Effective, Inheritable (наследуется) и Permitted (разрешено), соответственно.

В списке имен регистр символов не учитывается. Специальное имя **all** указывает все возможности, что эквивалентно заданию полного списка возможностей. Безымянные возможности можно указывать номером. Это позволяет библиотеке libcap поддерживать возможности, которые еще не были выделены к моменту компиляции библиотеки.

Оператор = задает для указанного имени сначала сброс всех 3 параметров возможности, а затем их установку в соответствии с флагами (не обязательны для этого оператора). Например, all=p сбросит установку Effective и Inheritable для батарейного питания и установит Permitted, а cap\_fowner=ep установит в Effective и Permitted возможности смены владельца файла (override-file-ownership) и сбросит их в Inheritable.

При указании первым оператора = без списка возможностей предполагается воздействие на все возможности (**all**). Например, варианты all=, =, и **cap\_chown,<every-other-capability>=** будут эквивалентны.

Операторы + и - требуют явного указания впереди списка возможностей, а после оператора - явных флагов. Оператор + включает все указанные в списке возможности в соответствии с флагами, а оператор - отключает. Например, **all+p** включает для всех возможностей Permitted, а **cap\_fowner-i** отключает override-file-ownership в группе Inheritable.

Список действий может включать множество пар «оператор-флаги», действия применяются слева направо. Например, **cap\_fowner+p-i** эквивалентно **cap\_fowner+p cap\_fowner-i**, а **cap\_fowner+pe-i** эквивалентно **cap\_fowner+pe**.

Специальная строка возможностей - позволяет указать возможности, считываемые со стандартного устройства ввода (stdin). В таких случаях набор возможностей завершается пустой строкой.

Специальная строка **-r** служит для удаления набора возможностей файла. Отметим, что установка пустого набора возможностей **отличается** от удаления набора. Пустой набор может использоваться для предотвращения исполнения файла с привилегиями, несмотря на то, что превалирующие наборы Ambient и Inheritable позволили бы это.

Флаг **-q** служит для сокращения выводимых программой сведений.

### Список возможностей

Для выполнения проверки прав доступа в традиционных системах UNIX процессы делят на две категории - привилегированные (идентификатор эффективного пользователя равен 0, как у root), и непривилегированные (ID эффективного пользователя не равен 0). Для привилегированных процессов проверки прав в ядре не выполняются, а для не привилегированных процессов выполняется полная проверка на основе возможностей процесса (обычно, эффективные UID и GID, а также и список дополнительных групп).

В ядре Linux, начиная с версии 2.2, все привилегии, обычно связываемые с суперпользователем, разделены на несколько частей, называемых возможностями (capability), которые можно разрешать и запрещать независимо. Возможности являются атрибутами потоков (нитей). Ниже приведен список возможностей, реализованных в Linux с указанием операции и поведении, разрешаемых ими. В скобках указаны версии ядра, начиная с которых каждая возможность поддерживается.

#### **CAP\_AUDIT\_CONTROL (2.6.11)**

Включение и выключение аудит ядра, изменение фильтрующих правил аудита, чтение состояния аудита и фильтрующие правила.

#### **CAP\_AUDIT\_READ (3.16)**

Чтение протокола аудита через многоадресный сокет netlink.

#### **CAP\_AUDIT\_WRITE (2.6.11)**

Запись данных в журнал аудита ядра.

**CAP\_BLOCK\_SUSPEND (3.5)**

Возможности, которые могут приводить к блокированию приостановки системы (**epoll EPOLLWAKEUP**, **/proc/sys/wake\_lock**).

**CAP\_BPF (5.8)**

Привилегированные операции BPF (**bpf**, **bpf\_helpers**). Возможность добавлена в Linux 5.8 для разгрузки перегруженной возможности CAP\_SYS\_ADMIN.

**CAP\_CHECKOUT\_RESTORE (5.9)**

- Обновление **/proc/sys/kernel/ns\_last\_pid** (**pid\_namespaces**);
- использование свойства **set\_tid** в **clone3**;
- чтение содержимого символической ссылки **/proc/[pid]/map\_files** для других процессов.

Возможность добавлена в Linux 5.9 для переноса функциональности **checkpoint/restore** из CAP\_SYS\_ADMIN.

**CAP\_CHOWN**

Произвольные изменения UID и GID для файлов.

**CAP\_DAC\_OVERRIDE**

Пропуск проверки доступа к файлу на чтение, запись и выполнение (DAC или discretionary access control - избирательный контроль доступа).

**CAP\_DAC\_READ\_SEARCH**

- Пропуск проверки доступа к файлу на чтение и к каталогу на чтение и выполнение (просмотр);
- вызов **open\_by\_handle\_at**;
- использование **linkat** с флагом **AT\_EMPTY\_PATH** для создания ссылки на файл, заданным дескриптором.

**AP\_FOWNER**

- Пропуск проверки доступа для операций, которые обычно требуют совпадения UID файловой системы процесса и UID файла (например, **chmod**, **utime**), исключая операции, охватываемые **CAP\_DAC\_OVERRIDE** и **CAP\_DAC\_READ\_SEARCH**, системы процесса и UID файла (например, **chmod**, **utime**), исключая операции, охватываемые **CAP\_DAC\_OVERRIDE** и **CAP\_DAC\_READ\_SEARCH**;
- изменение флагов inode (**ioctl\_iflags**) у произвольных файлов;
- установка списков контроля доступа (ACL) для произвольных файлов;
- игнорирование закрепляющего бита при удалении файла;
- установка **O\_NOATIME** для произвольных файлов в **open** и **fcntl**.

**CAP\_FSETID**

Позволяет не сбрасывать биты режима set-user-ID и set-group-ID при изменении файла любыми дополнительными GID вызывающего процесса.

**CAP\_IPC\_LOCK**

Блокировка памяти (**mlock**, **mlockall**, **mmap**, **shmctl**).

**CAP\_IPC\_OWNER**

Отмена проверки доступа для операций с объектами System V IPC.

**CAP\_KILL**

**ioctl** с операцией **KDSIGACCEPT**.

**CAP\_LEASE**

(2.4) Установка аренды для произвольных файлов (**fcntl**).

**CAP\_LINUX\_IMMUTABLE**

Установка флагов inode **FS\_APPEND\_FL** и **FS\_IMMUTABLE\_FL** (**ioctl\_iflags**).

**CAP\_MAC\_ADMIN**

(2.6.25) Изменение MAC или состояния. Реализована в Smack Linux Security Module (LSM).

**CAP\_MAC\_OVERRIDE**

(2.6.25) Замещение мандатного контроля доступа (MAC). Реализована в Smack LSM.

**CAP\_MKNOD**

(2.4) Создание специальных файлов с помощью **mknod**.

**CAP\_NET\_ADMIN**

Различные сетевые операции:

- настройка интерфейса;
- управление IP МСЭ, трансляцией адресов и ведением учёта;
- изменение таблиц маршрутизации;
- привязка к любому адресу для организации прозрачного прокси;
- назначение типа обслуживания (TOS);
- сброс статистики драйвера;
- управление режимом захвата (**promiscuous**);
- включение групповой рассылки (**multicasting**);
- использование **setsockopt** для включения параметров сокета **SO\_DEBUG**, **SO\_MARK**, **SO\_PRIORITY** (для приоритетов вне диапазона 0 - 6), **SO\_RCVBUFFORCE** и **SO\_SNDBUFFORCE**.

**CAP\_NET\_BIND\_SERVICE**

Привязка сокета к привилегированным портам домена Internet (номера портов меньше 1024).

**CAP\_NET\_BROADCAST**

(не используется) Позволяет выполнять широковещание с сокета и прослушивание многоадресных рассылок.

**CAP\_NET\_RAW**

Использование сокетов RAW и PACKET и привязка к любому адресу для создания прозрачного прокси.

**CAP\_PERFMON (5.9)**

Управление механизмами монитора, включая:

- вызов **perf\_event\_open**;
- реализация операций BPF, влияющих на производительность.

Возможность добавлена в Linux 5.9 для разгрузки перегруженной возможности CAP\_SYS\_ADMIN.

**CAP\_SETGID**

Произвольные действия с GID процесса и списком дополнительных GID, подмена GID при передаче возможностей сокета через доменные сокеты UNIX, запись отображения идентификатора пользователя в пользовательское пространство (**user\_namespaces**).

**CAP\_SETFCAP**

(2.6.24) Устанавливает произвольные возможности для файла.

**CAP\_SETPCAP**

Если поддерживаются возможности для файла, позволяет добавлять любую возможность из ограничивающего набора вызывающего потока в наследуемый набор, отменять возможности из ограничивающего набора (с помощью `prctl` с операцией `PR_CAPBSET_DROP`), изменять флаги `securebits`.

Если возможности для файла не поддерживаются (до 2.6.24), позволяет предоставлять и отменять любую возможность из списка разрешённых вызывающего или любого другого процесса (свойство недоступно с ядрами, поддерживающими возможности для файлов, поскольку `CAP_SETPCAP` в таких ядрах имеет другую семантику).

**CAP\_SETUID**

Произвольные действия с UID процесса (**setuid**, **setreuid**, **setresuid**, **setfsuid**), подмена UID при передаче возможностей сокета через доменные сокеты UNIX, запись отображения идентификатора пользователя в пользовательское пространство (**user\_namespaces**).

**CAP\_SYS\_ADMIN**

- Управление системой, включая **quotactl**, **mount**, **umount**, **swapon**, **swapoff**, **sethostname**, **setdomainname**;
- привилегированные операции `syslog` (начиная с Linux 2.6.37, для этих операций нужно использовать `CAP_SYSLOG`);
- команды `VM86_REQUEST_IRQ vm86`;
- функциональность `checkpoint/restore`, обеспечиваемая `CAP_CHECKPOINT_RESTORE` (предпочтительна);
- функциональность `BPF`, обеспечиваемая `CAP_BPF` (предпочтительна);
- функциональность мониторинга производительности, обеспечиваемая `CAP_PERFMON` (предпочтительна);
- операции `IPC_SET` и `IPC_RMID` над произвольными объектами System V IPC;
- перезапись ограничения ресурса `RLIMIT_NPROC`;
- операции над расширенными атрибутами `trusted` и `security (xattr)`;
- использование **lookup\_dcookie**;
- использование **ioprio\_set** для назначения классов планирования ввода-вывода `IOPRIO_CLASS_RT` и (до Linux 2.6.25) `IOPRIO_CLASS_IDLE`;
- подмена PID при передаче возможностей сокета через доменные сокеты UNIX;
- превышение системного ограничения (`/proc/sys/fs/file-max`) на количество открытых файлов в системных вызовах, открывающих файлы (например, **accept**, **execve**, **open**, **pipe**);
- использование флагов `CLONE_*`, создающих новые пространства имён с помощью **clone** и **unshare** (начиная с Linux 3.8 для создания пользовательских пространств никаких возможностей не требуется);
- вызовы **perf\_event\_open**;
- доступ к информации о привилегированном событии `perf`;
- вызовы **setns** (требуется `CAP_SYS_ADMIN` в целевом пространстве имён);
- вызовы **fanotify\_init**;
- привилегированные операции `KEYCTL_CHOWN` и `KEYCTL_SETPERM` в **keyctl**;
- вызовы **ptrace** `PTRACE_SECCOMP_GET_FILTER` для получения дампа фильтров `seccomp` трассируемого;
- операция `MADV_HWPOISON` в **madvise**;
- использование `TIOCSTI` в **ioctl** для вставки символов во входную очередь терминала, отличного от управляющего терминала вызывающего;
- устаревший системный вызов **nfsservctl**;
- устаревший системный вызов **bdflush**;
- привилегированные операции **ioctl** над блочными устройствами;
- привилегированные операции **ioctl** над файловой системой;
- привилегированные операции **ioctl** над устройством `/dev/random` (**random**);
- установка фильтров **seccomp** без начальной установки атрибута потока `no_new_privs`;
- изменение правил разрешения и запрета для групп управления устройствами;
- операция **ptrace** `PTRACE_SECCOMP_GET_FILTER` для получения дампа фильтров `seccomp` трассируемого;
- операция **ptrace** `PTRACE_SETOPTIONS` для приостановки защиты `seccomp` трассируемого (т. е., флаг `PTRACE_O_SUSPEND_SECCOMP`);
- административные операции над многими драйверами устройств;
- изменение значений приоритета `autogroup` путем записи в `/proc/[pid]/autogroup`.

**CAP\_SYS\_BOOT**

Использование **reboot** и **kexec\_load**.

**CAP\_SYS\_CHROOT**

Использование **chroot** и смена примонтированных пространств имен с помощью **setns**.

**CAP\_SYS\_MODULE**

Загрузка и выгрузка модулей ядра (**init\_module** и **delete\_module**), в ядрах до версии 2.6.25 - отзыв возможностей из системного набора `Bounded`.

**CAP\_SYS\_NICE**

- Снижение приоритета процесса (**nice**, **setpriority**) и изменение приоритета произвольных процессов;
- задание политики планирования в реальном масштабе времени для вызывающего процесса, а также политики планирования и приоритетов для произвольных процессов (**sched\_setscheduler**, **sched\_setparam**, **shed\_setattr**);
- привязка к ЦП для произвольных процессов (**sched\_setaffinity**);
- назначение класса планирования ввода-вывода и приоритета для произвольных процессов (**ioprio\_set**);
- применение **migrate\_pages** к произвольным процессам для их переноса на произвольные узлы;
- применение **move\_pages** к произвольным процессам;
- использование флага `MPOL_MF_MOVE_ALL` в **mbind** и **move\_pages**.

**CAP\_SYS\_PACCT**

Использование **acct**.

**CAP\_SYS\_PTRACE**

- Трассировка любого процесса с помощью **ptrace**;
- применение **get\_robust\_list** к произвольным процессам;
- перенос данных в память произвольного процесса (или из нее) с помощью **process\_vm\_readv** и **process\_vm\_writev**;

- изучение процессов с помощью **kcmp**.

**CAP\_SYS\_RAWIO**

- Операции ввода-вывода из портов (**iopl** и **ioperm**);
- доступ к **/proc/kcore**;
- использование операции **FIBMAP** в **ioctl(2)**;
- создание устройств для доступа к специальным регистрам x86 (**MSR**, **msr**);
- обновление **/proc/sys/vm/mmap\_min\_addr**;
- отображение памяти по адресам меньше значения, заданного в **/proc/sys/vm/mmap\_min\_addr**;
- отображение файлов в **/proc/bus/pci**;
- доступ к **/dev/mem** и **/dev/kmem**;
- выполнение различных команд устройств SCSI;
- определённые операции с устройствами **hpsa** и **cciss**;
- некоторые специальные операции с другими устройствами.

**CAP\_SYS\_RESOURCE**

- Использование зарезервированного пространства файловых систем **ext2**;
- вызовы **ioctl**, управляющие журналом **ext3**;
- превышение ограничений дисковой квоты;
- увеличение ограничений по ресурсам (**setrlimit**);
- перезапись ограничений ресурсов **RLIMIT\_NPROC**;
- превышение максимального числа консолей при выделении консоли;
- превышение максимального числа раскладок;
- использование более 64hz прерывания из часов реального времени;
- установка значения **msg\_qbytes** очереди сообщений System V больше **/proc/sys/kernel/msgmnb** (**msgop** и **msgctl**);
- разрешение **RLIMIT\_NOFILE** ограничивать число файловых дескрипторов, передаваемых в обход, при передаче другому процессу через доменный сокет UNIX (**unix**);
- переопределение **/proc/sys/fs/pipe-size-max** при назначении вместимости канала с помощью команды **F\_SETPPIPE\_SZ** в **fcntl**;
- использование **F\_SETPPIPE\_SZ** для увеличения вместимости канала сверх **/proc/sys/fs/pipe-max-size**;
- переопределение **/proc/sys/fs/mqueue/queues\_max**, **/proc/sys/fs/mqueue/msg\_max**, **/proc/sys/fs/mqueue/msg-size\_max** при создании очередей сообщений POSIX (**mq\_overview**);
- операция **prctl PR\_SET\_MM()**;
- установка в **/proc/[pid]/oom\_score\_adj** значения меньше последнего, заданного процессом с помощью **CAP\_SYS\_RESOURCE**.

**CAP\_SYS\_TIME**

Настройка системных часов (**settimeofday**, **stime**, **adjtimex**) и часов реального времени (аппаратных).

**CAP\_SYS\_TTY\_CONFIG**

Использование **vhangup(2)** и привилегированные операции **ioctl** с виртуальными терминалами.

**CAP\_SYSLOG (2.6.37)**

Привилегированные операции **syslog**, просмотр адресов ядра в **/proc** и других интерфейсах, когда значение **/proc/sys/kernel/kptr\_restrict = 1** (**kptr\_restrict** в **proc**).

**CAP\_WAKE\_ALARM (3.0)**

Вызов чего-либо при пробуждении системы (установка таймеров **CLOCK\_REALTIME\_ALARM** и **CLOCK\_BOOTTIME\_ALARM**).

**Группы (наборы) возможностей потока**

Каждый поток (нить) имеет несколько групп (наборов) возможностей из числа перечисленных выше. Группа может быть пустой.

**Permitted – разрешенные возможности**

Ограничивающее надмножество набора Effective с возможностями, которые поток может предполагать. Этот набор также ограничивает список возможностей, которые могут быть добавлены в наследуемый набор потоком без **CAP\_SETPCAP** в своём наборе Effective.

Если поток отменяет возможность в своём наборе Permitted, он не сможет вернуть ее (без вызова **execve** из программы с **set-user-ID-root** или программы, чьи возможности для файла позволяют это).

**Inheritable – наследуемые возможности**

Этот набор возможностей, сохраняемых при вызове **execve**. Наследуемые возможности остаются таковыми при выполнении любой программы и добавляются в набор Permitted, если у выполняющейся программы установлены соответствующие биты в наследуемом наборе для файлов.

Поскольку наследуемые возможности обычно не сохраняются после **execve**, если выполнение происходит не от суперпользователя, приложениям, которым нужно выполнять вспомогательные программы с расширенными возможностями, нужно применять внешние возможности (Ambient), описанные ниже.

**Bounding**

(на уровне потока, начиная с Linux 2.6.25) Механизм, который может использоваться для ограничения возможностей, предоставляемых вызовом **execve**.

Начиная с Linux 2.6.25, этот набор задается на уровне потока, а в более ранних версиях это атрибут системы для всех потоков.

**Effective**

Набор возможностей, используемый ядром при выполнении проверок прав для потока.

**Ambient - внешние (охватывающие) возможности**

(начиная с Linux 4.3) Набор возможностей, сохраняемый после **execve** для непривилегированных программ. Для набора внешних возможностей соблюдается правило, в соответствии с которым возможность не может быть внешней, если она не входит в оба набора Permitted и Inheritable.

Набор внешних возможностей можно изменять с помощью **prctl**. Этот набор автоматически сужается при сужении соответствующего набора Permitted и Inheritable.

Выполнение программы, меняющей UID или GID из-за установленного бита **set-user-ID** или **set-group-ID**, а также программы, которая может устанавливать возможности для файла, сбрасывает набор Ambient. Внешние возможности добавляются к набору Permitted и назначаются набору Effective при вызове **execve**. Если внешние

возможности расширяют набор Permitted или Effective при вызове **execve**, это не вызывает режим защищенного исполнения (*secure-execution*), описанный в **ld.so**.

Потомки (ветви) потока, созданные с помощью **fork**, наследуют копии наборов родительских возможностей.

Поток может менять свои возможности с помощью **capset**.

Начиная с версии 2.6.24 файл `/proc/sys/kernel/cap_last_cap` показывает максимальное численное значение возможностей, поддерживаемых работающим ядром. Это позволяет определить старший бит, который можно установить в наборе возможностей.

## Возможности для файлов

Начиная с ядра 2.6.24 поддерживается связывание наборов возможностей с исполняемым файлом с помощью утилиты **setcap**. Наборы возможностей хранятся в расширенном атрибуте (**setxattr**, **xattr**) с именем `security.capability`. Для записи в этот атрибут требуется возможность `CAP_SETFCAP`. Наборы возможностей файла вместе с наборами возможностей потока определяют реальные возможности потока после вызова **execve**.

Для файла поддерживается три набора возможностей.

### **Permitted** (ранее *forced*)

Возможности, автоматически предоставляемые потоку, независимо отнаследуемых им возможностей.

### **Inheritable** (ранее *allowed*)

Пересечение (AND) этого набора с набором наследуемых потоком возможностей определяет для потока набор Permitted после вызова **execve**.

### **Effective**

По сути, это не набор возможностей, а просто бит, при установке которого вызов **execve** переносит все вновь разрешенные для потока возможности в новый набор Effective. При сброшенном бите ни одна из вновь разрешенных вызовом **execve** возможностей потока не присутствует в новом наборе Effective.

Установка бита предполагает, что любая из возможностей Effective или Inheritable для файла, которая заставляет поток обрести соответствующую возможность Permitted в результате вызова **execve** (Правила преобразования), также будет включена в набор Effective. Поэтому при включении возможности для файла (**setcap**, **cap\_set\_file**, **cap\_set\_fd**) установка бита Effective для любой возможности должна сопровождаться установкой этого флага для всех других возможностей, где установлен флаг Permitted или Inheritable.

## Правила преобразования

При выполнении **execve** ядро рассчитывает новые возможности процесса, как показано ниже.

```
P'(ambient)      = (файл привилегированный) ? 0 : P(ambient)
P'(permitted)    = (P(inheritable) & F(inheritable)) | (F(permitted) & P(bounding)) | P'(ambient)
P'(effective)    = F(effective) ? P'(permitted) : P'(ambient)
P'(inheritable) = P(inheritable)      [т. е. не меняется]
P'(bounding)    = P(bounding)         [т. е. не меняется]
```

где P() указывает значение возможности потока до вызова **execve**, P'() - возможности после вызова, F() - набор возможностей для файла.

Отметим некоторые детали, относящиеся к правилам преобразования возможностей:

- набор Ambient появился лишь в Linux 4.3 и при преобразовании набора внешних возможностей в процессе **execve** привилегированным считается файл, имеющий возможности или установленный бит `set-user-ID` или `set-group-ID` bit `set`;
- до Linux 2.6.25 набор Bounding был системным атрибутом для всех потоков и это значение использовалось для расчетов при вызове **execve** как P(bounding).

Отметим, что при описанном выше преобразовании возможности файла можно игнорировать (считать пустыми) по тем же причинам, которые служат для игнорирования битов `set-user-ID` и `set-group-ID`. Возможности файлов игнорируются при загрузке ядра с опцией `no_file_caps` option.

Текст подготовлен на основании материалов `man` из Mageia 8.

Николай Малых

[nmalykh@protocols.ru](mailto:nmalykh@protocols.ru)