
Subject: [PATCH 1/3] Signal semantics for /sbin/init
Posted by [Sukadev Bhattiprolu](#) on Fri, 31 Aug 2007 20:29:49 GMT
[View Forum Message](#) <> [Reply to Message](#)

(This is Oleg's patch with my pid ns additions. Compiled and unit tested on 2.6.23-rc3-mm1 with other patches in this set. I think Oleg will send this to akpm separately - including it here for easier review of other patches in this set)

Currently, /sbin/init is protected from unhandled signals by the "current == child_reaper(current)" check in get_signal_to_deliver(). This is not enough, we have multiple problems:

- this doesn't work for multi-threaded inits, and we can't fix this by simply making this check group-wide.
- /sbin/init and kernel threads are not protected from handle_stop_signal(). Minor problem, but not good and allows to "steal" SIGCONT or change ->signal->flags.
- /sbin/init is not protected from __group_complete_signal(), sig_fatal() can set SIGNAL_GROUP_EXIT and block exec(), kill sub-threads, set ->group_stop_count, etc.

Also, with support for multiple pid namespaces, we need an ability to actually kill the sub-namespace's init from the parent namespace. In this case it is not possible (without painful and intrusive changes) to make the "should we honor this signal" decision on the receiver's side.

Hopefully this patch (adds 43 bytes to kernel/signal.o) can solve these problems.

Notes:

- Blocked signals are never ignored, so init still can receive a pending blocked signal after sigprocmask(SIG_UNBLOCK). Easy to fix, but probably we can ignore this issue.
- this patch allows us to simplify de_thread() playing games with pid_ns->child_reaper.

(Side note: the current behaviour of things like force_sig_info_fault() is not very good, init should not ignore these signals and go to the endless loop. Exit + panic is imho better, easy to change)

Oleg.

kernel/signal.c | 47 ++++++-----
1 file changed, 33 insertions(+), 14 deletions(-)

Index: 2.6.23-rc3-mm1/kernel/signal.c

=====

--- 2.6.23-rc3-mm1.orig/kernel/signal.c 2007-08-30 13:43:03.000000000 -0700

+++ 2.6.23-rc3-mm1/kernel/signal.c 2007-08-31 00:02:50.000000000 -0700

@@ -26,6 +26,7 @@

#include <linux/freezer.h>

#include <linux/pid_namespace.h>

#include <linux/nsproxy.h>

+#include <linux/hardirq.h>

#include <asm/param.h>

#include <asm/uaccess.h>

@@ -39,11 +40,35 @@

static struct kmem_cache *sigqueue_cache;

+static int sig_init_ignore(struct task_struct *tsk)
+{

-static int sig_ignored(struct task_struct *t, int sig)
+ // Currently this check is a bit racy with exec(),
+ // we can _simplify_ de_thread and close the race.
+ if (likely(!is_container_init(tsk->group_leader)))
+ return 0;
+
+ if (!in_interrupt())
+ return 0;
+
+ return 1;
+}

+
+static int sig_task_ignore(struct task_struct *tsk, int sig)
+{
- void __user * handler;
+ void __user * handler = tsk->sigband->action[sig-1].sa.sa_handler;
+
+ if (handler == SIG_IGN)
+ return 1;
+
+ if (handler != SIG_DFL)
+ return 0;
+
+ return sig_kernel_ignore(sig) || sig_init_ignore(tsk);
+}

```

+static int sig_ignored(struct task_struct *t, int sig)
+{
+/*
+ * Tracers always want to know about signals..
+ */
@@ -58,10 +83,7 @@ static int sig_ignored(struct task_struct
+ if (sigismember(&t->blocked, sig))
+ return 0;

- /* Is it explicitly or implicitly ignored? */
- handler = t->sighand->action[sig-1].sa.sa_handler;
- return handler == SIG_IGN ||
- (handler == SIG_DFL && sig_kernel_ignore(sig));
+ return sig_task_ignore(t, sig);
}

/*
@@ -568,6 +590,9 @@ static void handle_stop_signal(int sig,
+ */
+ return;

+ if (sig_init_ignore(p))
+ return;
+
+ if (sig_kernel_stop(sig)) {
+ /*
+ * This is a stop signal. Remove SIGCONT from all queues.
+ */
@@ -1863,12 +1888,6 @@ relock:
+ if (sig_kernel_ignore(signr)) /* Default is nothing. */
+ continue;

- /*
- * Global init gets no signals it doesn't want.
- */
- if (is_global_init(current))
- continue;
-
+ if (sig_kernel_stop(signr)) {
+ /*
+ * The default action is to stop all threads in
+ */
@@ -2320,6 +2339,7 @@ int do_sigaction(int sig, struct k_sigac
+ k = &current->sighand->action[sig-1];

+ spin_lock_irq(&current->sighand->siglock);
+
+ if (oact)
+ *oact = *k;

```

```
@ @ -2338,8 +2358,7 @ @ int do_sigaction(int sig, struct k_sigac
* (for example, SIGCHLD), shall cause the pending signal to
* be discarded, whether or not it is blocked"
*/
- if (act->sa.sa_handler == SIG_IGN ||
- (act->sa.sa_handler == SIG_DFL && sig_kernel_ignore(sig))) {
+ if (sig_task_ignore(current, sig)) {
    struct task_struct *t = current;
    sigemptyset(&mask);
    sigaddset(&mask, sig);
```

Containers mailing list
Containers@lists.linux-foundation.org
<https://lists.linux-foundation.org/mailman/listinfo/containers>

Subject: Re: [PATCH 1/3] Signal semantics for /sbin/init
Posted by [Oleg Nesterov](#) on Sat, 01 Sep 2007 11:02:21 GMT
[View Forum Message](#) <> [Reply to Message](#)

On 08/31, sukadev@us.ibm.com wrote:

```
>
> -static int sig_ignored(struct task_struct *t, int sig)
> + // Currently this check is a bit racy with exec(),
> + // we can _simplify_ de_thread and close the race.
> + if (likely(!is_container_init(tsk->group_leader)))
> + return 0;
> +
> + if (!in_interrupt())
> + return 0;
```

I don't understand why you are trying to mix this patch with pid_ns changes.

We don't need in_interrupt() check unless we use current to decide if the signal goes from the parent namespace.

And in fact, I'd personally prefer to use "is_global_init()" for this patch, because it hopefully can fix the problems we have even without namespaces. This also matches the current check in get_signal_to_deliver().

Oleg.

Containers mailing list
Containers@lists.linux-foundation.org
<https://lists.linux-foundation.org/mailman/listinfo/containers>

Subject: Re: [PATCH 1/3] Signal semantics for /sbin/init
Posted by [Sukadev Bhattiprolu](#) on Mon, 03 Sep 2007 15:56:09 GMT
[View Forum Message](#) <> [Reply to Message](#)

Oleg Nesterov [oleg@tv-sign.ru] wrote:

| On 08/31, sukadev@us.ibm.com wrote:

| >
| > -static int sig_ignored(struct task_struct *t, int sig)
| > + // Currently this check is a bit racy with exec(),
| > + // we can _simplify_ de_thread and close the race.
| > + if (likely(!is_container_init(tsk->group_leader)))
| > + return 0;
| > +
| > + if (!in_interrupt())
| > + return 0;
|

| I don't understand why you are trying to mix this patch with pid_ns changes.

| We don't need in_interrupt() check unless we use current do decide if the
| signal goes from the parent namespace.

| And in fact, I'd personally prefer to use "is_global_init()" for this patch,
| because it hopefully can fix the problems we have even without namespaces.
| This also matches the current check in get_signal_to_deliver().

Sorry. I wasn't paying enough attention to this patch and including it
only for reference. Was planning to replace this with your final patch.
Or do you want me to fix the two bugs and resend ?

| Oleg.

Containers mailing list
Containers@lists.linux-foundation.org
<https://lists.linux-foundation.org/mailman/listinfo/containers>

Subject: Re: [PATCH 1/3] Signal semantics for /sbin/init
Posted by [Oleg Nesterov](#) on Mon, 03 Sep 2007 16:45:48 GMT
[View Forum Message](#) <> [Reply to Message](#)

On 09/03, sukadev@us.ibm.com wrote:

>
> Oleg Nesterov [oleg@tv-sign.ru] wrote:
> | On 08/31, sukadev@us.ibm.com wrote:
> | >
> | > -static int sig_ignored(struct task_struct *t, int sig)
> | > + // Currently this check is a bit racy with exec(),
> | > + // we can _simplify_ de_thread and close the race.

```
> | > + if (likely(!is_container_init(tsk->group_leader)))
> | > + return 0;
> | > +
> | > + if (!in_interrupt())
> | > + return 0;
> |
> | I don't understand why you are trying to mix this patch with pid_ns changes.
> |
> | We don't need in_interrupt() check unless we use current do decide if the
> | signal goes from the parent namespace.
> |
> | And in fact, I'd personally prefer to use "is_global_init()" for this patch,
> | because it hopefully can fix the problems we have even without namespaces.
> | This also matches the current check in get_signal_to_deliver().
>
> Sorry. I wasn't paying enough attention to this patch and including it
> only for reference. Was planning to replace this with your final patch.
> Or do you want me to fix the two bugs and resend ?
```

Sorry! I didn't have any time for the kernel hacking last days. There are some other minor (and not related) problems with the blocked signals which I'd like to check before doing the final patch.

Please feel free to do what you think right. I am going to KS this night, and I will be completely offline till september 10. Any chance we could delay this a bit? In any case, patches 2-3 should not depend on any further possible changes in this patch.

Oleg.

Containers mailing list
Containers@lists.linux-foundation.org
<https://lists.linux-foundation.org/mailman/listinfo/containers>
