Join GitHub today	Dismiss
GitHub is home to over 20 million developers working together to host and	
review code, manage projects, and build software together.	\diamond
 Sign up	

New issue

copy: avoid useless comparison for non-directory entries #6656

圮 Cor	nversatio	on 22	- Commits 1		+4	↓ - 5 ••••
	rher	r tzog cor	nmented on Aug 22		Reviewers evverx	Ţ.
	The acro entry This kern devi	compari oss multip y is a dire also fixe nel older ices num	son of the device number is meant to not follow directories ole filesystems, so the test is only relevant when the ectory. es the test-copy test when /tmp is an overlayfs with a than Linux 4.12 since overlayfs reported inconsistent bers. See https://bugs.debian.org/854400 for some details.		Assignees No one assigned Labels	
	¢ ()	evverx	y: avoid useless comparison for non-directory entries	✓ fba31f5View changes	Projects None yet	
		src/b 394	<pre>asic/copy.c - else if (S_ISDIR(buf.st_mode)) 201 in the state of (O_IODID(buf.st_mode)) { </pre>		Milestone No milestone	
			<pre>391 + else if (S_ISDIR(but.st_mode)) { 392 + if (buf.st_dev != original_device) 393 + continue;</pre>		B participants	
			evverx 29 days ago • edited Member Thank you. I'm not sure where this if should be placed. fd_copy_directory s copying as soon as it detects a filesystem boundary, while cpone-file-system an empty directory before stopping copying.	tops m creates		
			rhertzog 29 days ago Right. I made my change to be minimal while keeping the current behaviour. But y come over such a mount point, it's true that the underlying filesystem has at least directory for the mount to be possible. So arguably the behaviour of cpone-fil is more logical. Shall I update my patch to behave like cp ?	when you an empty le-system		
			evverx 29 days ago Member @rhertzog , I think that the patch fixes the bug and should be updated, but I'd rath someone else to confirm that skipping device nodes and bind-mounted files is not	ner wait for t intentional.		
			rhertzog 29 days ago The stat call on device nodes does not report the device number of the device itse filesystem hosting the device node, so this one is fine. You are right though that the the behaviour for bind-mounted files which would be copied with my patch while t skipped currently. That said copy_tree is only used by systemd-tmpfiles to set temporary files and by machined to copy files between host and containers. In be don't believe that we have any specific requirement to ignore bind mounted files.	elf but of the his changes hey are up the oth cases, I		
			evverx 29 days ago Member			

I'm sorry. Indeed, I was wrong about device nodes. I should have used mknod to check the behaviour of stat .

I don't believe that we have any specific requirement to ignore bind mounted files.

Neither do I.

	-		Pro-
1	8	3	9 11
-	3	-	3
(Internet			

poettering commented 24 days ago

The comparison of the device number is meant to not follow directories across multiple filesystems, so the test is only relevant when the entry is a directory.

I don't follow? On Linux mounts may be on either files and directories, and often are. Hence limiting any such checks to directories will make us miss any file mounts. This patch hence looks not OK.

poettering commented 24 days ago

Owner

Owner

On traditional UNIX the way to detect file system boundaries is via comparing st_dev. This is implemented in numerous tools, and as a fallback in systemd too. overlayfs breaks with that UNIX API if it changes st_dev within the same file system, but the right place to fix that really appears to be overlayfs, instead of making all the numerous tools work around overlayfs' peculiarities on this. I am not sure what overlayfs' strategy on this is, but it appears to be quite a steep compat break on their side...

poettering added **needs-discussion util-lib** labels 24 days ago \bigcirc



poettering commented 24 days ago

Owner

Owner

Owner

(oh, and we have similar checks in path_is_mount_point() and various other places iirc. Any such patch if accepted would need to cover all cases not just one specific one — but again I am not convinced this is really the right way to go...)



rhertzog commented 24 days ago

@poettering What can you mount on a file except another file?

In a copy_tree() function, you want to avoid infinite copying through (directory) mount points loops, so it's a safety measure to skip the mounted directories. But a mounted file does not pose any risk, does it?

That's why I believe it's OK to do the change here. But it would not make sense to change the more generic path_is_mount_point() .

overlayfs improved already in Linux 4.12 to have fewer such edge cases but the vast majority of the code out there does not care about st_dev on files. cp -rx does not for instance.



poettering commented 24 days ago

@poettering What can you mount on a file except another file?

IIRC the semantics on Linux are that:

- 1. dirs can be only mounted on dirs, and
- 2. symlinks not at all and
- 3. everything else on everything else (i.e. also device nodes on files, and files on device nodes, and any other weird shit)



poettering commented 24 days ago

In a copy_tree() function, you want to avoid infinite copying through (directory) mount points loops, so it's a safety measure to skip the mounted directories. But a mounted file does not pose any risk, does it?

Well, while "same-fs" checks are useful to avoid bind mount loops they have other uses too. For example, they often are conceptual boundaries. For example, you copy a fully set up OS image, then /proc, /sys/, ... and so on, are all mounts, which you want to avoid, but everything else should be copied just fine. and if somebody mounts /proc/core to /root/core, then it should also be avoided.

I am still not convinced that adding such a work-around for one specific fs that departs so strongly from accepted UNIX behaviour is the right way to go...

It's weird being on the side of arguing for UNIX here, instead of the other side, but here I am ;-)



hmm, interesting to know that those tools only do this check for directories. But it still feels hackish to special case dirs for that I must say... not sure what to do on this... I am still tempted to say that the overlayfs people really need to pass out valid st_dev values, and evreything else is compat breakage... And it appears they are even aware of the issue...

Maybe a more digestable version would be to check the fs magic value of the relevant subdir, to special case this. It's awful, but then at least we do the broken logic only on broken overlayfs...



rhertzog commented 23 days ago

If the "fix/work-around" is not desired, as far as I am concerned, I would be also happy if the failing test was just skipped when we detect overlayfs (and maybe Linux < 3.10). I'm not sure how to achieve this though.



	Member
m not sure that people who use C in tmpfiles.d/*.conf expect systemd-tr nat is different from cp -r -x or rsync -rone-file-system. In addition, r oundaries is not documented, so that might probably be surprising.	mpfiles to work in a way not crossing file system
Also, I think that not creating empty files and directories is still a bug.	
evverx commented 22 days ago	Member
Also, I think that not creating empty files and directories is still a bug.	
On the other hand, rsync doesn't create empty directories whenone-file- seems that some people might expect systemd-tmpfiles to do the same when	system is repeated, so it C is used.
oettering commented 22 days ago	Owner
If the "fix/work-around" is not desired, as far as I am concerned, I would be alwas just skipped when we detect overlayfs (and maybe Linux < 3.10). I'm not though.	so happy if the failing test sure how to achieve this retty poor choice for /tmp for
all its semantics. Most apps expect a fully featured, fast is there, and overlayis is all kinds of weird shortcomings and performance weirdnesses (since copy up ne	sn't really that at all, it has eds to happen all the time)
all its semantics. Most apps expect a fully featured, fast is there, and overlayis is all kinds of weird shortcomings and performance weirdnesses (since copy up ne hertzog commented 21 days ago	sn't really that at all, it has eeds to happen all the time)
all its semantics. Most apps expect a fully featured, fast is there, and overlayts is all kinds of weird shortcomings and performance weirdnesses (since copy up ne hertzog commented 21 days ago The build bots are using overlayfs in order to be able to easily discard changes f n which the build is triggered. Arguably we could put a tmpfs on /tmp but then w backages storing/generating large amount of data on /tmp. The build itself happe directory outside of the build chroot which is thus not under the control of overlay o not use /tmp but a temporary directory withing the build-tree itself.	sn't really that at all, it has eeds to happen all the time) from the initial clean chroot re could hit ENOSPC for ens in bind mount with a yfs so the last solution is
all its semantics. Most apps expect a fully featured, fast is there, and overlayis is all kinds of weird shortcomings and performance weirdnesses (since copy up ne hertzog commented 21 days ago The build bots are using overlayfs in order to be able to easily discard changes f n which the build is triggered. Arguably we could put a tmpfs on /tmp but then w backages storing/generating large amount of data on /tmp. The build itself happed directory outside of the build chroot which is thus not under the control of overlay o not use /tmp but a temporary directory withing the build-tree itself. But I think this is a bit off-topic, I know how to avoid the problem I just thought n systemd.	sn't really that at all, it has eeds to happen all the time) from the initial clean chroot re could hit ENOSPC for ens in bind mount with a yfs so the last solution is that it would best addressed
All its semantics. Most apps expect a fully featured, fast is there, and overlayis is all kinds of weird shortcomings and performance weirdnesses (since copy up ne hertzog commented 21 days ago The build bots are using overlayfs in order to be able to easily discard changes f in which the build is triggered. Arguably we could put a tmpfs on /tmp but then w backages storing/generating large amount of data on /tmp. The build itself happed directory outside of the build chroot which is thus not under the control of overlay o not use /tmp but a temporary directory withing the build-tree itself. But I think this is a bit off-topic, I know how to avoid the problem I just thought in systemd.	sn't really that at all, it has eeds to happen all the time) from the initial clean chroot re could hit ENOSPC for ens in bind mount with a yfs so the last solution is that it would best addressed Owner



evverx commented 20 days ago

Member

I think it would be great if systemd-tmpfiles copied files in a similar way to rsync. On the other hand, nobody appears to have ever pointed systemd-tmpfiles to /proc or noticed missing empty files and directories, so everything seems to be fine :-)