Subject: TCP: time wait bucket table overflow - memory leak? Posted by nksupport on Tue, 14 Jul 2009 08:45:35 GMT View Forum Message <> Reply to Message

Hi guys.

We've set up a brand new openVZ server. It's hardly ever loaded in normal mode.

At times when one of the VEs gets overloaded we start getting "TCP: time wait bucket table overflow" and the kmemsize bean counter grows.

What i expect from the node in this case is to kill the failed process or even the entire VE.

Instead, the node's LA spikes up to hundreds and then the entire node just dies.

vzctl stop times out. vzctl stop --fast times out. kill -9 `init of the VE` fails.

This looks like a memory leak to me.

I tried the solution to the error message described at http://bugzilla.openvz.org/show\_bug.cgi?id=460 - it did not work for me. Anyway, i doubt it's the root cause of our problem - looks like just one of the symptoms.

The main problem is that when a VE hits a counter (probably only kmemsize, but i'm not sure whether other limits trigger the same problem), the node dies itself instead of killing a VE. That's not what you'd expect from an encapsulated virtual server.

The node's normal production load is 0.05 to 0.40.

I've tried setting both loose and strict memory UBC limits - it didn't change anything. The current UBC limits are drakonian.

The server's 2xXeon L5410 (8 cores total) with 8G RAM running Centos 5.3 x64.

The kernel's 2.6.18-128.1.1.el5.028stab062.3 #1 SMP Sun May 10 18:54:51 MSD 2009 x86\_64

rpm -qa | grep vz ovzkernel-2.6.18-128.1.1.el5.028stab062.3 vzyum-2.4.0-11 vzrpm43-python-4.3.3-7\_nonptl.6 vzquota-3.0.12-1 ovzkernel-devel-2.6.18-128.1.1.el5.028stab062.3 vzctl-3.0.23-1 vzctl-lib-3.0.23-1 vzpkg-2.7.0-18 The VEs are different, centos and debian. I have confirmed the same behaviour on three different VEs running different OS releases: one of them hits the limit, node dies.

I've attached some debug output, hope someone can find a clue - so far i could not. I could really use a hand, thanks!

## File Attachments

1) vzstats, downloaded 632 times

2) vzctl\_strace, downloaded 595 times

3) sysctl, downloaded 590 times

4) vzctl\_fast\_strace, downloaded 619 times

5) dmesg, downloaded 624 times

Subject: Re: TCP: time wait bucket table overflow - memory leak? Posted by maratrus on Tue, 14 Jul 2009 15:57:26 GMT View Forum Message <> Reply to Message

Hello,

you gave a very good description of the problem.

Quote:

What i expect from the node in this case is to kill the failed process or even the entire VE.

But there is no such a killer in OpenVZ. Parallels Virtuozzo Containers contains subsystem which does what you have said but in OpenVZ you can adjust the behavior of the container via user\_beancounters i.e. if barriers and limits are really huge there is no a restrictive force which prevents a containers from consuming a lot of resources.

Quote:

Instead, the node's LA spikes up to hundreds and then the entire node just dies.

vzctl stop times out. vzctl stop --fast times out. kill -9 `init of the VE` fails.

Could you please say what means "the entire node just dies"? Are you able to ping the node? Are you able to invoke commands?

What is the reason of LA being so big?

Does the node perform CPU-consuming operations? What is the CPU state when "the node is being died"?

Or there are a lot of input-output operations being invoked?

Subject: Re: TCP: time wait bucket table overflow - memory leak? Posted by nksupport on Wed, 15 Jul 2009 14:36:47 GMT

maratrus wrote

But there is no such a killer in OpenVZ. Parallels Virtuozzo Containers contains subsystem which does what you have said but in OpenVZ you can adjust the behavior of the container via user\_beancounters i.e. if barriers and limits are really huge there is no a restrictive force which prevents a containers from consuming a lot of resources.

Yes there is! A process that failed to allocate memory will onviously be killed, isn't that right? I am talking exactly about UBC in my post. I'm using UBC to limit the VEs. Like i said, the current UBC limits are drakonian. I want all processes that try to overuse memory to be killed and OpenVZ is supposed to do it.

maratrus wrote

Could you please say what means "the entire node just dies"? Are you able to ping the node? Are you able to invoke commands?

What is the reason of LA being so big?

Does the node perform CPU-consuming operations? What is the CPU state when "the node is being died"?

Or there are a lot of input-output operations being invoked?

You know, it's hard to debug a server at LA of 500 I can still ping it, but it's just too slow. Random processes from VEs are in top, it behaves just like it would with unlimited UBC, i.e. instead of killing or throttling the runaway process it gives it more and more CPU. Eventually the node's ssh dies and this ends the show. Anyway - i attached my /proc/user\_beanconters and performance stats, please have a look if i missed anything.

Subject: Re: TCP: time wait bucket table overflow - memory leak? Posted by maratrus on Wed, 15 Jul 2009 16:27:40 GMT View Forum Message <> Reply to Message

Hi,

Quote:

I want all processes that try to overuse memory to be killed and OpenVZ is supposed to do it.

Where did you get such information from?

Quote:

You know, it's hard to debug a server at LA of 500

But it cannot be an instantaneous process. Is there a process that begins to consume a lot of CPU? What VE does it belong to? Subject: Re: TCP: time wait bucket table overflow - memory leak? Posted by nksupport on Fri, 17 Jul 2009 16:29:09 GMT View Forum Message <> Reply to Message

hi maratrus. Thanks for your help so far

maratrus wrote on Wed, 15 July 2009 19:27 Quote: I want all processes that try to overuse memory to be killed and OpenVZ is supposed to do it.

Where did you get such information from?

openvz refuses resource allocation when the process requests a resource that is over limit. That's a fact, right?

a process that failed to allocate RAM will just die - that's common sense. I'm not really telling that there's an openvz feature that kills the process. Normal processes - mysql, apache, bash, common stuff - will die when they can't allocate memory, fork, etc. Maybe some custom badly written code won't, but commonly used software does.

Quote:Quote: You know, it's hard to debug a server at LA of 500

But it cannot be an instantaneous process. Is there a process that begins to consume a lot of CPU? What VE does it belong to?

Just a random common process belonging to the VE that has hit the kmemsize limit. It could be apache, mysql, named, bash, ssh - and it's not like a single process remains in top. All processes of this VE begin rotating in the node's top. This may sound shady, but please take into account that i only have several minutes to collect data before ssh dies.

Subject: Re: TCP: time wait bucket table overflow - memory leak? Posted by nksupport on Fri, 17 Jul 2009 20:31:26 GMT View Forum Message <> Reply to Message

funny thing. I wasn't actually right about "the failed VE's processes popping up in top".

top - 22:32:30 up 7 days, 8:04, 2 users, load average: 49.88, 39.72, 23.21 Tasks: 390 total, 2 running, 387 sleeping, 0 stopped, 1 zombie Cpu0 : 1.0%us, 7.0%sy, 0.0%ni, 91.7%id, 0.3%wa, 0.0%hi, 0.0%si, 0.0%st Cpu1 : 1.0%us, 6.0%sy, 0.0%ni, 93.0%id, 0.0%wa, 0.0%hi, 0.0%si, 0.0%st Cpu2 : 0.0%us, 0.0%sy, 0.0%ni, 0.0%id, 0.0%wa, 0.0%hi, 0.0%si, 0.0%st Cpu3 : 2.6%us, 13.8%sy, 0.0%ni, 83.6%id, 0.0%wa, 0.0%hi, 0.0%si, 0.0%st Cpu4 : 2.0%us, 12.8%sy, 0.0%ni, 85.2%id, 0.0%wa, 0.0%hi, 0.0%si, 0.0%st Cpu5 : 1.7%us, 12.2%sy, 0.0%ni, 85.2%id, 0.0%wa, 0.0%hi, 0.0%si, 0.0%st Cpu6 : 1.3%us, 10.6%sy, 0.0%ni, 87.8%id, 0.3%wa, 0.0%hi, 0.0%si, 0.0%st Cpu7 : 1.3%us, 8.2%sy, 0.0%ni, 90.5%id, 0.0%wa, 0.0%hi, 0.0%si, 0.0%st Mem: 8164444k total, 8119288k used, 45156k free, 149272k buffers This has happened just now. I've killed the applications running on the failed VE. NOTHING is using the CPU at all. LA is growing. I'm floored. top only shows its own "top" process - and that's natural... there's nothing launched except it. the LA is already above 50. I can't restart the failed VE, can't raise the kmemsize or anything. The error about table overflow is there in dmesg. I am now nearly sure it's a kernel bug - i can't imagine any other reason.

This may be the last time i'm facing this - i have loads of RAM and have been raising memory limits for all failing VEs. It's been over a week of stable work, the next time will probably happen in August if it will happen at all. that entire situation just sucks.

Subject: Re: TCP: time wait bucket table overflow - memory leak? Posted by maratrus on Mon, 20 Jul 2009 06:23:18 GMT View Forum Message <> Reply to Message

Hi,

Quote:

I am now nearly sure it's a kernel bug - i can't imagine any other reason.

You can file a new bugreport whenever you want. Please don't hesitate doing it. I just want to clarify the situation because I don't understand clearly what your problem consists of.

The last "top" output doesn't contain anything frightening from my point of view. A "huge" load average is just a consequence of a "great" number of process. A good idea is to obtain the status of these processes. You can do it with help of "ps" utility. Please, examine their states.

Quote:

I can't restart the failed VE, can't raise the kmemsize or anything.

Why do you want to increase kmemsize? The only failcounters I can see on proivided vzstats output are those concerning with privvmpages not kememsize. What process are running in a failed VE?

You have mentioned that ssh dies when this occurs. Don't you think that this problem relates to network? Do you have a direct access to that server to confirm that the node completely dies?

Do you have any specific settings, for example nfs inside VE?

Please the next time it occurs please gather alt-sysrq-\* information from the problem node.

Alt-sysrq-

1) "m" - for memory info dump

2) "p" - for registers - several times, please, twice the number of CPUs

3) "a" - for scheduler stat - 3 times

4) "w" - another scheduler info - 3 times

5) "t" - for all processes calltraces. Warning - this is a resource consuming operation. At least - twice.

To gather it you probably have to install serial console http://wiki.openvz.org/Remote\_console\_setup#Serial\_console

Subject: Re: TCP: time wait bucket table overflow - memory leak? Posted by nksupport on Mon, 20 Jul 2009 06:56:20 GMT View Forum Message <> Reply to Message

Quote:

I just want to clarify the situation because I don't understand clearly what your problem consists of.

me neither, that's why i am attempting to clarify it before filing.

Quote: The last "top" output doesn't contain anything frightening from my point of view.

A "huge" load average is just a consequence of a "great" number of process. A good idea is to obtain the status of these processes. You can do it with help of "ps" utility. Please, examine their states.

It is freaking frightening because i can't understand it! There are no "great number" of processes. ps doesn't show anything unusual. The top output i provided shows a completely idle system what are those tons of active processes? I will try to find something in the process table next time, but i had a brief look before - there's nothing strange there.

Quote:Why do you want to increase kmemsize? The only failcounters I can see on proivided vzstats output are those concerning with privvmpages not kememsize.

i posted UBC taken during normal operation. The UBC i posted is just the info about limits set in my system. Like i wrote, kmemsize counters start growing when i encounter this issue: i hope you can trust me this much

Quote: What process are running in a failed VE?

like i wrote, during the last test i killed most of the apps within the failed VE. There was pretty much nothing running except for init and a few basic daemons. I'll try to grab ps next time.

Quote:You have mentioned that ssh dies when this occurs. Don't you think that this problem relates to network? Do you have a direct access to that server to confirm that the node completely dies?

man, the server behaves exactly like you'd expect from an overloaded server i just can't find the reason for this overload.

it's not a networking issue. Local login stops working as well - it just times out since the system

becomes too slow to process password or rsa auth. The actual question is why the load grows infinitely: this is our problem.

Quote:Do you have any specific settings, for example nfs inside VE? No, there's no crap. 5 basic LAMP VEs.

Will try to collect extra data - looks like we have nothing useful yet... Thanks anyway!

Subject: Re: TCP: time wait bucket table overflow - memory leak? Posted by maratrus on Mon, 20 Jul 2009 10:26:30 GMT View Forum Message <> Reply to Message

Quote:

The top output i provided shows a completely idle system - what are those tons of active processes?

Not only a great number of active processes but also processes in "D" state are taken into account when "load average" is calculated.

Quote:

like i wrote, during the last test i killed most of the apps within the failed VE. There was pretty much nothing running except for init and a few basic daemons. I'll try to grab ps next time.

Moreover you can try to stop a problem VE and get alt-sysrq-\* output at that moment to find out where "vzctl" process sticks.

Quote:

The actual question is why the load grows infinitely: this is our problem.

What kind out of load you have in mind? Load average numbers?

Subject: Re: TCP: time wait bucket table overflow - memory leak? Posted by nksupport on Mon, 20 Jul 2009 10:29:55 GMT View Forum Message <> Reply to Message

yes, LA.

Subject: Re: TCP: time wait bucket table overflow - memory leak?

## Quote:

Tasks: 390 total, 2 running, 387 sleeping, 0 stopped, 1 zombie

sleeping tasks may wait semaphores.

Subject: Re: TCP: time wait bucket table overflow - memory leak? Posted by nksupport on Mon, 17 Aug 2009 13:56:06 GMT View Forum Message <> Reply to Message

hi Maratrus.

Quote:Please the next time it occurs please gather alt-sysrq-\* information from the problem node.

Please see attached. I didn't find anything interesting, maybe you know where to look. The server was overloaded during the dump, had to powercycle in at the end.

File Attachments
1) kernel.log.gz, downloaded 508 times

Subject: Re: TCP: time wait bucket table overflow - memory leak? Posted by atxadmin on Tue, 25 Aug 2009 11:28:25 GMT View Forum Message <> Reply to Message

Can anyone help us get this problem resolved as the server went down an additional 2 times within the last week.

RAM modules have been replaced as well so that is also not the cause of the problem.

Any help is greatly appreciated.

Subject: Re: TCP: time wait bucket table overflow - memory leak? Posted by seanfulton on Thu, 27 Aug 2009 13:46:43 GMT View Forum Message <> Reply to Message

I am not sure if this will help but here is my \$.02:

I have had something similar happen when we over-allocated the RAM on the HN. We had 6G in the machine and had four VEs that collectively had access to 18G. The HN can't stuff 18G of crap into a 6G bag so it keeps going until it dies. The solution for us was to run ve-split for four (in your case five) containers and start with that. Then adjust resources \*carefully\*.

Another issue we have had (and are still having) has to do with I/O activity on the HN during nightly backups. We tar each VE nightly (used to use vzdump, but it uses tar as well, so there was no functional difference in this respect). The tar used to be done to an NFS mount, now we use ssh/dd to pipe it to the backup server.

Anyway, tar -czv \* knocks the machine to its knees. The VE it is backing up will effectively lock. On a busy mail server (inside a VE), this causes all sorts of damage to ongoing mail connections, so bad that we now have to restart the VE after every nightly backup.

In my experience, OpenVZ is very good about enforcing limits of activities within a VE, but can not seem to contain activities on the HN itself. It has been suggested that this is a kernel bug in the 2.6.18 kernel, but no solutions have been found or suggested.

One solution recommended was to change the scheduler on the drive from cpq to deadline like this:

For a drive that is /dev/sda, use: # cat /sys/block/sda/queue/scheduler # cat /sys/block/sda/queue/scheduler

This still caused the load to grow, but the VEs no longer completely lock up during the process, which is a plus.

I am still looking for a more lasting solution to this last issue, but that's not your problem.

Hopefully this experience will point you in a useful direction. Good luck. I know it's frustrating.

sean

Subject: Re: TCP: time wait bucket table overflow - memory leak? Posted by nksupport on Thu, 27 Aug 2009 17:10:28 GMT View Forum Message <> Reply to Message

hi Sean!

Thanks for your two cents - you keep our hope alive!

This isn't actually an overallocation - like i wrote in the first post, i already tried draconian limits using only 10% of the server. Anyway, your suggestion is good - we're using vzsplit.

As for tar - you should try making local copies instead of tunneling it. Tar to a local drive should show you much better performance - at least that's my experience. And no, there is no pattern in our failures, they're not caused by backups. However, your assumption was good.