



Sachin Malhotra [Follow](#)

Masters Student at University of Southern California |

Student Researcher @ISI | [edorado93.github.io](https://github.com/edorado93)

Mar 26, 2017 · 9 min read

## Load Testing HAProxy (Part 2)



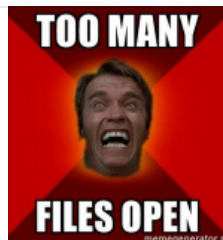
This is the second part in the 3 part series on performance testing of the famous TCP load balancer and reverse proxy, HAProxy. If you haven't

gone through the previous post, I would highly suggest you do so to get some sort of context.

### Load Testing HAProxy (Part-1)

Load Testing ? HAProxy ? If all this seems greek to you, don't worry. I will provide inline links to rea...

[medium.com](#)



This post will focus on the TCP Port exhaustion problem and how we can deal with it. In the last post we talked about how we can tune the kernel level and process level ulimit settings. This post is focussed on modifying the sysctl settings to get over the port exhaustion limits.

## SYSCTL Local Port Range and Orphaned Sockets

Port exhaustion is a problem that will cause TCP communications with other machines over the network to fail. Most of the times there is a single process that leads to this problem and restarting it will fix the issue, temporarily. It will however come back to bite in a few hours or days depending on the system load.

Port exhaustion does not mean that the ports actually get tired. Of course, that is not possible because the computer is not human and the ports are not capable of getting tired. The truth is much more insidious.

Port exhaustion simply means that the system does not have any more *ephemeral ports* left to communicate with other machines / servers.

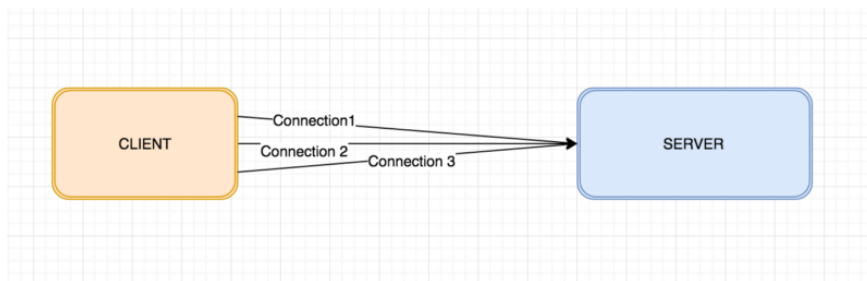
Before going further, *let us understand what constitutes a TCP connection and what really does an inbound and an outbound connection means.*

In majority of the cases whenever we talk about TCP connections and high scalability and ability to support concurrent connections, we usually refer to the number of inbound connections.

Say, the HAProxy is listening on port 443 for new inbound connections. If we say that the HAProxy can support X number of concurrent connections, what we really mean are X number of incoming connections and all of them are established on port 443 on the HAProxy machine.

If these connections are inbound for the HAProxy, then these have to be *outbound for the client machines where the connection originated*. Any sort of communication from the client requires them to initiate outbound connections to the servers.

*When a connection is established over TCP, a socket is created on both the local and the remote host. These sockets are then connected to create a socket pair, which is described by a unique 4-tuple consisting of the local IP address and port along with the remote IP address and port.*



3 TCP connections from client to server/proxy

If you understood the concept of the quadruple, you will realise that in an outbound connection or rather multiple outbound connections to the SAME backend server, 2 things always remain the same i.e. Destination IP and Destination Port. Assuming we are only taking into account a single client machine, the client IP will also remain the same.

This means that the number of outbound connections is dependent on the number of client ports that can be used for establishing the connection. While establishing an outbound connection, the source port is randomly selected from the ephemeral port range and this port gets freed up once the connection is destroyed. That's why such ports are called as ephemeral ports.

By default, the total number of local ephemeral ports available are around 28000.

```
sachinm@ip-192-168-0-100:~$ sysctl net.ipv4.ip_local_port_range
net.ipv4.ip_local_port_range = 32768    61000
```

Now you might be thinking that 28k is a pretty large number and what can possibly cause 28k connections to get used up at a single point of time? In order to understand this, **we have to understand the TCP connection lifecycle**.

During the TCP handshake, the connection state goes from

SYN\_SENT → SYN\_RECV → ESTABLISHED. Once the connection is in ESTABLISHED state, it means that the TCP connection is now active. *However, once the connection is terminated, the local port that was being used earlier does not become active immediately.*

The connection enters a state known as the *TIME\_WAIT* state for a period of 120 seconds before it is finally terminated. This is a kernel level setting that exists to allow any delayed or out of order packets to be ignored by the network.

*If you do the math, it won't take more than **230 concurrent connections per second** before the supposedly large limit of 28000 ephemeral ports on the system is reached. This limit is very easy to reach on proxies like HAProxy or NGINX because all the traffic is routed through them to the backend servers.*

When a connection enters the *TIME\_WAIT* state, it is known as an **orphaned socket** because the TCP socket in this case is not help by any socket descriptor but are still held by the system for the designated time i.e. 120 seconds by default.

## How to detect this?

Enough with all the theoretical stuff. Let's jump in and see how we can identify if this limit has been hit on the system. There are two commands I absolutely love to use to find out the number of TCP connections established on the system.

## ss (Socket Statistics)

The socket statistics command is a sort of replacement of the famous netstat command and is much faster than the netstat command in rendering information because it fetches the connections info directly from the kernel space. In order to get a hang of the different options supported by the ss command, check out

### 10 examples of Linux ss command to monitor network connections

In a previous tutorial we saw how to use the netstat command to get statistics on...

[www.binarytides.com](http://www.binarytides.com)

ss -ttn		
State	Local Address:Port	Remote Address:Port
ESTAB	192.168.1.2:33075	74.125.22.0
ESTAB	192.168.1.2:33076	74.125.22.0
ESTAB	192.168.1.2:33077	74.125.22.0
ESTAB	192.168.1.2:33078	74.125.22.0
ESTAB	192.168.1.2:33079	74.125.22.0
ESTAB	192.168.1.2:33080	74.125.22.0
ESTAB	192.168.1.2:33081	74.125.22.0
ESTAB	192.168.1.2:33082	74.125.22.0
ESTAB	192.168.1.2:33083	74.125.22.0
ESTAB	192.168.1.2:33084	74.125.22.0
ESTAB	192.168.1.2:33085	74.125.22.0
ESTAB	192.168.1.2:33086	74.125.22.0
ESTAB	192.168.1.2:33087	74.125.22.0
ESTAB	192.168.1.2:33088	74.125.22.0
ESTAB	192.168.1.2:33089	74.125.22.0
ESTAB	192.168.1.2:33090	74.125.22.0
ESTAB	192.168.1.2:33091	74.125.22.0
ESTAB	192.168.1.2:33092	74.125.22.0
ESTAB	192.168.1.2:33093	74.125.22.0
ESTAB	192.168.1.2:33094	74.125.22.0
ESTAB	192.168.1.2:33095	74.125.22.0
ESTAB	192.168.1.2:33096	74.125.22.0
ESTAB	192.168.1.2:33097	74.125.22.0
ESTAB	192.168.1.2:33098	74.125.22.0
ESTAB	192.168.1.2:33099	74.125.22.0
ESTAB	192.168.1.2:33100	74.125.22.0
ESTAB	192.168.1.2:33101	74.125.22.0
ESTAB	192.168.1.2:33102	74.125.22.0
ESTAB	192.168.1.2:33103	74.125.22.0
ESTAB	192.168.1.2:33104	74.125.22.0
ESTAB	192.168.1.2:33105	74.125.22.0
ESTAB	192.168.1.2:33106	74.125.22.0
ESTAB	192.168.1.2:33107	74.125.22.0
ESTAB	192.168.1.2:33108	74.125.22.0
ESTAB	192.168.1.2:33109	74.125.22.0
ESTAB	192.168.1.2:33110	74.125.22.0
ESTAB	192.168.1.2:33111	74.125.22.0
ESTAB	192.168.1.2:33112	74.125.22.0
ESTAB	192.168.1.2:33113	74.125.22.0
ESTAB	192.168.1.2:33114	74.125.22.0
ESTAB	192.168.1.2:33115	74.125.22.0
ESTAB	192.168.1.2:33116	74.125.22.0
ESTAB	192.168.1.2:33117	74.125.22.0
ESTAB	192.168.1.2:33118	74.125.22.0
ESTAB	192.168.1.2:33119	74.125.22.0
ESTAB	192.168.1.2:33120	74.125.22.0
ESTAB	192.168.1.2:33121	74.125.22.0
ESTAB	192.168.1.2:33122	74.125.22.0
ESTAB	192.168.1.2:33123	74.125.22.0
ESTAB	192.168.1.2:33124	74.125.22.0
ESTAB	192.168.1.2:33125	74.125.22.0
ESTAB	192.168.1.2:33126	74.125.22.0
ESTAB	192.168.1.2:33127	74.125.22.0
ESTAB	192.168.1.2:33128	74.125.22.0
ESTAB	192.168.1.2:33129	74.125.22.0
ESTAB	192.168.1.2:33130	74.125.22.0
ESTAB	192.168.1.2:33131	74.125.22.0
ESTAB	192.168.1.2:33132	74.125.22.0
ESTAB	192.168.1.2:33133	74.125.22.0
ESTAB	192.168.1.2:33134	74.125.22.0
ESTAB	192.168.1.2:33135	74.125.22.0
ESTAB	192.168.1.2:33136	74.125.22.0
ESTAB	192.168.1.2:33137	74.125.22.0
ESTAB	192.168.1.2:33138	74.125.22.0
ESTAB	192.168.1.2:33139	74.125.22.0
ESTAB	192.168.1.2:33140	74.125.22.0
ESTAB	192.168.1.2:33141	74.125.22.0
ESTAB	192.168.1.2:33142	74.125.22.0
ESTAB	192.168.1.2:33143	74.125.22.0
ESTAB	192.168.1.2:33144	74.125.22.0
ESTAB	192.168.1.2:33145	74.125.22.0
ESTAB	192.168.1.2:33146	74.125.22.0
ESTAB	192.168.1.2:33147	74.125.22.0
ESTAB	192.168.1.2:33148	74.125.22.0
ESTAB	192.168.1.2:33149	74.125.22.0
ESTAB	192.168.1.2:33150	74.125.22.0
ESTAB	192.168.1.2:33151	74.125.22.0
ESTAB	192.168.1.2:33152	74.125.22.0
ESTAB	192.168.1.2:33153	74.125.22.0
ESTAB	192.168.1.2:33154	74.125.22.0
ESTAB	192.168.1.2:33155	74.125.22.0
ESTAB	192.168.1.2:33156	74.125.22.0
ESTAB	192.168.1.2:33157	74.125.22.0
ESTAB	192.168.1.2:33158	74.125.22.0
ESTAB	192.168.1.2:33159	74.125.22.0
ESTAB	192.168.1.2:33160	74.125.22.0
ESTAB	192.168.1.2:33161	74.125.22.0
ESTAB	192.168.1.2:33162	74.125.22.0
ESTAB	192.168.1.2:33163	74.125.22.0
ESTAB	192.168.1.2:33164	74.125.22.0
ESTAB	192.168.1.2:33165	74.125.22.0
ESTAB	192.168.1.2:33166	74.125.22.0
ESTAB	192.168.1.2:33167	74.125.22.0
ESTAB	192.168.1.2:33168	74.125.22.0
ESTAB	192.168.1.2:33169	74.125.22.0
ESTAB	192.168.1.2:33170	74.125.22.0
ESTAB	192.168.1.2:33171	74.125.22.0
ESTAB	192.168.1.2:33172	74.125.22.0
ESTAB	192.168.1.2:33173	74.125.22.0
ESTAB	192.168.1.2:33174	74.125.22.0
ESTAB	192.168.1.2:33175	74.125.22.0
ESTAB	192.168.1.2:33176	74.125.22.0
ESTAB	192.168.1.2:33177	74.125.22.0
ESTAB	192.168.1.2:33178	74.125.22.0
ESTAB	192.168.1.2:33179	74.125.22.0
ESTAB	192.168.1.2:33180	74.125.22.0
ESTAB	192.168.1.2:33181	74.125.22.0
ESTAB	192.168.1.2:33182	74.125.22.0
ESTAB	192.168.1.2:33183	74.125.22.0
ESTAB	192.168.1.2:33184	74.125.22.0
ESTAB	192.168.1.2:33185	74.125.22.0
ESTAB	192.168.1.2:33186	74.125.22.0
ESTAB	192.168.1.2:33187	74.125.22.0
ESTAB	192.168.1.2:33188	74.125.22.0
ESTAB	192.168.1.2:33189	74.125.22.0
ESTAB	192.168.1.2:33190	74.125.22.0
ESTAB	192.168.1.2:33191	74.125.22.0
ESTAB	192.168.1.2:33192	74.125.22.0
ESTAB	192.168.1.2:33193	74.125.22.0
ESTAB	192.168.1.2:33194	74.125.22.0
ESTAB	192.168.1.2:33195	74.125.22.0
ESTAB	192.168.1.2:33196	74.125.22.0
ESTAB	192.168.1.2:33197	74.125.22.0
ESTAB	192.168.1.2:33198	74.125.22.0
ESTAB	192.168.1.2:33199	74.125.22.0
ESTAB	192.168.1.2:33200	74.125.22.0
ESTAB	192.168.1.2:33201	74.125.22.0
ESTAB	192.168.1.2:33202	74.125.22.0
ESTAB	192.168.1.2:33203	74.125.22.0
ESTAB	192.168.1.2:33204	74.125.22.0
ESTAB	192.168.1.2:33205	74.125.22.0
ESTAB	192.168.1.2:33206	74.125.22.0
ESTAB	192.168.1.2:33207	74.125.22.0
ESTAB	192.168.1.2:33208	74.125.22.0
ESTAB	192.168.1.2:33209	74.125.22.0
ESTAB	192.168.1.2:33210	74.125.22.0
ESTAB	192.168.1.2:33211	74.125.22.0
ESTAB	192.168.1.2:33212	74.125.22.0
ESTAB	192.168.1.2:33213	74.125.22.0
ESTAB	192.168.1.2:33214	74.125.22.0
ESTAB	192.168.1.2:33215	74.125.22.0
ESTAB	192.168.1.2:33216	74.125.22.0
ESTAB	192.168.1.2:33217	74.125.22.0
ESTAB	192.168.1.2:33218	74.125.22.0
ESTAB	192.168.1.2:33219	74.125.22.0
ESTAB	192.168.1.2:33220	74.125.22.0
ESTAB	192.168.1.2:33221	74.125.22.0
ESTAB	192.168.1.2:33222	74.125.22.0
ESTAB	192.168.1.2:33223	74.125.22.0
ESTAB	192.168.1.2:33224	74.125.22.0
ESTAB	192.168.1.2:33225	74.125.22.0
ESTAB	192.168.1.2:33226	74.125.22.0
ESTAB	192.168.1.2:33227	74.125.22.0
ESTAB	192.168.1.2:33228	74.125.22.0
ESTAB	192.168.1.2:33229	74.125.22.0
ESTAB	192.168.1.2:33230	74.125.22.0
ESTAB	192.168.1.2:33231	74.125.22.0
ESTAB	192.168.1.2:33232	74.125.22.0
ESTAB	192.168.1.2:33233	74.125.22.0
ESTAB	192.168.1.2:33234	74.125.22.0
ESTAB	192.168.1.2:33235	74.125.22.0
ESTAB	192.168.1.2:33236	74.125.22.0
ESTAB	192.168.1.2:33237	74.125.22.0
ESTAB	192.168.1.2:33238	74.125.22.0
ESTAB	192.168.1.2:33239	74.125.22.0
ESTAB	192.168.1.2:33240	74.125.22.0
ESTAB	192.168.1.2:33241	74.125.22.0
ESTAB	192.168.1.2:33242	74.125.22.0
ESTAB	192.168.1.2:33243	74.125.22.0
ESTAB	192.168.1.2:33244	74.125.22.0
ESTAB	192.168.1.2:33245	74.125.22.0
ESTAB	192.168.1.2:33246	74.125.22.0
ESTAB	192.168.1.2:33247	74.125.22.0
ESTAB	192.168.1.2:33248	74.125.22.0
ESTAB	192.168.1.2:33249	74.125.22.0
ESTAB	192.168.1.2:33250	74.125.22.0
ESTAB	192.168.1.2:33251	74.125.22.0
ESTAB	192.168.1.2:33252	74.125.22.0
ESTAB	192.168.1.2:33253	74.125.22.0
ESTAB	192.168.1.2:33254	74.125.22.0
ESTAB	192.168.1.2:33255	74.125.22.0
ESTAB	192.168.1.2:33256	74.125.22.0
ESTAB	192.168.1.2:33257	74.125.22.0
ESTAB	192.168.1.2:33258	74.125.22.0
ESTAB	192.168.1.2:33259	74.125.22.0
ESTAB	192.168.1.2:33260	74.125.22.0
ESTAB	192.168.1.2:33261	74.125.22.0
ESTAB	192.168.1.2:33262	74.125.22.0
ESTAB	192.168.1.2:33263	74.125.22.0
ESTAB	192.168.1.2:33264	74.125.22.0
ESTAB	192.168.1.2:33265	74.125.22.0
ESTAB	192.168.1.2:33266	74.125.22.0
ESTAB	192.168.1.2:33267	74.125.22.0
ESTAB	192.168.1.2:33268	74.125.22.0
ESTAB	192.168.1.2:33269	74.125.22.0
ESTAB	192.168.1.2:33270	74.125.22.0
ESTAB	192.168.1.2:33271	74.125.22.0
ESTAB	192.168.1.2:33272	74.125.22.0
ESTAB	192.168.1.2:33273	74.125.22.0
ESTAB	192.168.1.2:33274	74.125.22.0
ESTAB	192.168.1.2:33275	74.125.22.0
ESTAB	192.168.1.2:33276	74.125.22.0
ESTAB	192.168.1.2:33277	74.125.22.0
ESTAB	192.168.1.2:33278	74.125.22.0
ESTAB	192.168.1.2:33279	74.125.22.0
ESTAB	192.168.1.2:33280	74.125.22.0
ESTAB	192.168.1.2:33281	74.125.22.0
ESTAB	192.168.1.2:33282	74.125.22.0
ESTAB	192.168.1.2:33283	74.125.22.0
ESTAB	192.168.1.2:33284	74.125.22.0
ESTAB	192.168.1.2:33285	74.125.22.0
ESTAB	192.168.1.2:33286	74.125.22.0
ESTAB	192.168.1.2:33287	74.125.22.0
ESTAB	1	

The `netstat` command is a very famous command that provides information about all sorts of connections established on the machine's networking stack.

```
sudo netstat -anptl
```

This will show you the details about all the connections on the machine. The details include

- local address
- remote address
- connection state
- process pid

We can also use this to see if a single process has established 28k connections to an outbound server which gives us insights into the port exhaustion problem.

```
sachinm@ip-192-168-0-122:~$ sudo netstat -anptl | grep 443 | head
tcp      0      0 192.168.0.122:6018    192.168.0.168:443    ESTABLISHED 9758/vegeta
tcp      0      0 192.168.0.122:37607   192.168.0.168:443    ESTABLISHED 9758/vegeta
tcp      0      0 192.168.0.122:5195    192.168.0.168:443    ESTABLISHED 9758/vegeta
tcp      0      0 192.168.0.122:20232   192.168.0.168:443    ESTABLISHED 9758/vegeta
tcp      0      0 192.168.0.122:7894    192.168.0.168:443    ESTABLISHED 9758/vegeta
tcp      0      0 192.168.0.122:40408    192.168.0.168:443    ESTABLISHED 9758/vegeta
tcp      0      0 192.168.0.122:6831    192.168.0.168:443    ESTABLISHED 9758/vegeta
tcp      0      0 192.168.0.122:54080    192.168.0.168:443    ESTABLISHED 9758/vegeta
tcp      0      0 192.168.0.122:46377    192.168.0.168:443    ESTABLISHED 9758/vegeta
tcp      0      0 192.168.0.122:63188    192.168.0.168:443    ESTABLISHED 9758/vegeta
```

For eg:- the above image shows that a process with pid 9758 has established multiple connections with the foreign machine with IP 192.168.0.168 and port 443. As we can clearly see, on the source side of things, there are numerous ports being used.

```
sachinm@ip-192-168-0-122:~$ sudo netstat -anptl | grep  
'192.168.0.168:443' | cut -c69-79 | sort | uniq -c | sort -  
rn
```

```
5670 ESTABLISHED
```

This modified command will show the status of the different connections established with 192.168.0.168 on port 443. Currently there are 5670 connections. If this limit were to reach 28k, then you should look at options to increase the ephemeral port range on the machine.

Let's look at another interesting command that you can issue at the server end or the proxy end to find out how many inbound connections have been established and by which IPs. So for example check out the result of the below command

```
ss -tan 'sport = :443' | awk '{print $(NF)}' | sed  
's/:[^ ]*//g' | sort | uniq -c
```



```
sachin@ip-192-168-0-168:~$ ss -tan 'sport = :443' | awk '{print $(NF)}' | sed 's/:[^ ]*//g' | sort | uniq -c
1 * *
2342 192.168.0.122 192.168.0.168
2340 192.168.0.123 192.168.0.168
2344 192.168.0.124 192.168.0.168
2336 192.168.0.125 192.168.0.168
2338 192.168.0.126 192.168.0.168
1018 192.168.0.138 192.168.0.168
2347 192.168.0.141 192.168.0.168
2349 192.168.0.142 192.168.0.168
2338 192.168.0.170 192.168.0.168
2339 192.168.0.179 192.168.0.168
2339 192.168.0.18 192.168.0.168
2341 192.168.0.232 192.168.0.168
2337 192.168.0.244 192.168.0.168
2341 192.168.0.5 192.168.0.168
2338 192.168.0.58 192.168.0.168
2339 192.168.0.59 192.168.0.168
2340 192.168.0.68 192.168.0.168
```

This shows that there are about 14 different machines who have established around 2300 connections each with 192.168.0.168 and if you look at the command closely, we have filtered out results only for port 443.

Enough with finding the problem already. Let's dive straight into finding the solution(s) to this problem.

## What's the way out?



Well don't be afraid because sysctl just happens to be a friendly monster. There are many ways by which we can solve this problem.

## Approach #1

One of the most practical approaches to solve this problem and one that you most likely will or rather should end up doing is to increase the

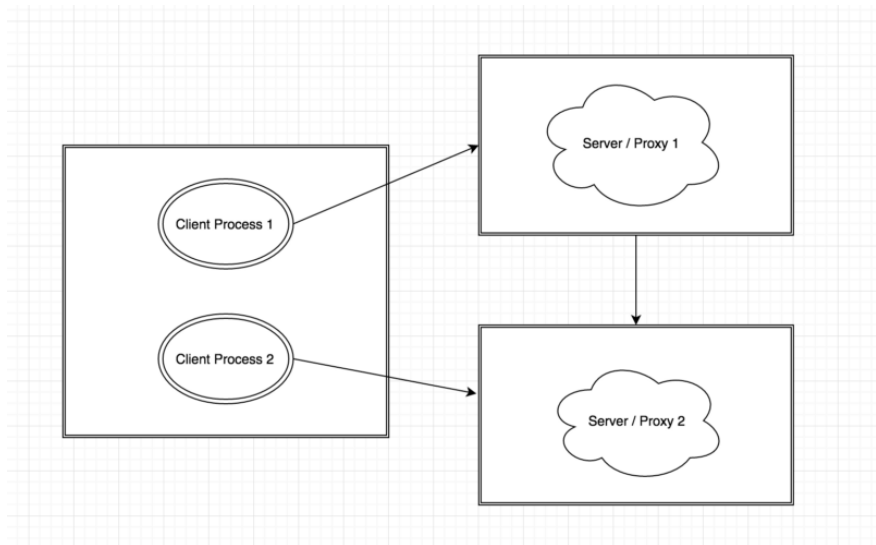
local ephemeral port range to the maximum possible value. As mentioned before, the default range is very small.

```
echo 1024 65535 > /proc/sys/net/ipv4/ip_local_port_range
```

This will increase the local port range to a bigger value. We cannot increase the range beyond this as there can only be a maximum of 65535 ports and the first 1024 are reserved for select services and purposes.

Note that you might still get bottleneck on this issue. However, instead of 28000 ports being used locally, it will be 64000 ports. Not a full proof solution but this is something that you can do to give you some breathing room.

Does this mean I can only get about 64k concurrent connections from a single client machine? The answer is NO.



In this scenario, a single client machine will be able to generate about 120k concurrent connections because both the processes are connecting to two different backend servers or proxies and hence different destination IPs.

## Approach #2

Another simple solution is to enable a Linux TCP option called ***tcp\_tw\_reuse***. This option enables the Linux kernel to reclaim a connection slot from a connection in TIME\_WAIT state and reallocate it to a new connection.

```
--> vim /etc/sysctl.conf
```

```
--> Add the following line in the end  
# Allow reuse of sockets in TIME_WAIT state for new  
connections  
# only when it is safe from the network stack's perspective.  
net.ipv4.tcp_tw_reuse = 1  
  
--> Reload sysctl settings  
sysctl -p
```

## Approach #3

Use more server ports. Till now we have talked about port exhaustion problems arising because in the quadruplet logic discussed before, the destination Ip, destination port and source Ip remain constant. The only thing that changes is the client ports.

However, if the server listens on two different ports instead of one, then we have twice the number of ephemeral ports available instead of one. This clubbed with the first approach gives you about 120k concurrent connections on a single machine.

You have to however take care that running the server on two ports—which essentially means running two servers on the same machine—does not have a huge impact on the hardware.

## Approach #4

In a real production scenario, you may have millions of concurrent users simultaneously hitting the system. But in a load testing scenario, these users are to be artificially generated by a client running on a machine.

Here again the 65k port limit comes to bite on the client side. The only way to overcome this from the client's perspective is to increase the number of client machines that are generating the load. As you will read the next part in this series you will find that we had to use about 14 different machines to generate the kind of load we wanted to test HAProxy.

## Putting it all together

There isn't one single configuration that will solve all your woes and work like a charm. It is always the combination of multiple things that work out in the end.

For us as a prerequisite to load testing HAProxy, we followed approach #1 and approach #2 and eventually approach #3 to generate a huge... huge load of **2 million concurrent connections** on a single HAProxy machine.

[Here's the final part of this series](#), where I'll put together all the components that went into generating this kind of load, the tunings we did and the learnings that came out of it.

Do let me know how this blog post helped you and stay tuned for the final part in this series of posts. Also, please recommend (♥) this post if you think this may be useful for someone.



