



# DPDK & Layer 4 Packet Processing

## TLDK – Transport Layer Development Kit

M Jay

Presentation March 1 2017

# Legal Disclaimer

## General Disclaimer:

© Copyright 2017 Intel Corporation. All rights reserved. Intel, the Intel logo, Intel Inside, the Intel Inside logo, Intel. Experience What's Inside are trademarks of Intel Corporation in the U.S. and/or other countries. \*Other names and brands may be claimed as the property of others.

## Technology Disclaimer:

Intel technologies' features and benefits depend on system configuration and may require enabled hardware, software or service activation. Performance varies depending on system configuration. No computer system can be absolutely secure. Check with your system manufacturer or retailer or learn more at [intel.com].

## Performance Disclaimers (include only the relevant ones):

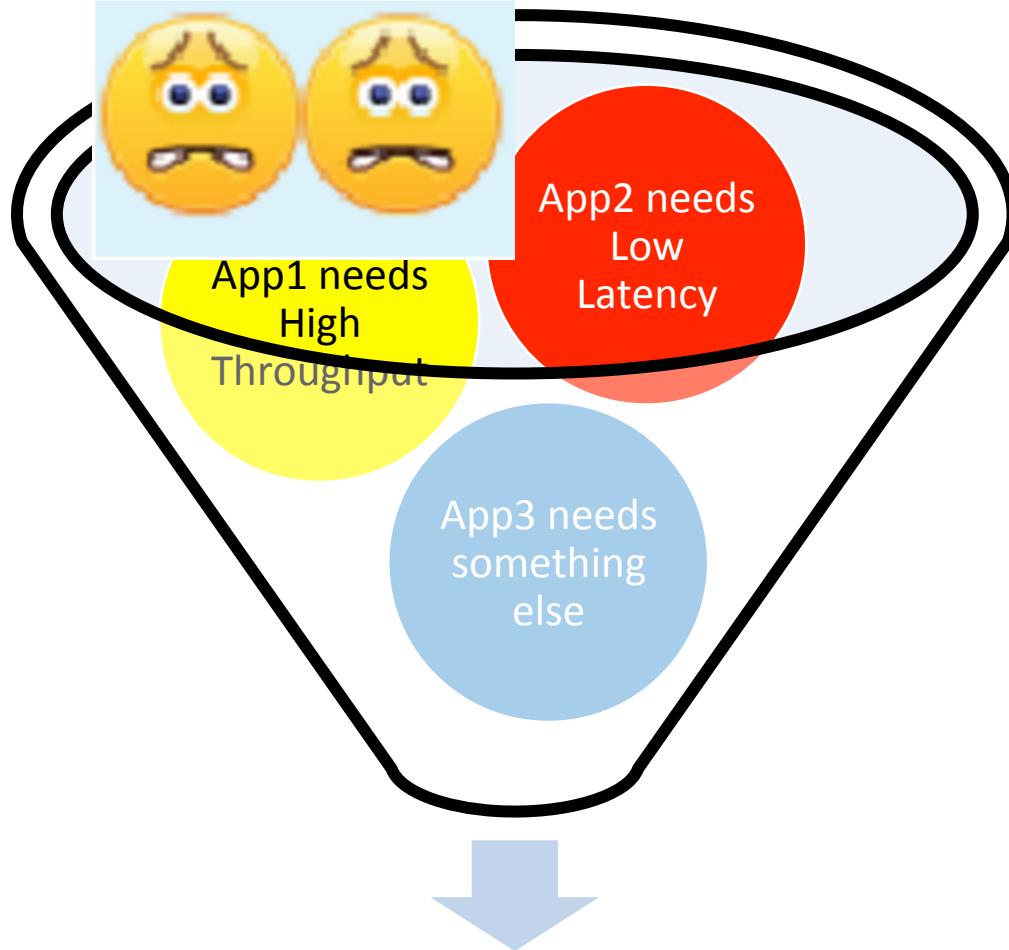
Cost reduction scenarios described are intended as examples of how a given Intel- based product, in the specified circumstances and configurations, may affect future costs and provide cost savings. Circumstances will vary. Intel does not guarantee any costs or cost reduction.

Results have been estimated or simulated using internal Intel analysis or architecture simulation or modeling, and provided to you for informational purposes. Any differences in your system hardware, software or configuration may affect your actual performance.



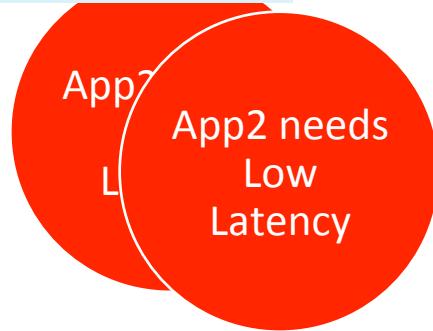
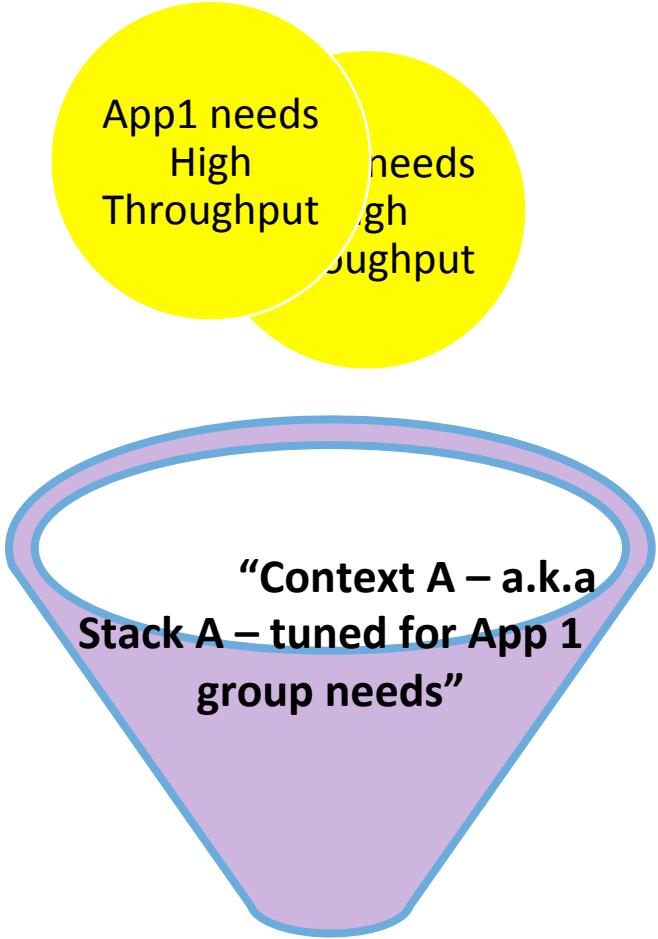
# What are the problems with traditional TCP/UDP stacks?

# One size “Kernel Stack” does not Fit All App needs

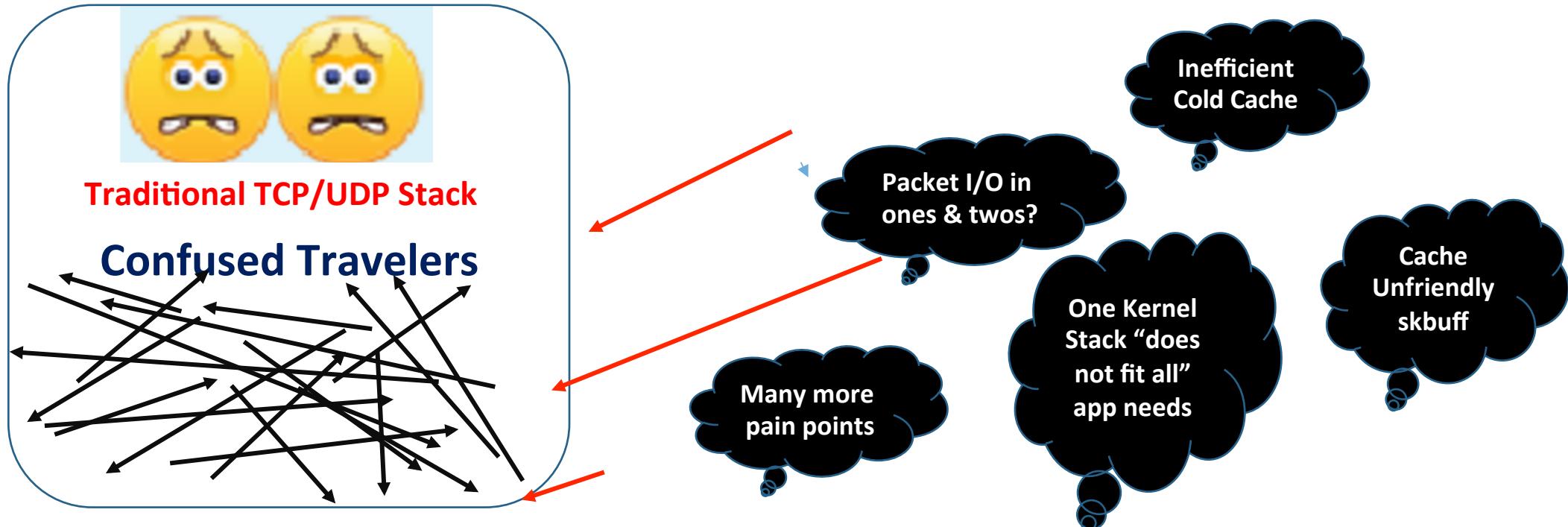


“one size does not fit all”  
Single Kernel Stack

# What if you can have multiple “user level” stacks?



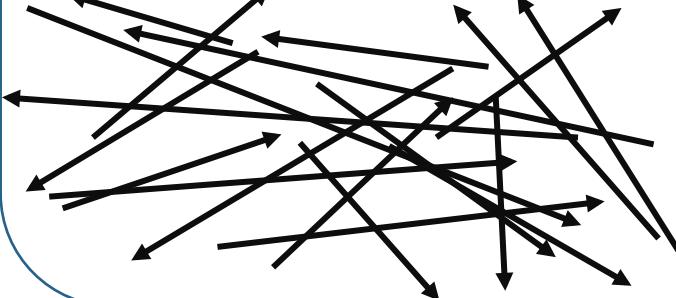
# What are the problems with traditional TCP/UDP stacks?





## Traditional TCP/UDP Stack

### Confused Travelers



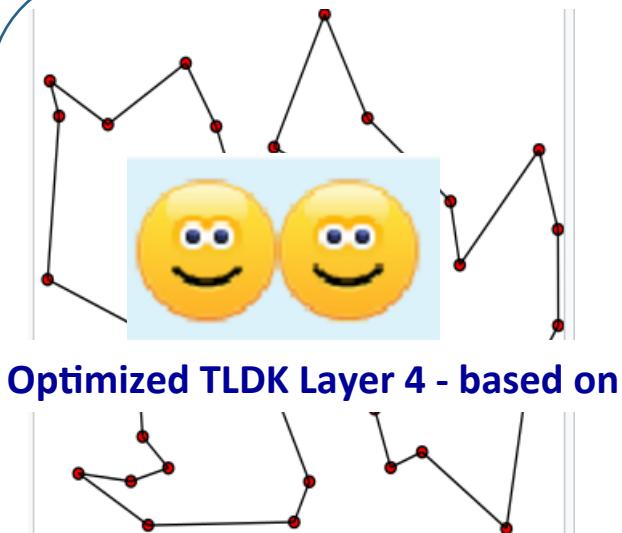
Packet I/O in  
ones & twos?

Inefficient  
Cold Cache

Many more  
pain points

One Kernel  
Stack “does  
not fit all”  
app needs

Cache  
Unfriendly  
skbuff



[https://en.wikipedia.org/wiki/Travelling\\_salesman\\_problem](https://en.wikipedia.org/wiki/Travelling_salesman_problem)

Warm I Cache  
& D Cache –  
Pull packets  
only when  
ready to  
process

User Level  
TLDK – no  
constraint of  
one stack

Cache Friendly  
mbuf,  
Dynamic Ring

Bulk  
Transfer

# TLDK – Any Sample Applications?

You know L2fwd, L3fwd

- What about **udpfwd?**

The Eternal Loop of /tldk/examples/udpfwd/  
main.c

\tldk\examples\udpfwd\main.c

```
/* launch all slave lcores. */
RTE_LCORE_FOREACH_SLAVE(i) {
    if (prm[i].be.lc != NULL || prm[i].fe.max_streams != 0)
        rte_eal_remote_launch(lcore_main, prm + i, i);
}

/* launch master lcore. */
i = rte_get_master_lcore();
if (prm[i].be.lc != NULL || prm[i].fe.max_streams != 0)
    lcore_main(prm + i);
```

```
static int
lcore_main(void *arg)
{
    int32_t rc;
    uint32_t lcore;
    struct lcore_prm *prm;

    prm = arg;
    lcore = rte_lcore_id();

    RTE_LOG(NOTICE, USER1, "%s(lcore=%u) start\n",
           __func__, lcore);

    rc = 0;

    /* lcore FE init. */
    if (prm->fe.max_streams != 0)
        rc = netfe_lcore_init(&prm->fe);

    /* lcore BE init. */
    if (rc == 0 && prm->be.lc != NULL)
        rc = netbe_lcore_setup(prm->be.lc);

    if (rc != 0)
        sig_handle(SIGQUIT);

    while (force_quit == 0) {
        netfe_lcore();
        netbe_lcore();
    }
}
```

Eternal Loop of  
/tldk/examples/udpfwd/main.c



# What Are The Main APIs ?

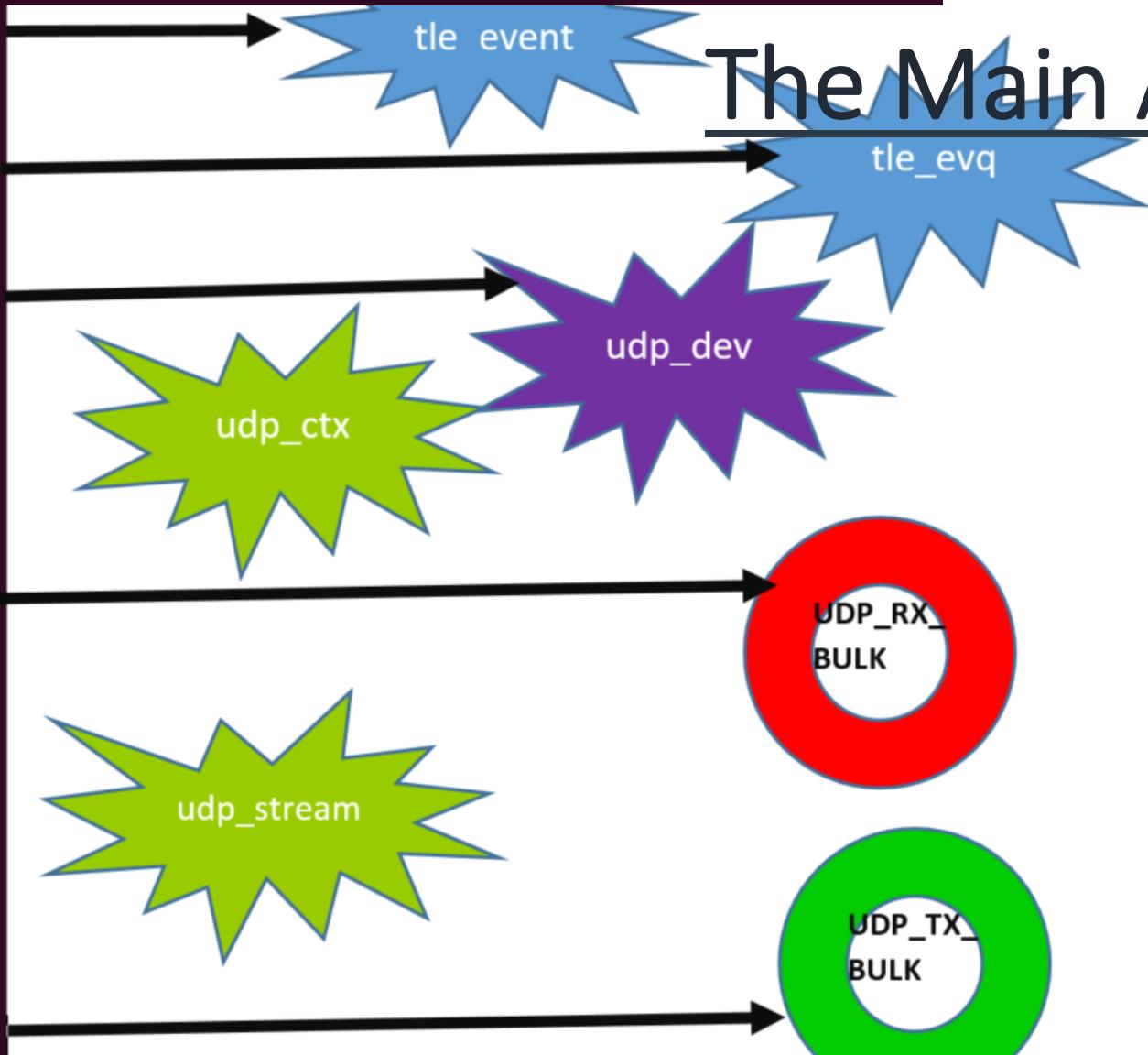
# 3 Key Modules

- 1) event.o
- 2) Udp\_ctl.o
- 3) Udp\_rxtx.o

```
0x000000000004302b0 0x3983 main.o
0x00000000000433c33 0x3a3 /home/user/tldk/x86_64-native-linuxapp-gcc/lib/libtle_udp.a(event.o)
0x00000000000433c7c tle_evq_create
0x00000000000433e22 tle_evq_destroy
0x00000000000433e3d tle_event_alloc
0x00000000000433f0d tle_event_free
0x00000000000433fd6 0x5796 /home/user/tldk/x86_64-native-linuxapp-gcc/lib/libtle_udp.a(udp_ctl.o)
0x00000000000436bde tle_udp_create
0x00000000000436e37 tle_udp_destroy
0x00000000000436ef5 tle_udp_ctx_invalidate
0x00000000000437096 tle_udp_add_dev
0x00000000000437e72 tle_udp_del_dev
0x000000000004388a4 tle_udp_stream_open
0x00000000000438b93 tle_udp_stream_close
0x000000000004396c3 tle_udp_stream_get_param
0x0000000000043976c 0x7be1 /home/user/tldk/x86_64-native-linuxapp-gcc/lib/libtle_udp.a(udp_rxtx.o)
0x0000000000043ce42 tle_udp_rx_bulk
0x0000000000043dd33 tle_udp_tx_bulk
0x0000000000043e380 tle_udp_stream_recv
0x00000000000440c3c tle_udp_stream_send
```

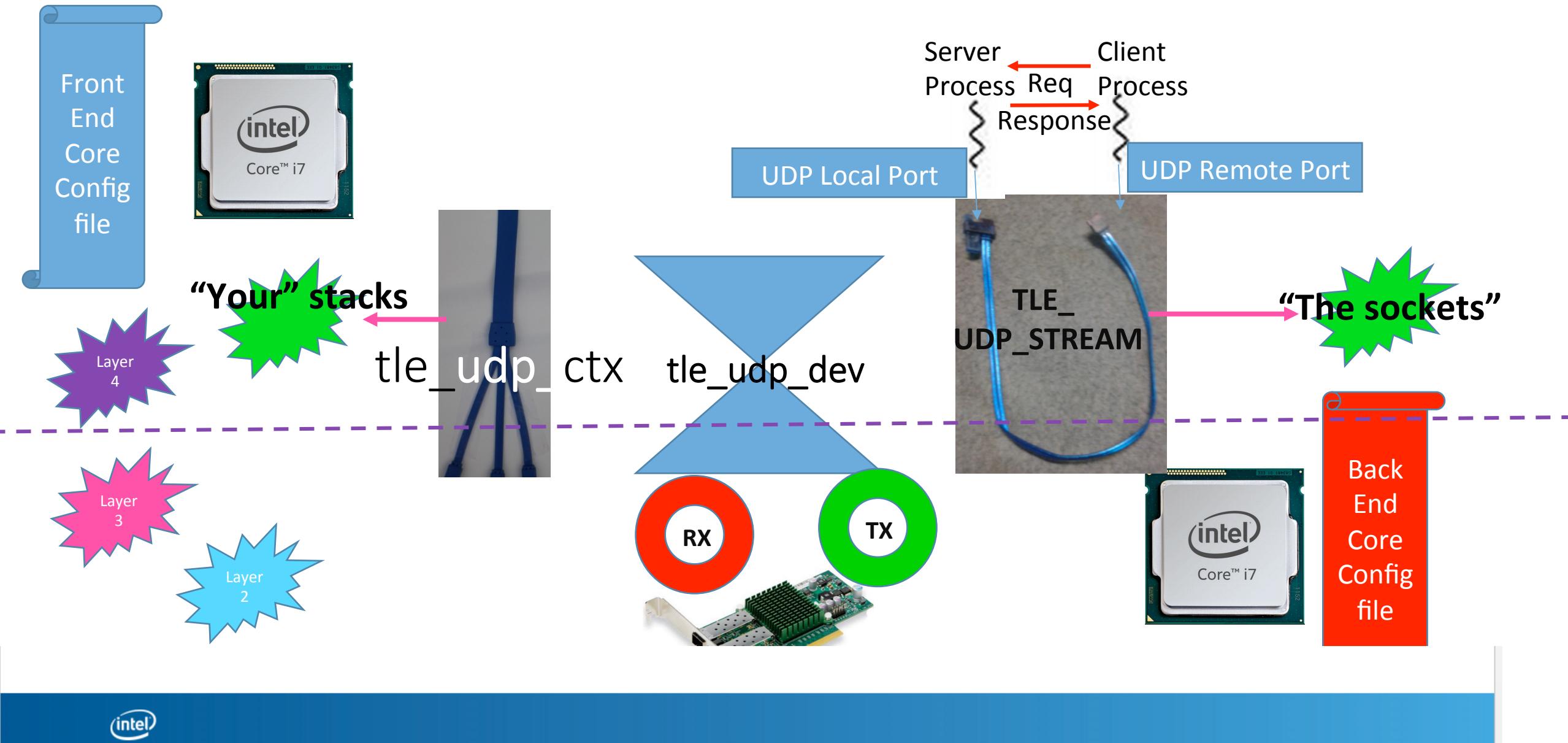
/tldk/x86\_64-native-linuxapp-gcc/lib/libtle\_udp.a(event.o)

- tle\_event\_free
- tle\_evq\_create
- tle\_evq\_destroy
- tle\_udp\_add\_dev
- tle\_udp\_create
- tle\_udp\_ctx\_invalidate
- tle\_udp\_del\_dev
- tle\_udp\_destroy
- tle\_udp\_rx\_bulk
- tle\_udp\_stream\_close
- tle\_udp\_stream\_get\_param
- tle\_udp\_stream\_open
- tle\_udp\_stream\_recv
- tle\_udp\_stream\_send
- tle\_udp\_tx\_bulk

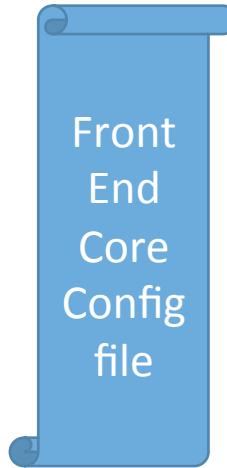


## The Main APIs

# Front End Core + Back End Core Pulling UDP Stack



# What does FE config file configure?



```
lcore=3,  
op=fwd,  
laddr=0.0.0.0,  
lport=11211,  
raddr=0.0.0.0,  
rport=0,  
fwladdr=:,:  
fwlport=0,  
fwraddr=2001:4860:b002::56,  
fwrport=11211
```

- 1) Read [Server]
- 2) Write [Client]
- 3) Forward [Gateway]

## FE config record format:

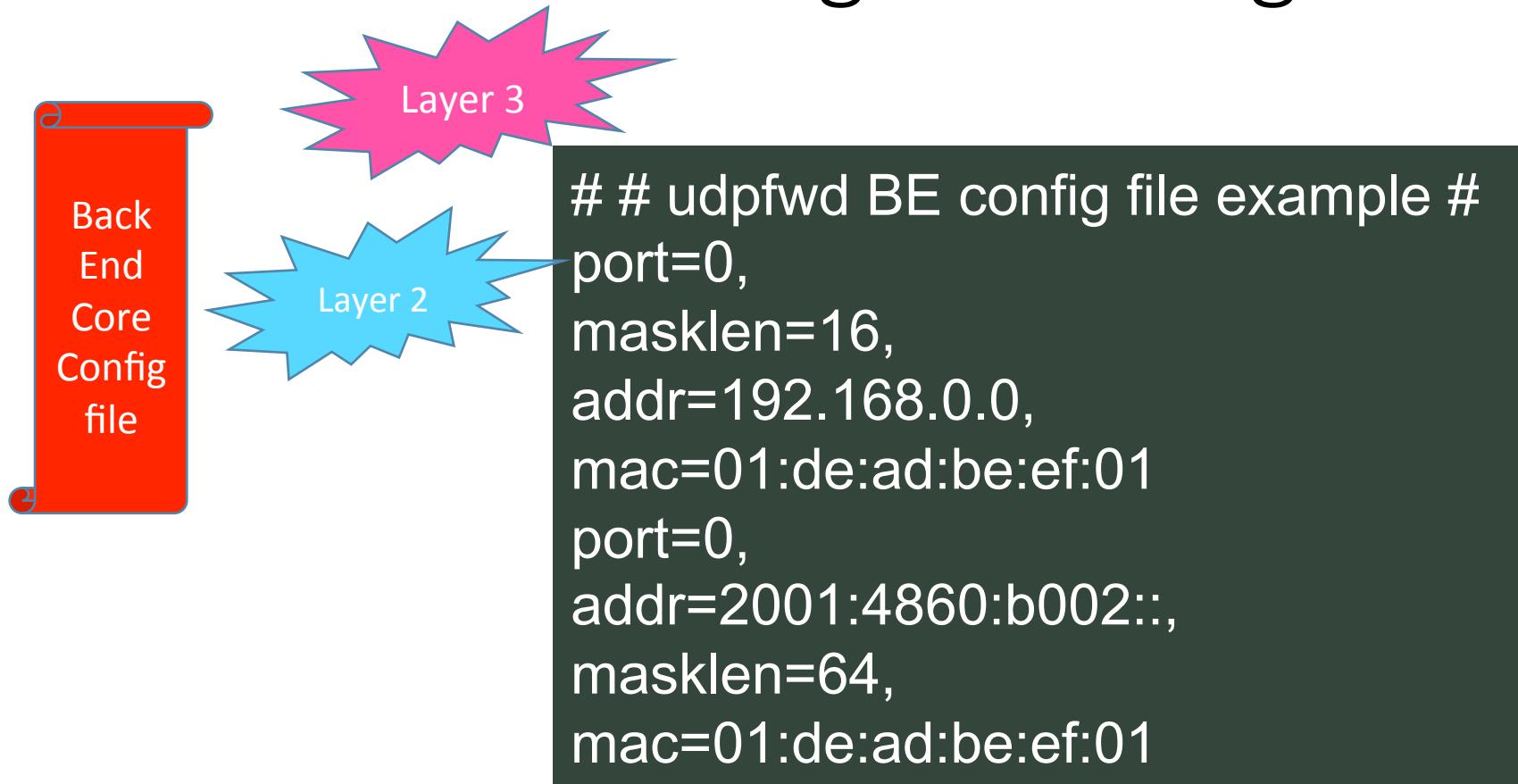
---

```
Icore=<uint>,op=<"rx|tx|echo|fwd">,\  
laddr=<ip>,lport=<uint16>,raddr=<ip>,rport=<uint16>,\  
[txlen=<uint>,fwladdr=<ip>,fwlport=<uint16>,fwraddr=<ip>,fwrport  
=<uint16>,\  
belcore=<uint>]
```

Icore - EAL Icore to manage that stream(s) in the FE.  
op - operation to perform on that stream:  
    "rx" - do receive only on that stream.  
    "tx" - do send only on that stream.  
    "echo" - mimic recvfrom(..., &addr);sendto(..., &addr);  
          on that stream.  
    "fwd" - forward packets between streams.

laddr - local address for the stream to open.  
lport - local port for the stream to open.  
raddr - remote address for the stream to open.  
rport - remote port for the stream to open.  
txlen - data length to send with each packet ("tx" mode only).  
fwladdr - local address for the forwarding stream(s) to open  
("fwd mode only).  
fwlport - local port for the forwarding stream(s) to open  
("fwd mode only).  
fwraddr - remote address for the forwarding stream(s) to open  
("fwd mode only).  
fwrport - remote port for the forwarding stream(s) to open  
("fwd mode only).  
belcore - EAL Icore to manage that stream(s) in the BE.

# What does BE config file configure?



In addition, can I do command line config. as well?

# Yes ! Command line Options configuring

```
As an example: udpfwd  
--lcores='3,6,8'  
-w 01:00.0 -- \ --promisc  
--rbufs 0x1000  
--sbufs 0x1000  
--streams 0x100  
\ --fecfg ./fe.cfg  
--becfg ./be.cfg  
\ port=0,  
lcore=6,  
lcore=8,  
rx_offload=0xf,  
tx_offload=0,\  
 ipv4=192.168.1.233,  
ipv6=2001:4860:b002::28
```

Will create TLDK UDP context on

lcore=6 and lcore=8 (BE lcore) to manage DPDK port 0.

Will assign IPv4 address 192.168.1.233 and  
IPv6 address 2001:4860:b002::28 to that port.

The following supported by

DPDK RX HW offloads:

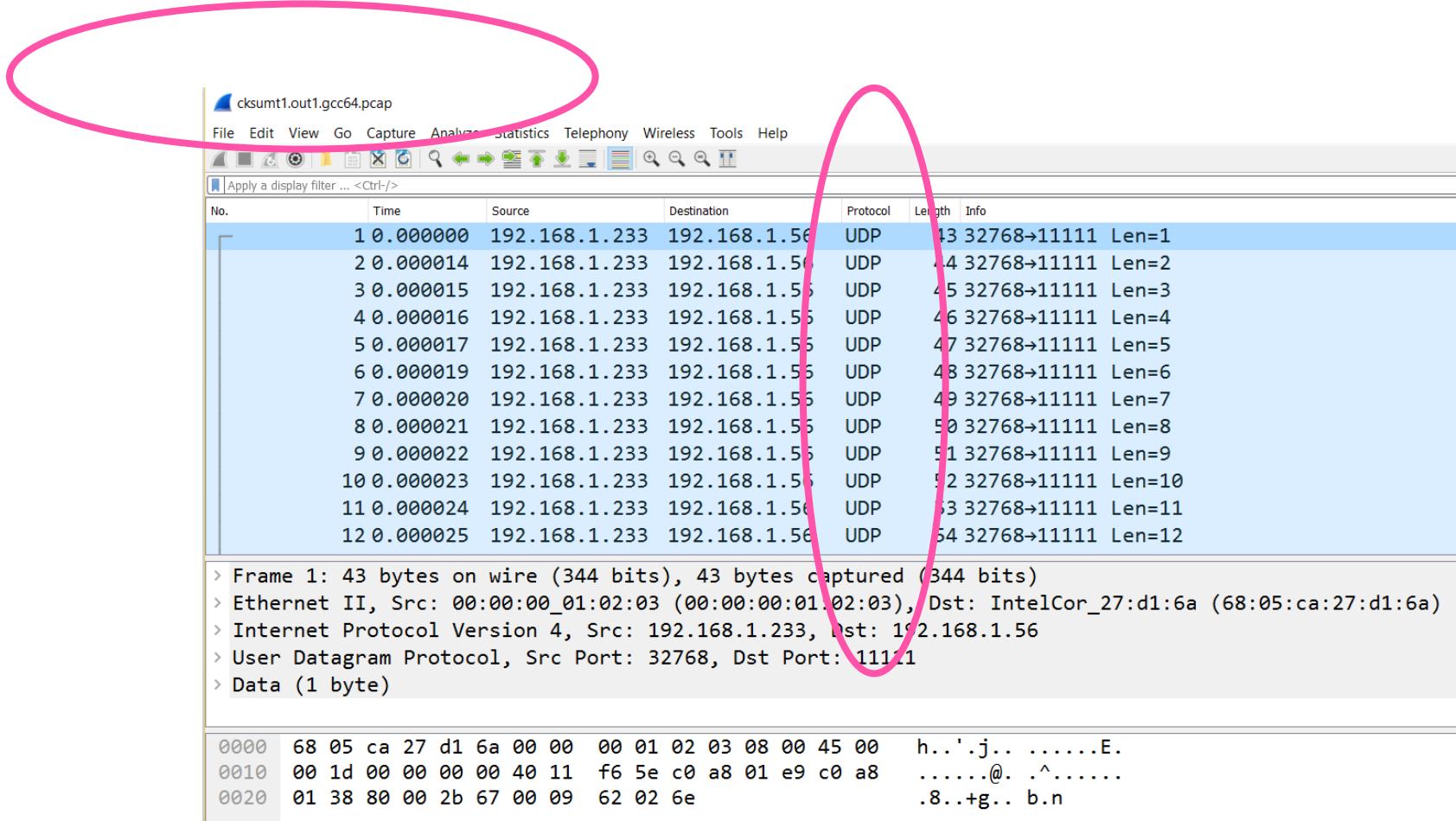
- DEV\_RX\_OFFLOAD\_VLAN\_STRIP,
- DEV\_RX\_OFFLOAD\_IPV4\_CKSUM,
- DEV\_RX\_OFFLOAD\_UDP\_CKSUM,
- DEV\_RX\_OFFLOAD\_

TCP\_CKSUM will be enabled on that port. No HW TX offloads will be enabled.

# Have you run udpfwd in echo mode?

- Do you have output – we can look at?

# Wireshark Output of the file created by udpfwd – echo mode



The screenshot shows the Wireshark interface displaying a capture named "cksum1.out1.gcc64.pcap". The main pane shows a list of 12 UDP frames. The first frame is highlighted with a pink oval. A second pink oval highlights the detailed information pane below the list, which provides a summary of Frame 1 and its protocol headers.

No.	Time	Source	Destination	Protocol	Length	Info
1	0.000000	192.168.1.233	192.168.1.56	UDP	43	32768→11111 Len=1
2	0.000014	192.168.1.233	192.168.1.56	UDP	44	32768→11111 Len=2
3	0.000015	192.168.1.233	192.168.1.56	UDP	45	32768→11111 Len=3
4	0.000016	192.168.1.233	192.168.1.56	UDP	46	32768→11111 Len=4
5	0.000017	192.168.1.233	192.168.1.56	UDP	47	32768→11111 Len=5
6	0.000019	192.168.1.233	192.168.1.56	UDP	48	32768→11111 Len=6
7	0.000020	192.168.1.233	192.168.1.56	UDP	49	32768→11111 Len=7
8	0.000021	192.168.1.233	192.168.1.56	UDP	50	32768→11111 Len=8
9	0.000022	192.168.1.233	192.168.1.56	UDP	51	32768→11111 Len=9
10	0.000023	192.168.1.233	192.168.1.56	UDP	52	32768→11111 Len=10
11	0.000024	192.168.1.233	192.168.1.56	UDP	53	32768→11111 Len=11
12	0.000025	192.168.1.233	192.168.1.56	UDP	54	32768→11111 Len=12

> Frame 1: 43 bytes on wire (344 bits), 43 bytes captured (344 bits)  
> Ethernet II, Src: 00:00:00\_01:02:03 (00:00:00:01:02:03), Dst: IntelCor\_27:d1:6a (68:05:ca:27:d1:6a)  
> Internet Protocol Version 4, Src: 192.168.1.233, Dst: 192.168.1.56  
> User Datagram Protocol, Src Port: 32768, Dst Port: 11111  
> Data (1 byte)

Hex	Dec	ASCII
0000	68 05 ca 27 d1 6a 00 00	h...'.j.. .....
0010	00 1d 00 00 00 40 11 f6 5e c0 a8 01 e9 c0 a8	.....@. .^.....
0020	01 38 80 00 2b 67 00 09 62 02 6e	.8..+g.. b.n

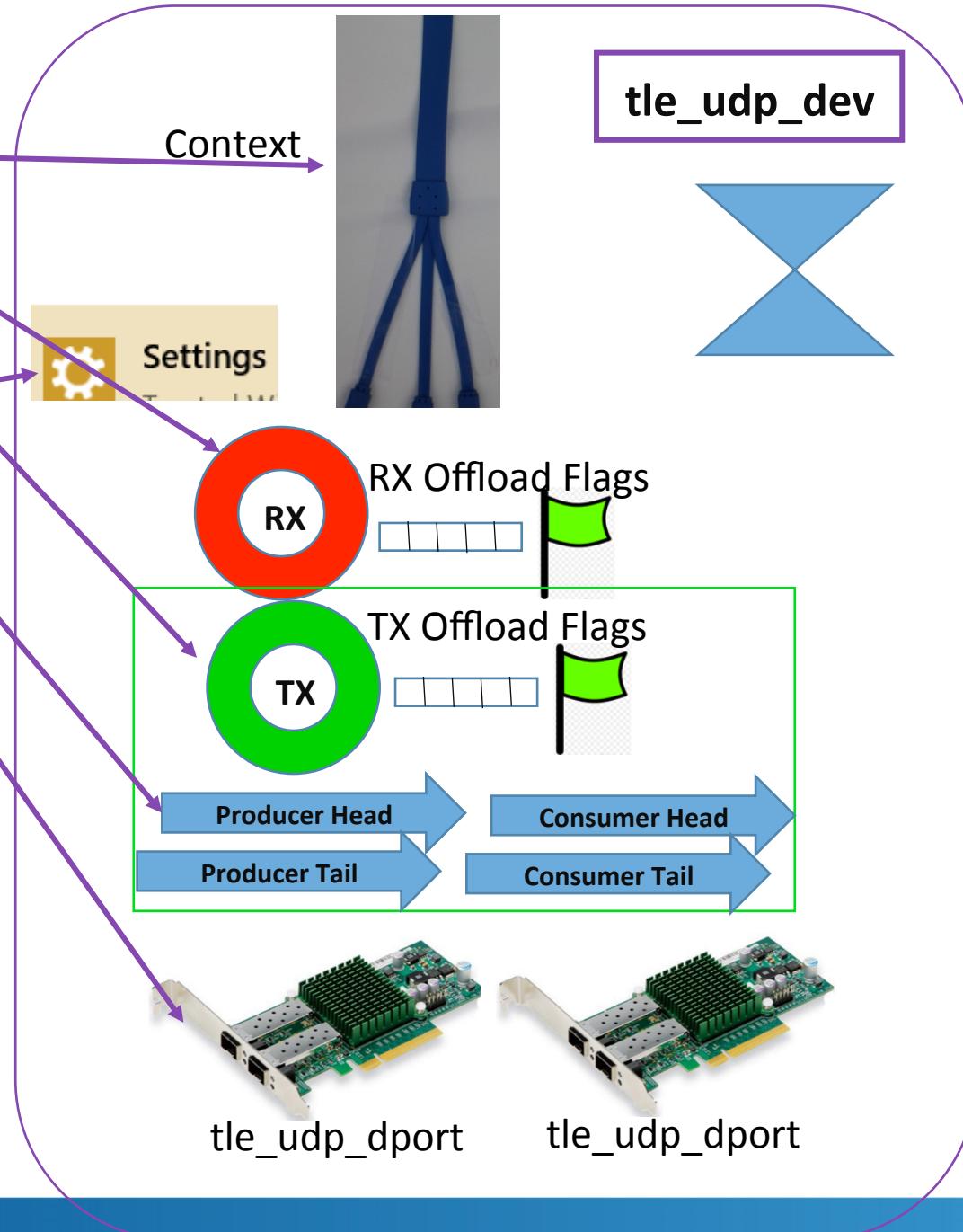
# tle\_udp\_dev



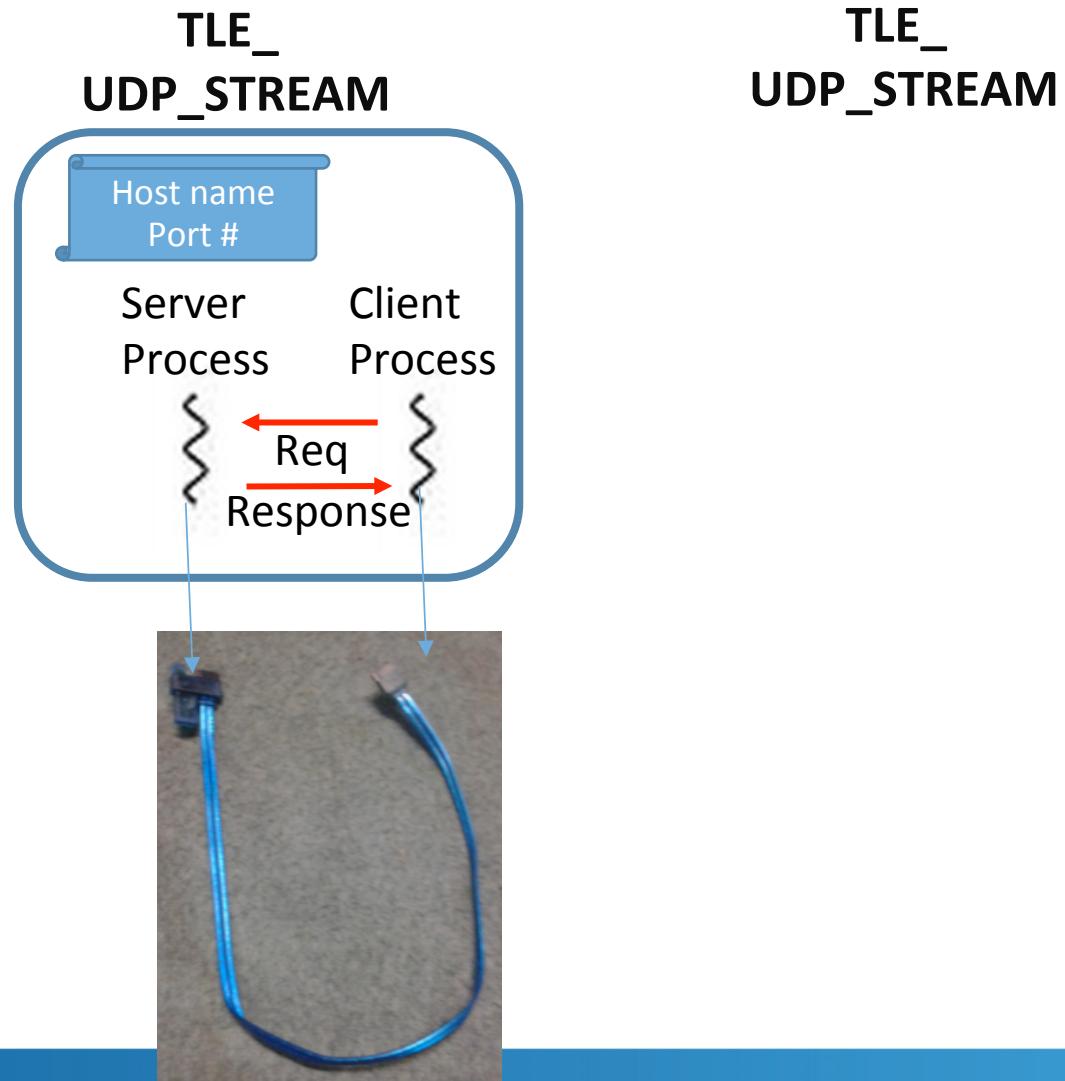
What is tle\_udp\_dev?

What it contains?

```
(gdb) ptype dev
type = struct tle_udp_dev {
    struct tle_udp_ctx *ctx;
    struct {
        uint64_t ol_flags[2];
    } rx;
    struct {
        uint64_t ol_flags[2];
        rte_atomic32_t packet_id[2];
        struct tle_dring dr;
    } tx;
    struct tle_udp_dev_param prm;
    struct tle_udp_dport *dp[2];
} *
```



# What is TLE\_UDP\_STREAM?



# What is TLE\_UDP\_STREAM?



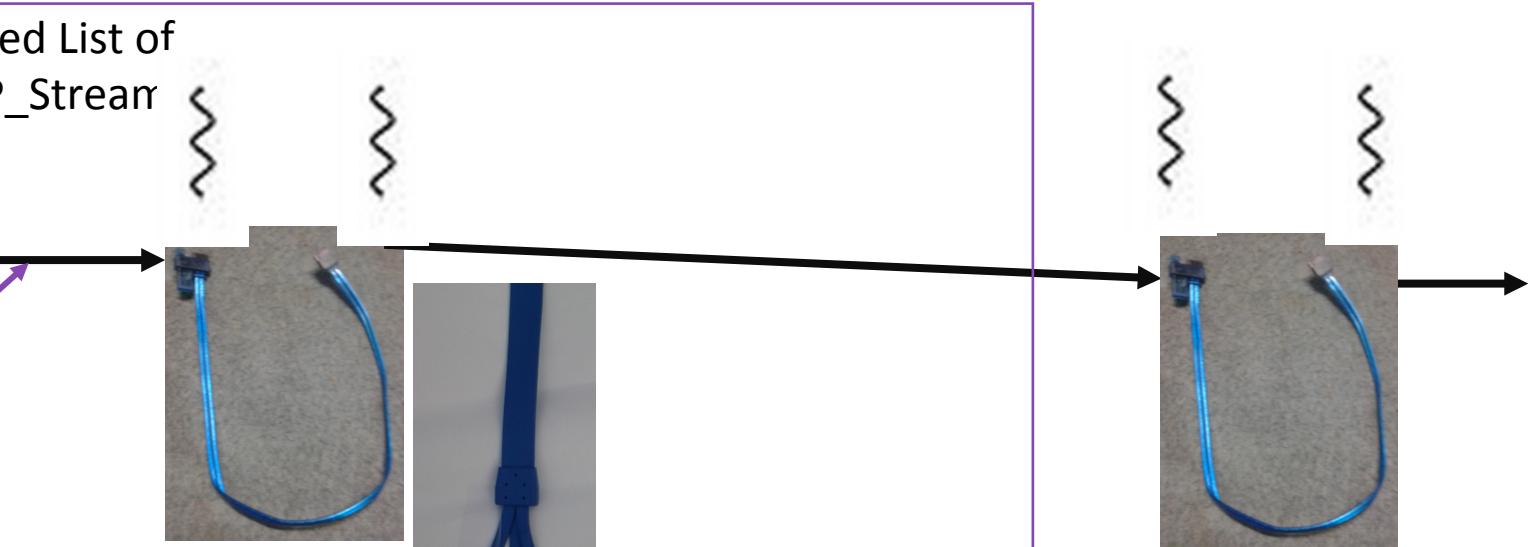
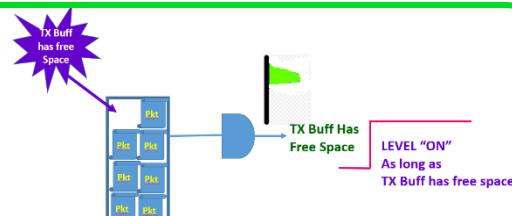
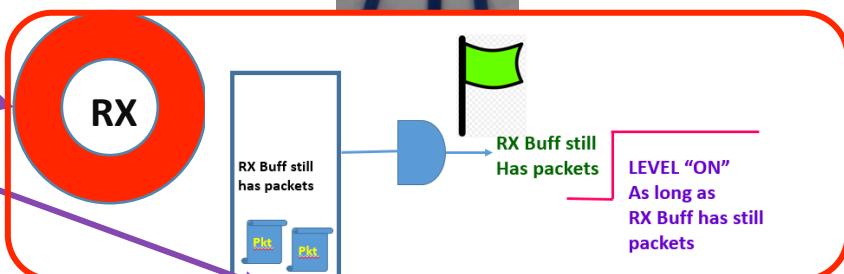
```
(gdb) ptype s
type = struct tle_udp_stream {
    struct {
        struct tle_udp_stream *stqe_next;
    } link;
    struct tle_udp_ctx *ctx;
    uint8_t type;
    struct {
        struct rte_ring *q;
        struct tle_event *ev;
        struct tle_udp_stream_cb cb;
        rte_atomic32_t use;
    } rx;
    union udp_ports port;
    union udp_ports pmsk;
    union {
        struct {...} ipv4;
        struct {...} ipv6;
    };
    struct {
        rte_atomic32_t use;
        struct {...} drb;
        struct tle_event *ev;
        struct tle_udp_stream_cb cb;
    } tx;
    struct tle_udp_stream_param prm;
}
```

Linked List of UDP\_Stream

Context

Settings

DNS – 53 DHCP – 67 ...

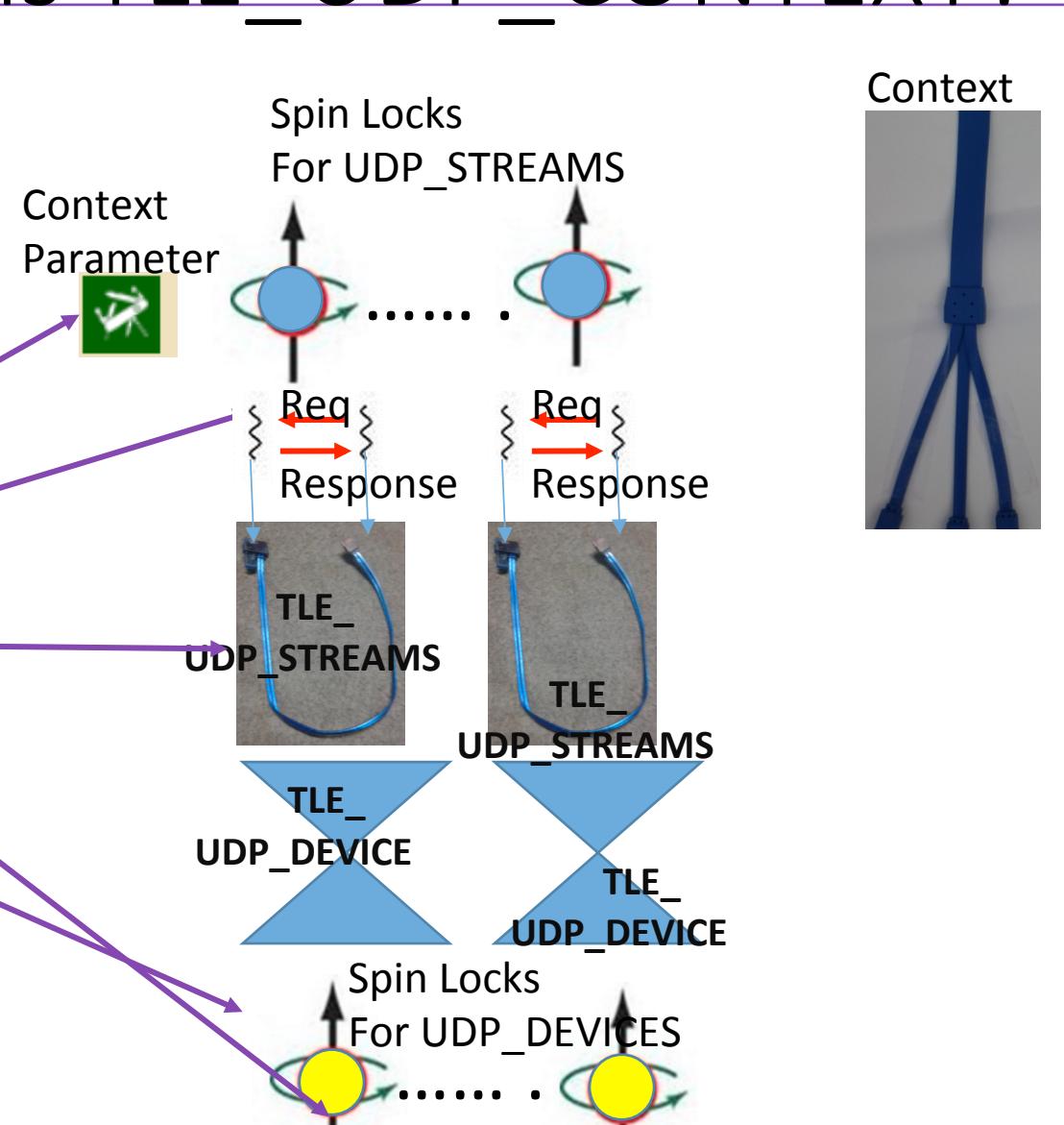


# What is TLE\_UDP\_CONTEXT?



Context

```
(gdb) ptype ctx
type = struct tle_udp_ctx {
    struct tle_udp_ctx_param prm;
    struct {
        rte_spinlock_t lock;
        uint32_t nb_free;
        struct {...} free;
        struct tle_udp_stream *buf;
    } streams;
    rte_spinlock_t dev_lock;
    uint32_t nb_dev;
    struct udp_pbm use[2];
    struct tle_udp_dev dev[32];
} *
```

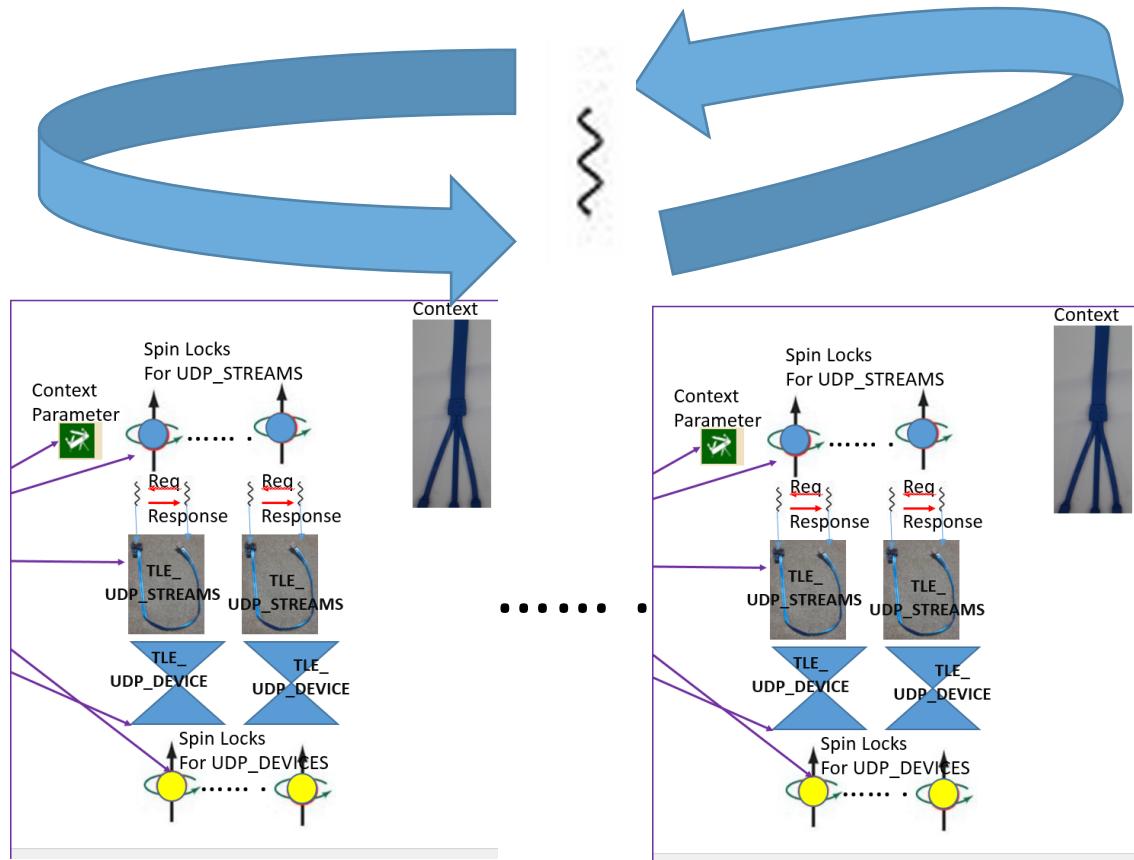


Context

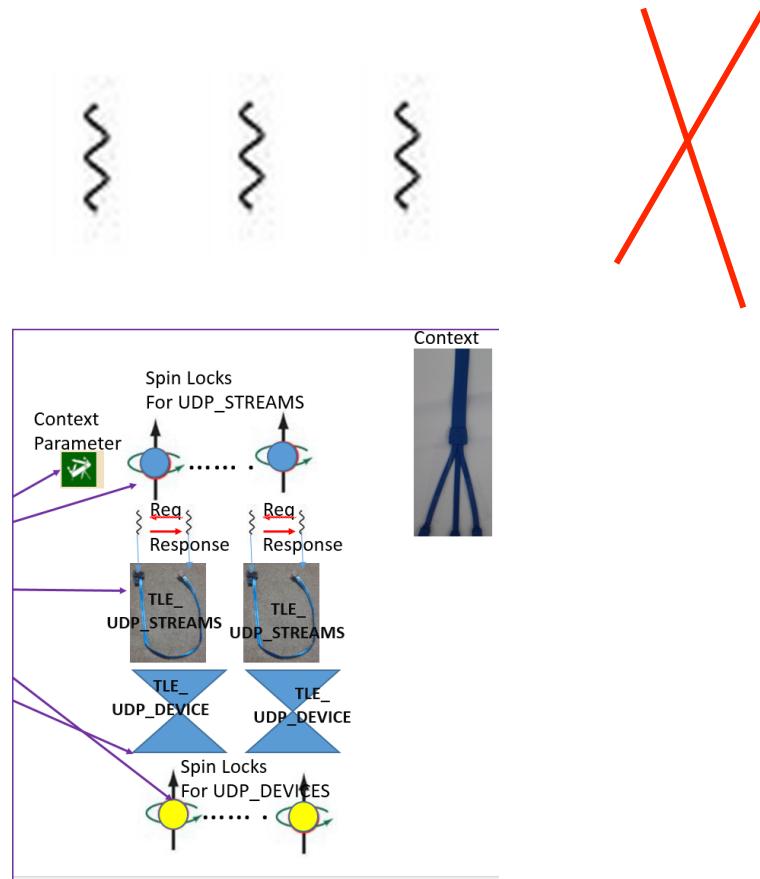


# Developer's Responsibility – Scope

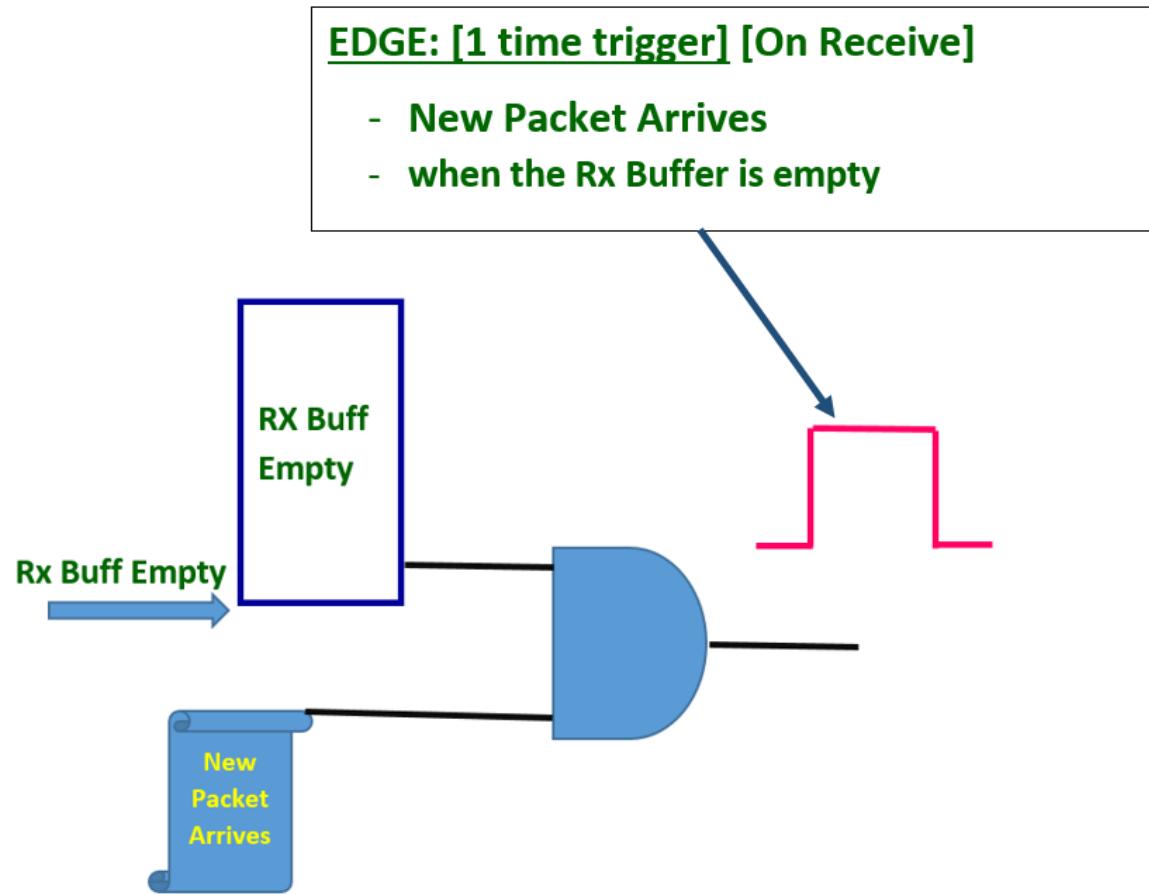
One Thread Driving  
Multiple Contexts Allowed As Is



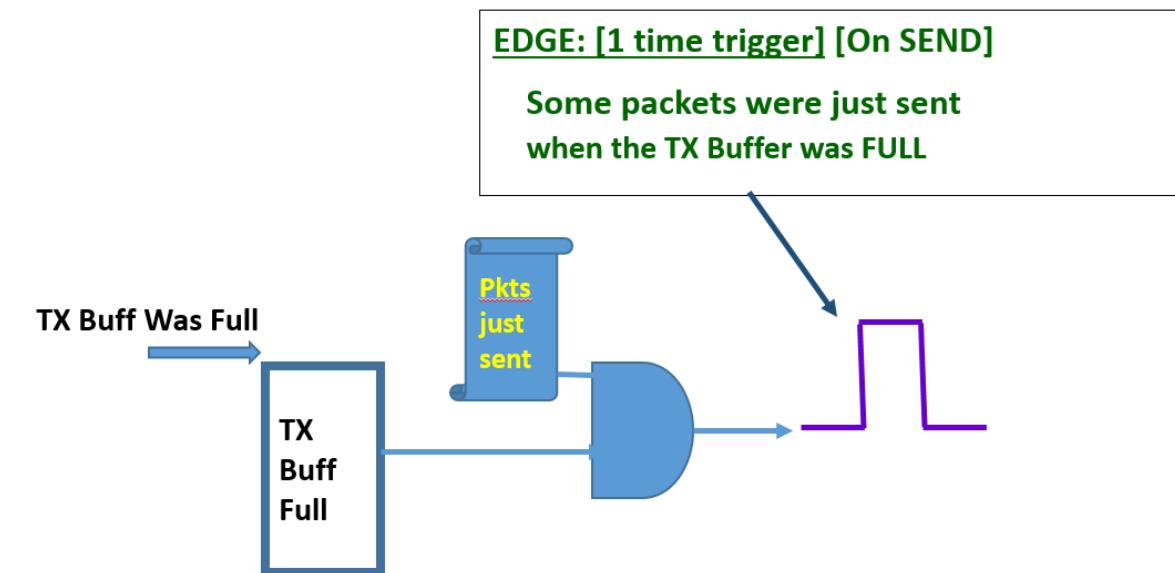
Multiple Threads Driving Single Context  
NOT Allowed As Is. You need to provide  
Explicit synchronization



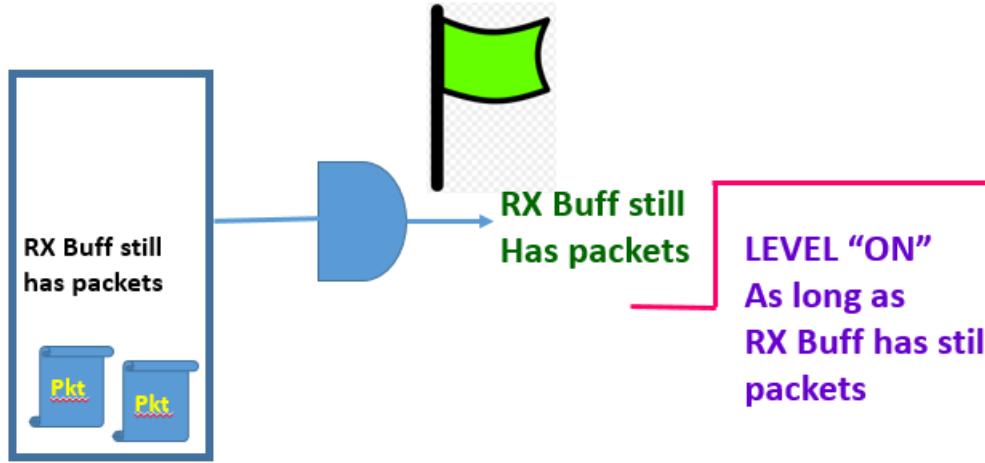
# Callback Model for RECEIVE



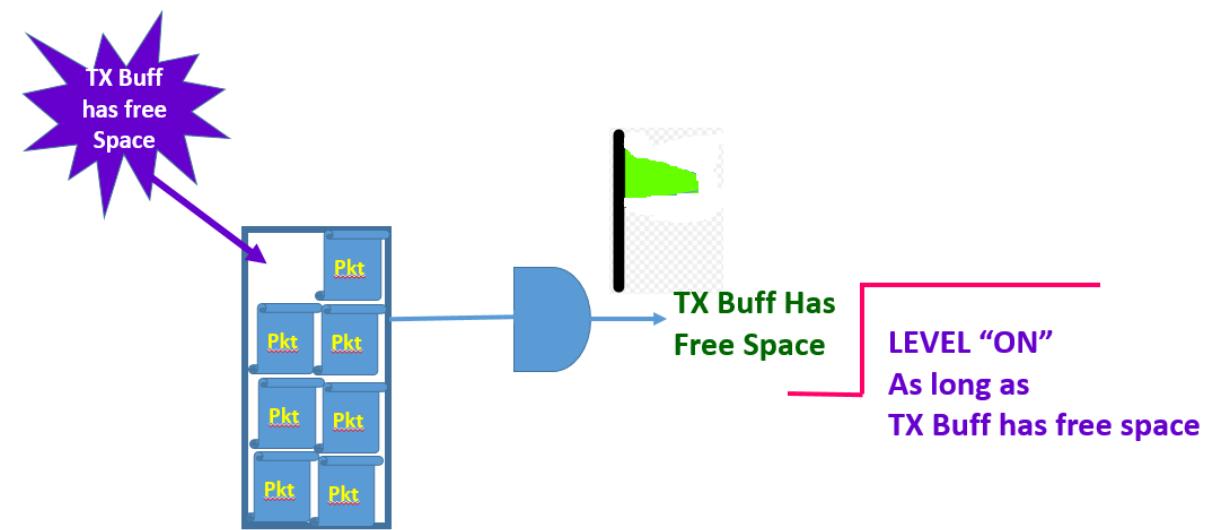
# Callback Model for TRANSMIT



# Event Model for RX



# Event Model for TX



# Can I use Callback for RX & Event notification for TX?

- Vice versa?

# Back Up

# Event/Callback & RX/TX - Combinations

Notification Choice	Event Notify Set for Send	Callback Set for Send	No Notification Set for Send
<b>Event Notify set for Rcv</b>	For both receive & Send, you get LEVEL TRIGGERED notification	You get LEVEL TRIGGERED For recv & EDGE TRIGGERED notification for Send	Only Receive gets notified – LEVEL TRIGGERED. No Notification for Send
<b>Callback set for Rcv</b>	You get EDGE TRIGGERED For recv & LEVEL TRIGGERED notification for Send	For both receive & Send, you get EDGE TRIGGERED notification	Only Receive gets notified – EDGE TRIGGERED. No Notification for Send
<b>No Notification set for Rcv</b>	Only Send gets notified – LEVEL TRIGGERED. No Notification for Rcv	Only Send gets notified – EDGE TRIGGERED. No Notification for Rcv	No Notification for both send & Receive

# Key gdb commands

ptype x	Data Structure Type of struct x
b <function name>	Set breakpoint
clear <function name>	Clear breakpoint @ function
info locals	
gdb -- arg ./<function><args>	Debugging
run	
continue	

# Key virsh & ssh commands

sudo virsh domifaddr <Host Name>		
Sudo virsh list --all		
ssh user@192.168....		
ip r add default via 192.168... dev ens8	After reboot of VMs, to resolve any connectivity issue	

# Key git & Misc commands

git log	Returns 40 bit IDs	
git checkout <40 bit ID>		
sed -i 's/\r\$//' filename	Windows control characters issue	

# Key Objdump commands

Objdump - - syms <file.o>	To read the symbols	
Objdump -S <file.o>	Disassembly mixed with source + Assembly	Objdump -S <file.o>
Objdump -h <file.o>	To get details about file	Txt, data, bss, stack
Objdump -x <file.o>	For all headers in detail	

## Index of /~sgtatham/putty/0.60/x86

Name	Last modified	Size	Description
<a>Parent Directory</a>		-	
<a>pageant.exe</a>	2007-04-29 14:02	132K	
<a>pageant.exe.DSA</a>	2007-04-29 14:02	65	
<a>pageant.exe.RSA</a>	2007-04-29 14:02	152	
<a>plink.exe</a>	2007-04-29 14:02	276K	
<a>plink.exe.DSA</a>	2007-04-29 14:02	65	
<a>plink.exe.RSA</a>	2007-04-29 14:02	152	
<a>pscp.exe</a>	2007-04-29 14:02	288K	
<a>pscp.exe.DSA</a>	2007-04-29 14:02	65	
<a>pscp.exe.RSA</a>	2007-04-29 14:02	152	
<a>psftp.exe</a>	2007-04-29 14:02	300K	
<a>psftp.exe.DSA</a>	2007-04-29 14:02	65	
<a>psftp.exe.RSA</a>	2007-04-29 14:02	152	
<a>putty-0.60-installer.exe</a>	2007-04-29 14:02	1.7M	
<a>putty-0.60-installer.exe.DSA</a>	2007-04-29 14:02	65	
<a>putty-0.60-installer.exe.RSA</a>	2007-04-29 14:02	152	
<a>putty.exe</a>	2007-04-29 14:02	444K	
<a>putty.exe.DSA</a>	2007-04-29 14:02	65	
<a>putty.exe.RSA</a>	2007-04-29 14:02	152	
<a>putty.zip</a>	2007-04-29 14:02	1.4M	
<a>putty.zip.DSA</a>	2007-04-29 14:02	65	
<a>putty.zip.RSA</a>	2007-04-29 14:02	152	
<a>puttygen.exe</a>	2007-04-29 14:02	168K	

# Install puTTY

- Click the URL
- <http://the.earth.li/~sgtatham/putty/0.60/x86/>
- Double Click putty-0.60-installer.exe option
- You have downloaded on your Download folder
- Right Click and Run As Administrator
- You have installed PuTTY

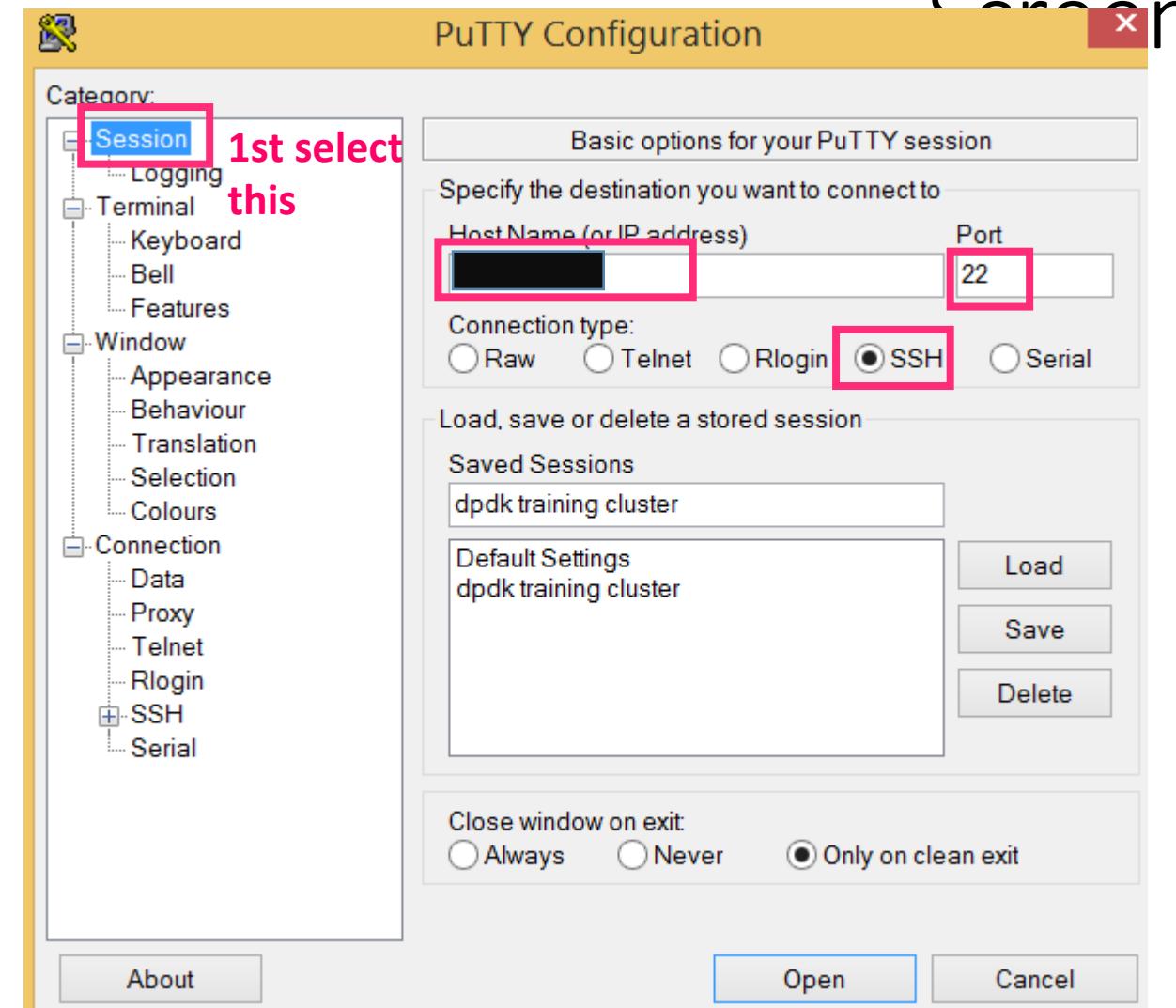


PuTTY Config: → Session  
1 / 2

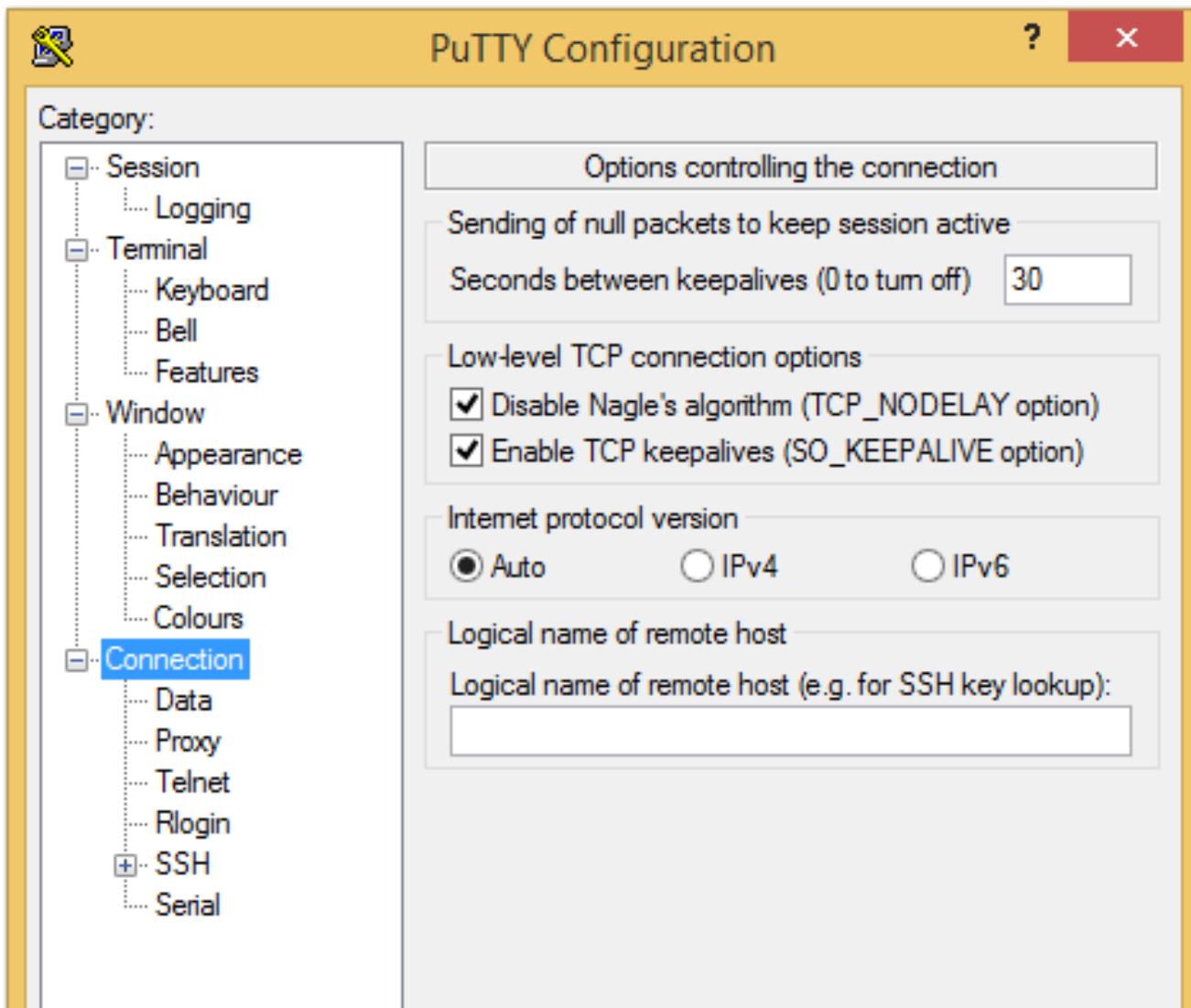
1) Host name

2) Port 22

3) Connection Type SSH



# Keep Alive – 30 sec; Enable TCP Keepalives



If you are outside

# PuTTY Config: Connection --> Proxy

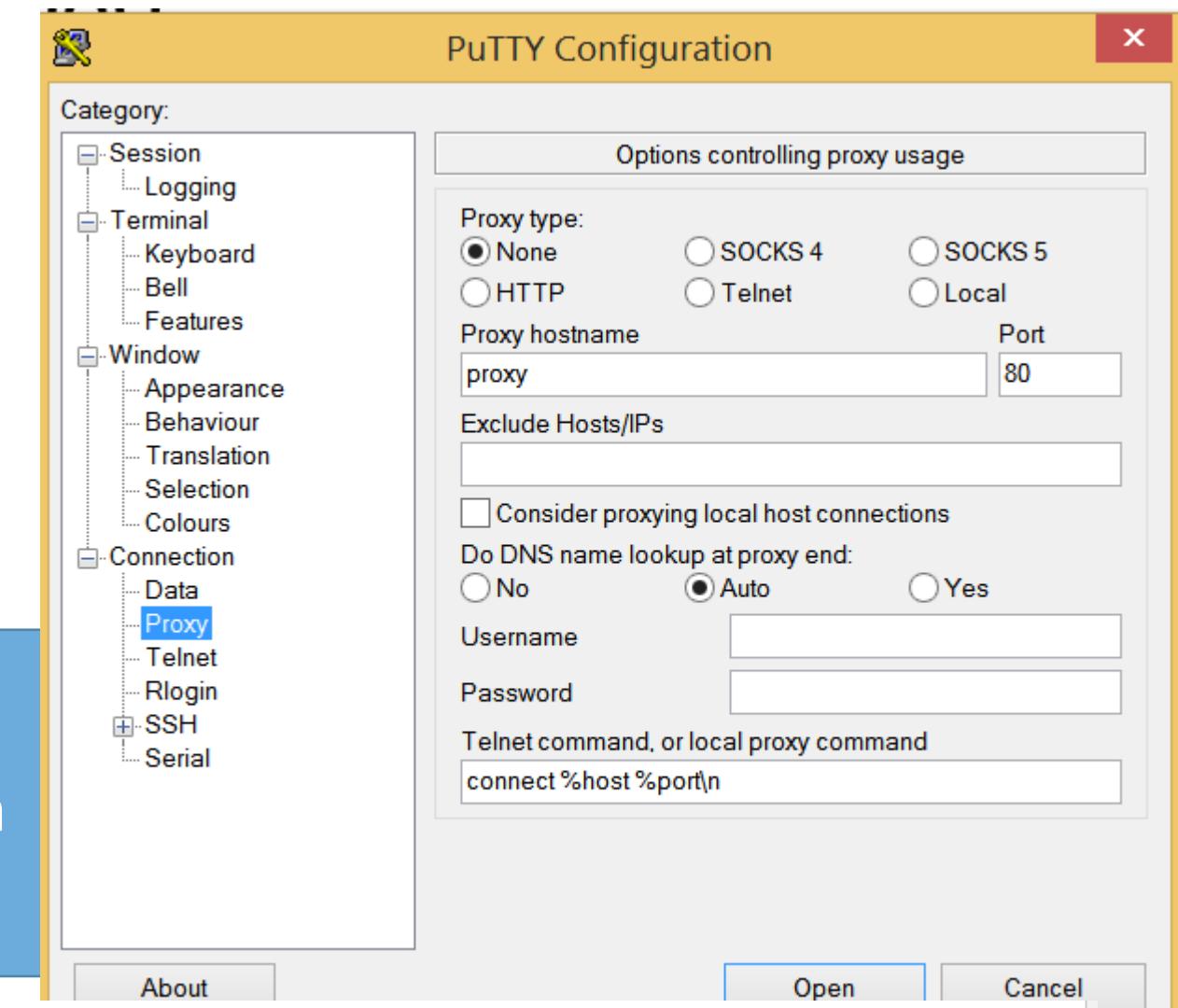
Screen 2 / 2

1) Proxy Type: None

2) Proxy host name: Leave it blank

3) Do DNS name lookup at proxy end:

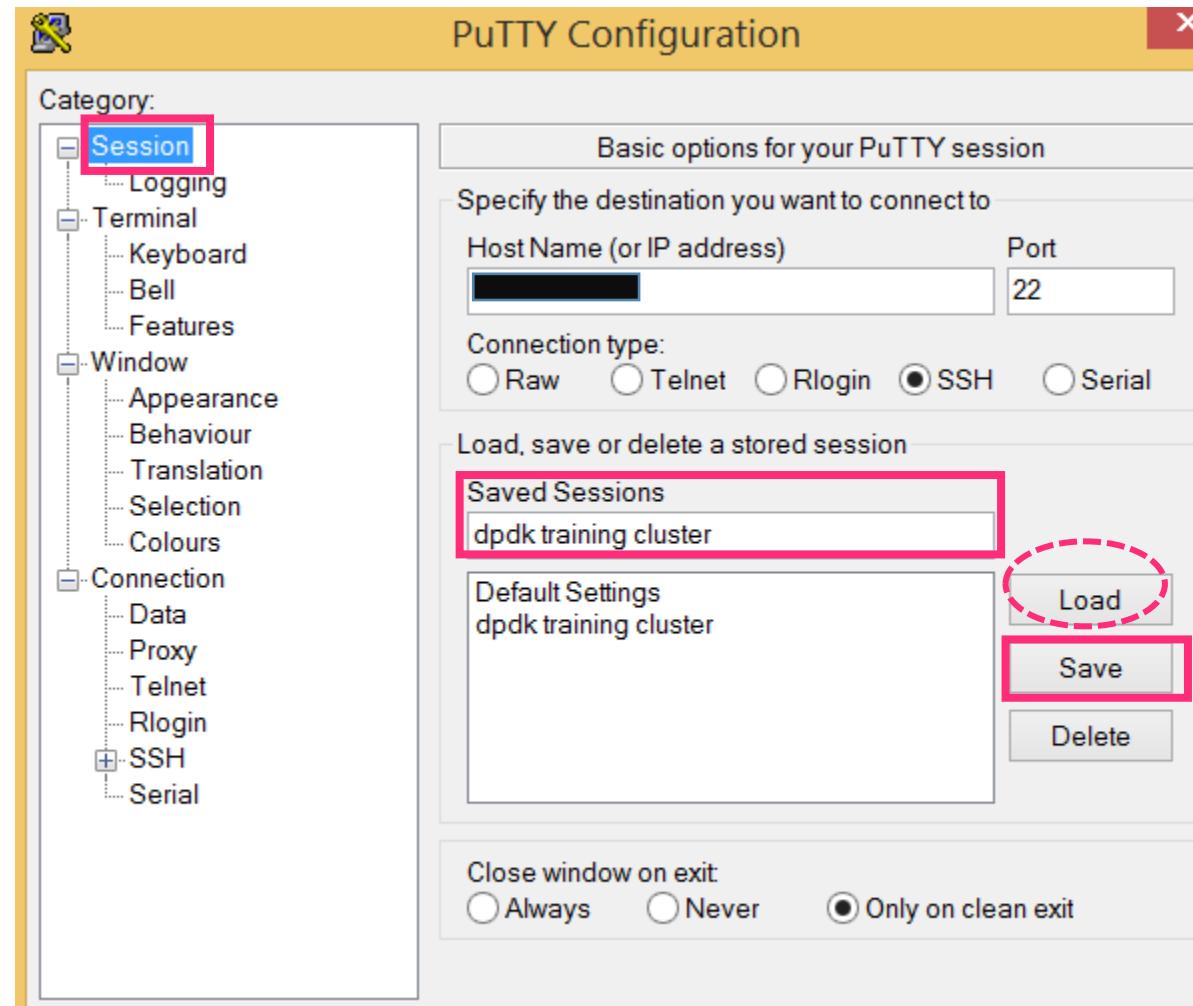
Auto



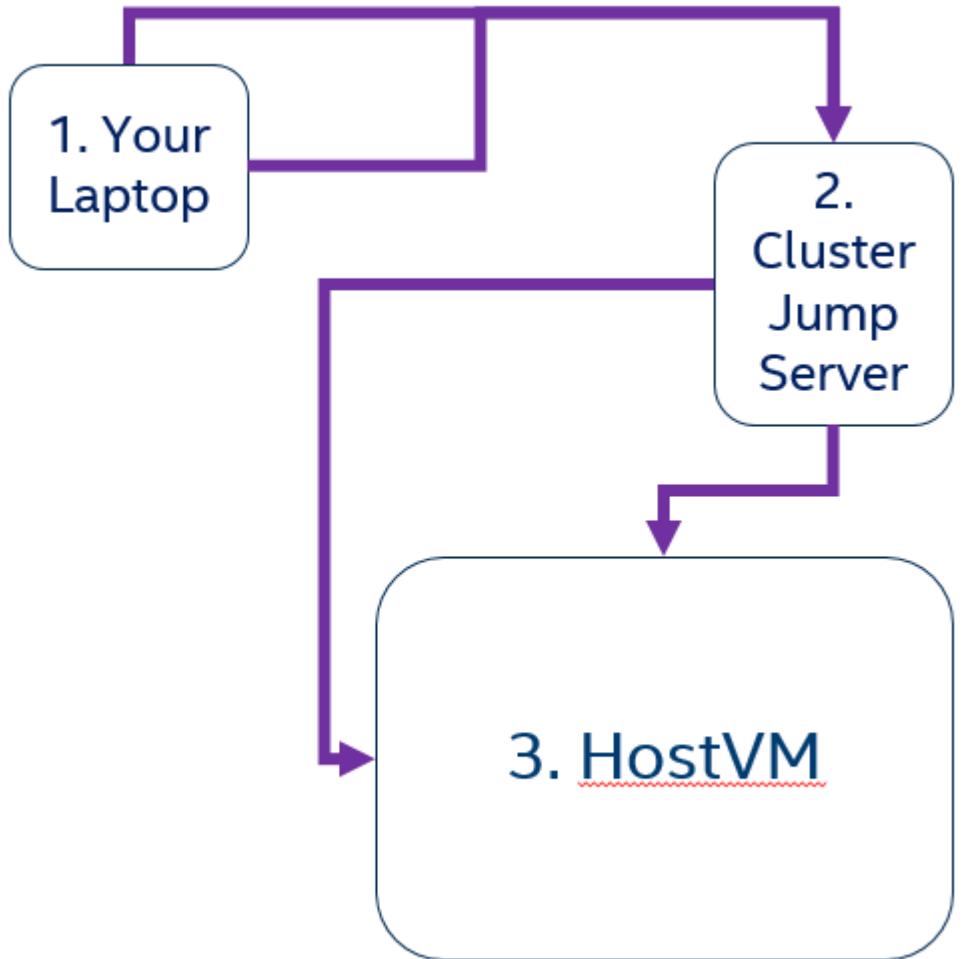
▪ Next Step is To Save The Configuration

# 1) Name The Configuration. 2) Save It

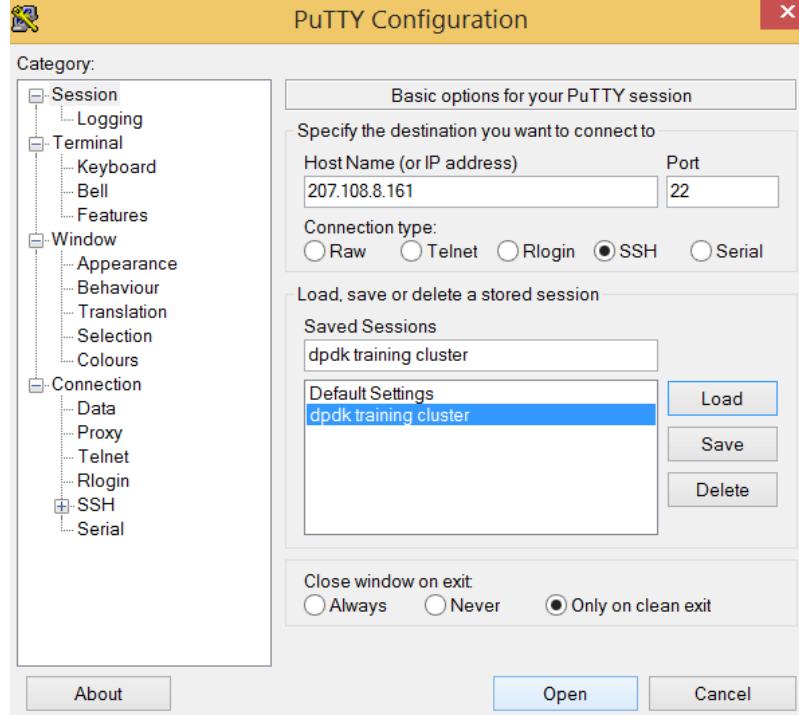
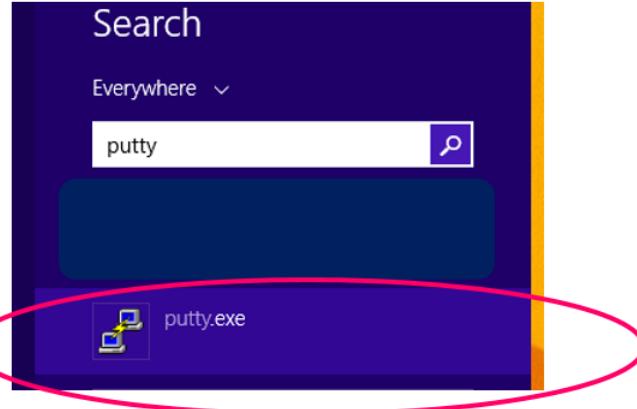
- Click “Session” on the top left corner.
- It will take you to the 1<sup>st</sup> screen - shown here .
- In Saved Sessions Box, Enter a name, e.g., dpdk training cluster
- Press Save Button.
- From now onwards, you can load the saved session when starting PuTTY



# Getting To The HostVM-<m> through Cluster Jump Server



# Start PuTTY



- Find PuTTY that you just installed
- Click on icon “putty.exe”. You will get the PuTTY Config Screen
- Select the Saved Session, shown here e.g., dpdk training cluster.
- Press Load Button. Press Open Button.
- You will get ssh session shown in 3<sup>rd</sup> screen asking username
- Username: student<m> <m> is given to you. For example it may be student19 or student25 or any other student<m>.
- Ask for your specific <m> and use only that. This will avoid overlapping with other teams’ <m>.
- Password is same as username. For example, if your username is student19, then password is student19
- Repeat the above steps so that you have many connections to the jump server

# How to login and Connect to HostVM-<m>



```
student20@dlogin1:~  
login as: student20  
student20@207.108.8.161's password:  
Last login: Fri Nov 25 16:23:04 2016 from 134.134.139.76  
[student20@dlogin1 ~]$ ssh HostVM-█
```

- **Login as student<m>    password student<m>**

- The next step is to connect to HostVM-<m> assigned to you.
- In case not given, Ask for your own <m>
- Note: We are using 1:1 – same student name and same HostVM- name
- Type in the following command

- **ssh HostVM-<m>**

- Username: password Password: password

- **sudo su –**

- **cd /home/user/dpdk\_hands\_on**

- **Vi commands.txt**

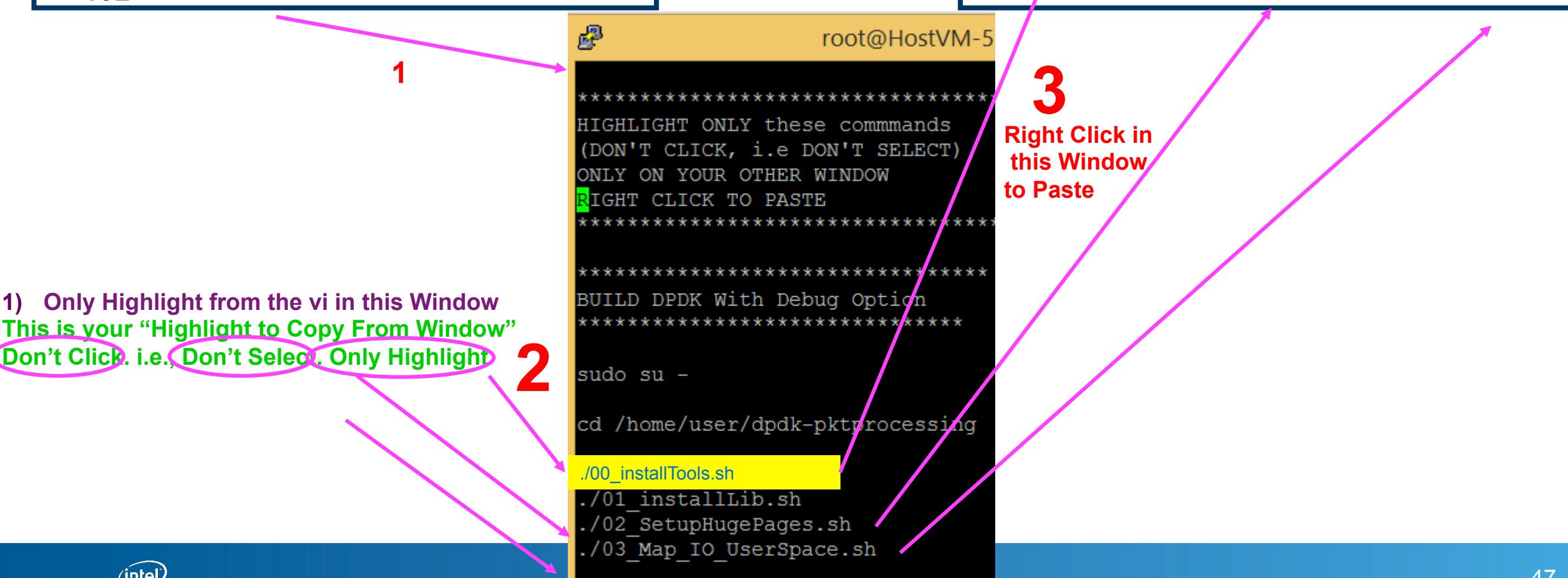
# Open 2 Windows

## A) FIRST, In the 1<sup>st</sup> Window, enter these commands

```
ssh HostVM-<M>  
Password is password  
sudo su -  
cd /  
vi Cpy_PasteFromMe.txt
```

## B) NEXT, in the 2<sup>nd</sup> Window, enter these commands

```
ssh HostVM-<M>  
Password is password  
sudo su -  
cd /
```



# Cpy\_Paste From Me.txt

# \*\*\*\*\* TEST 1 \*\*\*\*\*

```
#Cpy_PasteFromMe.txt
```

```
sudo su -
```

```
cd /
```

```
./00_ReserveHugePage.sh
```

```
source 01_SetEnvironmentVariables.sh
```

```
cd $RTE_SDK
```

```
*****
```

```
**** Please be patient
```

```
*** The following commands will take time to start
```

```
*****
```

```
make config T=$RTE_TARGET
```

```
make
```

```
make install T=$RTE_TARGET
```

```
cd $GTEST_DIR
```

```
cmake CMakeLists.txt
```

```
make
```

```
cd $TLDK_ROOT
```

```
make clean
```

```
make all O=$RTE_TARGET
```

```
# Shall we take a look at key data structures?
```

```
# For that let us use gdb and run Dynamic Ring Buffer Tests
```

```
# *****
```

```
# *****
```

```
gdb --arg ./x86_64-native-linuxapp-gcc/app/test_dring
```

```
# Set BreakPoints
```

```
b drb_dump
```

```
b tle_dring_dump
```

```
run
```

```
#Examine dynamic ring
```

```
ptype dr
```

```
# Examine dynamic buffer
```

```
ptype db
```

```
continue
```

```
*****
```

STOP! OBSERVE the tle\_dring structure

WHAT DO YOU SEE?

Producer / Consumer pointers with Head and Tail Pointers.

WHAT VALUE YOU SEE for Head and Tail Pointers?

\*\*\*\*\*

continue

continue

\*\*\*

STOP: Now what value you see for head and tail pointers

\*\*\*\*\*

continue

# type continue till you see [Inferior 1 (process <> exited normally)]

quit

# \*\*\*\*\* TEST 2 \*\*\*\*\*

# Now run the same dring test without gdb

# \*\*\*\*\*

./x86\_64-native-linuxapp-gcc/app/test\_dring

```
*****
#*****
#*****TEST 3 *****
# Shall we run the Layer 4 UDP Application? Echo mode?
#*****
```

```
cd $TLDK_ROOT
./02_run_UDP_Echo_Test.sh
```

```
# Check that all packet in output pcap file contain valid
udp checksum:
```

How do you “repair” the build problem with gtest?

# How To Repair Build error?

- unset GMOCK\_DIR
- unset GTEST\_DIR
- cd \$TLDK\_ROOT
- cd make clean
- Make all O=\$RTE\_TARGET
- cd googletest\_feb2017
- cmake CMakeLists.txt
- Make
- export GTEST\_DIR=/googletest\_feb2017/googletest
- export GMOCK\_DIR=/googletest\_feb2017/googlemock
- cd \$TLDK\_ROOT
- make clean
- make all O=\$RTE\_TARGET