# Benchmarking and profiling DNS systems with modern Linux tools

Robert Edmonds (edmonds@fsi.io) Farsight Security, Inc.

### Intro

- Demonstrate Linux tools useful for DNS benchmarking
  - trafgen, ifpps (from the netsniff-ng toolkit)
  - tc ("traffic control")
  - perf\_events
  - mpstat
- Use output to generate graphs
- Profile hotspots revealed by benchmarking

#### Benchmarking

Slide 2 of 28

## Hardware

- Entry level hardware, two years old
- Receiver:
  - Intel Xeon E3-1245v3 (3.4 GHz, quad core)
- Sender:
  - Intel Core i3-4130 (3.4 GHz, dual core)
- No HyperThreading, SpeedStep, Turbo Boost
- Intel Ethernet Server Adapter I350-T2 (PCIe cards)
  - Dedicated 1 Gbit/s link, directly connected

#### Benchmarking

Slide 3 of 28

## Software

- Linux kernel 4.2
  - Defaults, no attempts to tune networking stack
- DNS servers
  - BIND 9.10.3 (ISC)
  - Knot DNS 2.0.1 (CZ.NIC)
  - Unbound 1.5.4 (NLnet Labs)
  - dns-echo 0.1 (Ray Bellis / ISC)

#### Benchmarking

Slide 4 of 28

# Software

- Benchmarking tools
  - netsniff-ng toolkit
    - trafgen
    - ifpps
  - tc
  - mpstat
  - perf\_events

Benchmarking

Slide 5 of 28

# trafgen (netsniff-ng toolkit)

- High speed userspace network traffic generator
- Uses packet(7) TX\_RING interface of the Linux kernel
- Low-level packet description language
- Sample DNS trafgen config on next two slides
- Easily generates >1Mq/s with two CPU cores

```
# Ethernet header
0xa0, 0x36, 0x9f, 0x28, 0x96, 0x11, # Destination MAC
0xa0, 0x36, 0x9f, 0x28, 0x95, 0xb1, # Source MAC
const16(0x0800), # Ethernet protocol: IPv4
# IP header
               # IP version (4b), header length (4b)
0b01000101,
Θ,
                   # TOS
const16(71), # Total length
drnd(2),
                  # Fragment ID (randomized)
Ob00000000, 0, # Flags (3b), fragment offset (13b)
                   # IP TTL
64,
17,
                  # IP protocol: UDP
csumip(14, 33), # IP checksum
192, 168, 0, 1, # Source IP address
192, 168, 0, 53, # Destination IP address
# UDP header
drnd(2),
                  # UDP source port (randomized)
const16(53053), # UDP destination port
const16(51), # UDP length
const16(0), # UDP checksum (disabled)
```

# DNS header							
drnd(2),	<pre># ID (randomized)</pre>						
const16(0x0100),	<pre># Flags: RD=1</pre>						
const16(1),	# Question count						
const16(0),	# Answer count						
const16(0),	<pre># Authority count</pre>						
const16(1),	<pre># Additional count</pre>						
<pre># DNS question section</pre>							
0x03, "WWW",	#						
0x06, <mark>"google</mark> ",	#						
0x03, "com",	#						
0x00,	<pre># Question name: "www.google.com."</pre>						
const16(28),	<pre># Question type: AAAA</pre>						
const16(1),	<pre># Question class: IN</pre>						
<pre># DNS additional sec</pre>	ction						
0×00,	# "."						
const16(41),	# <b>OPT</b>						
const16(4096),	<pre># UDP payload size</pre>						
0x00,	# Extended RCODE						
0x00,	# EDNS version						
0x00, 0x00	# Z						
const16(0)	<pre># Data length</pre>						

# trafgen (netsniff-ng toolkit)

- Send packets on *iface* until stopped:
  - # trafgen --qdisc-path -i trafgen.cfg -o \${iface}
- This generates >1Mq/s
- Need a way to control the query rate

# tc ("traffic control")

- Linux kernel facility that can shape bits going out an interface
- E.g., allow **x** Mbit/sec out *iface*:
  - # tc qdisc replace dev \${iface} root handle 1: tbf \
     burst 20480 limit 20480 rate \${x}Mbit
- For each benchmark run, ramp x from 25 to 750 Mbps in 25 Mbps increments

Benchmarking

Slide 10 of 28

# tc ("traffic control")

- tc limits are specified by bandwidth (bits/sec), not packets/sec
- Query messages are identical, packets are all the same size (85 bytes)
- Therefore:
  - 250 Mbit/s  $\rightarrow$  368 Kq/s
  - 500 Mbit/s  $\rightarrow$  736 Kq/s
  - 750 Mbit/s  $\rightarrow$  1.10 Mq/s

#### Benchmarking

Slide 11 of 28

# ifpps (netsniff-ng toolkit)

- Network interface packet statistics
- Queries /proc, low overhead
- "Top"-like interface or CSV output

### mpstat

- Record system idle % during the benchmark run
- 100% minus idle % is the overall system load
  - 0% idle  $\rightarrow$  system is 100% loaded

Slide 13 of 28

# Orchestration

- Home grown Python script:
  - Start server on receiver
  - Clamp sender to 1 Mbit/s (tc)
  - Start trafgen on sender
  - Ramp sender from 25 Mbit/s .. 750 Mbit/s
    - Record network Rx/Tx load (ifpps)
    - Record system load on receiver (mpstat)
  - Stop server, trafgen
- Repeat for each DNS server

#### Benchmarking

Slide 14 of 28

# **Performance graphs**

- Responses/second graphed against system load
  - R/s taken from ifpps
  - System load taken from mpstat
- Want a nice line with low slope and evenly spaced points topping out at 100% system load
- <100% response rate will cause a cluster of points</p>

#### BIND 9.10.3







### Unbound 1.5.4







### dns-echo 0.1





# perf\_events

- Kernel-based lightweight profiler
- Profiles the whole system: userspace and kernel
- "perf top": real-time "top"-like output, but at the level of individual functions
- "perf record", "perf report", "perf diff": take multiple snapshots throughout the benchmark run, diff results
- Call-graph visualization
  - Compile with -fno-omit-frame-pointer

## perf\_events

Record a 3 second snapshot, with call-graph data:

# perf record --call-graph fp -a -o perf.data sleep 3

Start interactive tool to browse the profile:

# perf report

Benchmarking

Slide 25 of 28

Samples:	35K of	event	'cycles'	, Event	count (approx.	): 2539992999	3	
Childr	en 🤤	Self	Command		Shared Object		Symbol 🔶	
- 67.8	2% 0	.00%	named		libpthread-2.	19.so	[.] start_t	
- sta	rt_threa	d						
-	run							
- 98.57% client_request								
- 76.94% ns_query_start								
- 92.93% query_find								
- 44.63% ns_client_send								
- 86.32% client_send								
- 31.01% dns_message_rendersection								
- 99.27% dns_rdataset_towiresorted								
- 95.33% towiresorted.isra.3								
- 67.42% dns_name_towire								
- 44.60% dns_compress_findglobal								
57.17% dns_name_hash								
26.86% dns_name_getlabelsequence								
10.65%memcmp_sse4_1								
3.19% dns_name_caseequal								
1.07% memcmp@plt								
	1.06%netif_receive_skb_core							
or a hi	aher lev	el ove	erview, t	rv: per	f reportsort	comm,dso		

Benchmarking

C

Slide 26 of 28



Benchmarking

Slide 27 of 28

## Thanks!

Benchmarking

Slide 28 of 28