

Ze života Knot DNS

Daniel Salzman • daniel.salzman@nic.cz • 10. 11. 2021



Knot DNS – stručný úvod

- Autoritativní DNS server
 - Primární nebo sekundární server
 - Vysoký výkon zpracování DNS dotazů
- Pokročilá podpora DNSSEC
 - Jednoduché použití
 - Výkonné podepisování
- A mnoho dalšího
 - www.knot-dns.cz



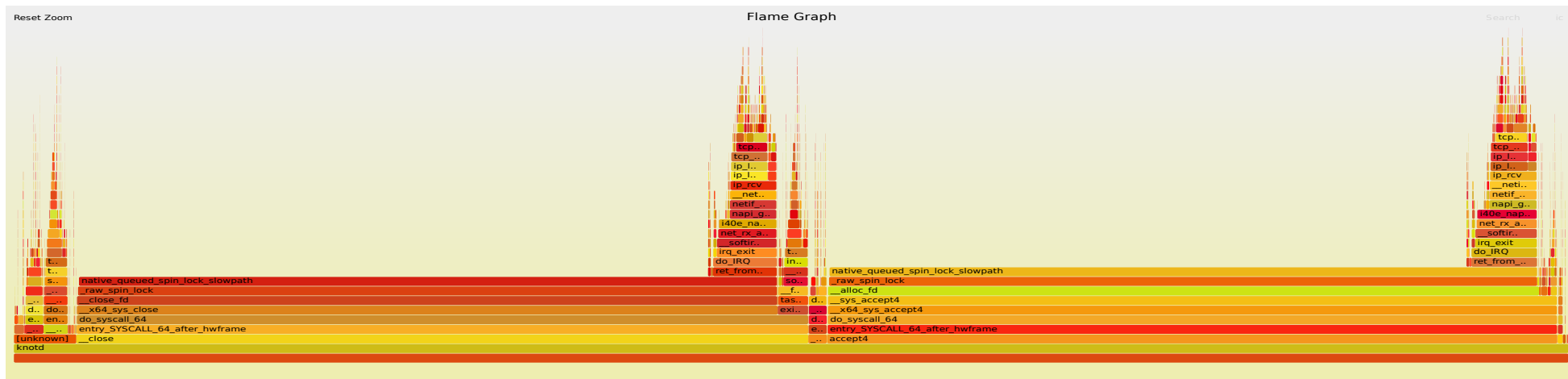
Výkon zpracování DNS dotazů

- Výkon nelze charakterizovat jedním číslem, protože závisí na spoustě faktorů
 - Velikost a struktura zón (TLD vs ROOT vs. jednoduchá zóna)
 - Počet zón (TLD vs. DNS hosting)
 - Zabezpečení DNSSEC (algoritmus DNSKEY, NSEC vs. NSEC3)
 - Typ DNS dotazu (pozitivní vs. negativní odpověď, DNSSEC)
 - Konfigurace démona (dodatečné moduly, počet adres,...)
- Stále narážíme na další scénáře nasazení, pro které je třeba optimalizovat implementaci
- Moderní hardware nabízí nové možnosti ale i výzvy
 - CPU se spoustou jader
 - Síťové karty s vysokou propustností
- Operační systém
 - Možnosti a nastavení síťového stacku (SO_REUSEPORT, recvmsg/sendmsg,...)
 - Přibývající bezpečnostní opravy v OS často zhoršují výkon aplikací



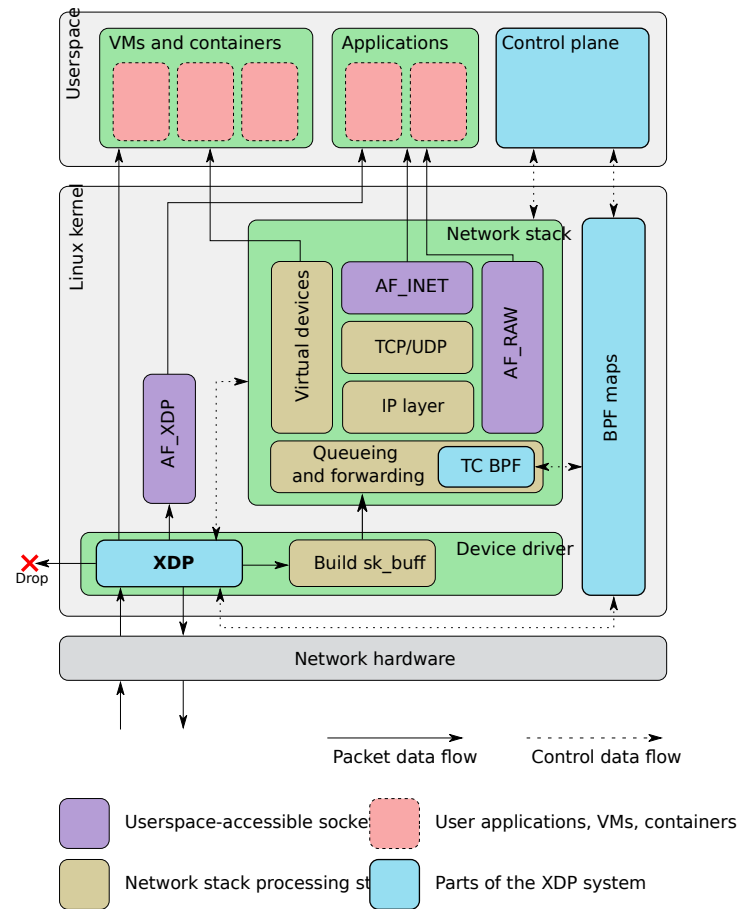
Problémy konvenčního síťového stacku

- Mnoho složitých vrstev zpracování
- Kopírování paketů mezi kernelem a aplikací
- Sdílená tabulka deskriptorů TCP spojení mezi vlákny procesu



Režim XDP – popis

- Podpora pro UDP v Knot DNS 3.0 (pro TCP v 3.1)
- Klasifikace paketů v XDP (eBPF program)
 - Nativní režim – vyžaduje podporu ovladače
 - Emulovaný režim – nízký výkon, nevhodné do produkce
- Využití socketů AF_XDP (Linux 4.18+)
 - Předání paketů z XDP přímo do paměti aplikace v user space (Zero-Copy od Linux 5.0+)
- Nevyžaduje úpravu síťového ovladače
- Obcházení síťového stacku OS pro vybrané pakety
- Extra oprávnění procesu jsou třeba pouze při startu
- Ostatní provoz není dotčen a je zpracován normálně



"Kernel diagram" by Toke Høiland-Jørgensen licensed under CC-BY-SA

Výkon zpracování DNS dotazů po UDP

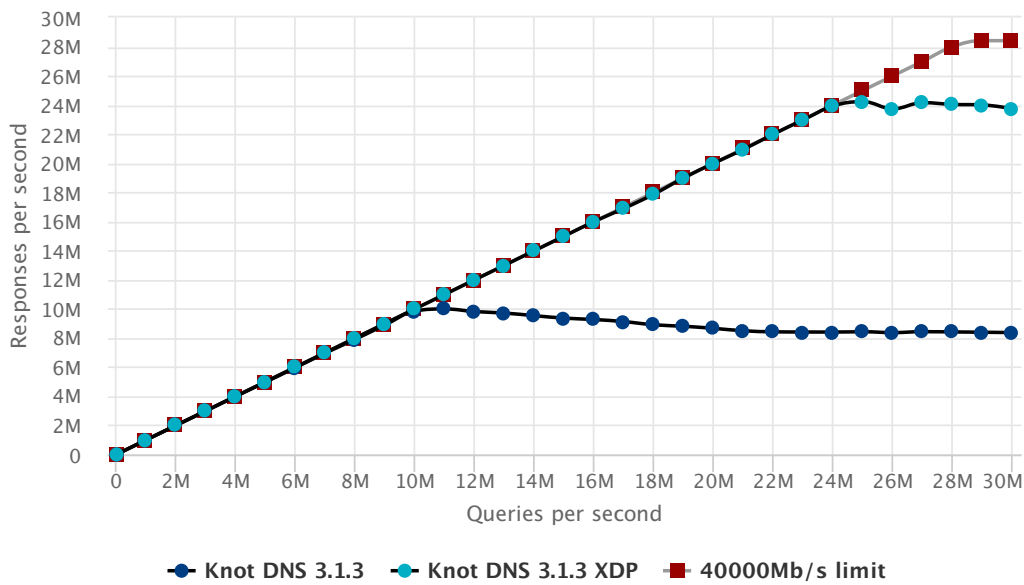
- htop: 64 cores, 30M qps, XDP

```
1[|||||] 100.0% 17[|||||] 100.0% 33[|||||] 100.0% 49[|||||] 100.0%
2[|||||] 100.0% 18[|||||] 100.0% 34[|||||] 100.0% 50[|||||] 100.0%
3[|||||] 100.0% 19[|||||] 100.0% 35[|||||] 100.0% 51[|||||] 100.0%
4[|||||] 100.0% 20[|||||] 100.0% 36[|||||] 100.0% 52[|||||] 100.0%
5[|||||] 100.0% 21[|||||] 100.0% 37[|||||] 100.0% 53[|||||] 100.0%
6[|||||] 100.0% 22[|||||] 100.0% 38[|||||] 100.0% 54[|||||] 100.0%
7[|||||] 100.0% 23[|||||] 100.0% 39[|||||] 100.0% 55[|||||] 100.0%
8[|||||] 100.0% 24[|||||] 100.0% 40[|||||] 100.0% 56[|||||] 100.0%
9[|||||] 100.0% 25[|||||] 100.0% 41[|||||] 100.0% 57[|||||] 100.0%
10[|||||] 100.0% 26[|||||] 100.0% 42[|||||] 100.0% 58[|||||] 100.0%
11[|||||] 100.0% 27[|||||] 100.0% 43[|||||] 100.0% 59[|||||] 100.0%
12[|||||] 100.0% 28[|||||] 100.0% 44[|||||] 100.0% 60[|||||] 100.0%
13[|||||] 100.0% 29[|||||] 100.0% 45[|||||] 100.0% 61[|||||] 100.0%
14[|||||] 100.0% 30[|||||] 100.0% 46[|||||] 100.0% 62[|||||] 100.0%
15[|||||] 100.0% 31[|||||] 100.0% 47[|||||] 100.0% 63[|||||] 100.0%
16[|||||] 100.0% 32[|||||] 100.0% 48[|||||] 100.0% 64[|||||] 100.0%
Mem[|||||] 2.91G/62.4G Tasks: 29; 1 running
Swp[|||||] 0K/2.00G Load average: 22.99 17.91 8.96
Uptime: 19:28:29
```

- htop: 64 cores, 30M qps

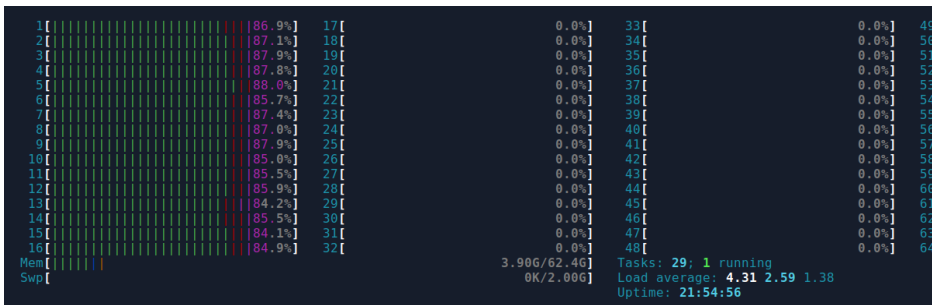
```
1[|||||] 100.0% 17[|||||] 100.0% 33[|||||] 100.0% 49[|||||] 100.0%
2[|||||] 100.0% 18[|||||] 100.0% 34[|||||] 100.0% 50[|||||] 100.0%
3[|||||] 100.0% 19[|||||] 100.0% 35[|||||] 100.0% 51[|||||] 100.0%
4[|||||] 100.0% 20[|||||] 100.0% 36[|||||] 100.0% 52[|||||] 100.0%
5[|||||] 100.0% 21[|||||] 100.0% 37[|||||] 100.0% 53[|||||] 100.0%
6[|||||] 100.0% 22[|||||] 100.0% 38[|||||] 100.0% 54[|||||] 100.0%
7[|||||] 100.0% 23[|||||] 100.0% 39[|||||] 100.0% 55[|||||] 100.0%
8[|||||] 100.0% 24[|||||] 100.0% 40[|||||] 100.0% 56[|||||] 100.0%
9[|||||] 100.0% 25[|||||] 100.0% 41[|||||] 100.0% 57[|||||] 100.0%
10[|||||] 100.0% 26[|||||] 100.0% 42[|||||] 100.0% 58[|||||] 100.0%
11[|||||] 100.0% 27[|||||] 100.0% 43[|||||] 100.0% 59[|||||] 100.0%
12[|||||] 100.0% 28[|||||] 100.0% 44[|||||] 100.0% 60[|||||] 100.0%
13[|||||] 100.0% 29[|||||] 100.0% 45[|||||] 100.0% 61[|||||] 100.0%
14[|||||] 100.0% 30[|||||] 100.0% 46[|||||] 100.0% 62[|||||] 100.0%
15[|||||] 100.0% 31[|||||] 100.0% 47[|||||] 100.0% 63[|||||] 100.0%
16[|||||] 100.0% 32[|||||] 100.0% 48[|||||] 100.0% 64[|||||] 100.0%
Mem[|||||] 1.83G/62.4G Tasks: 29; 1 running
Swp[|||||] 0K/2.00G Load average: 15.99 15.00 7.35
Uptime: 19:27:18
```

Response Rate
Linux 5.13.0, TLD, (2021-11-01)

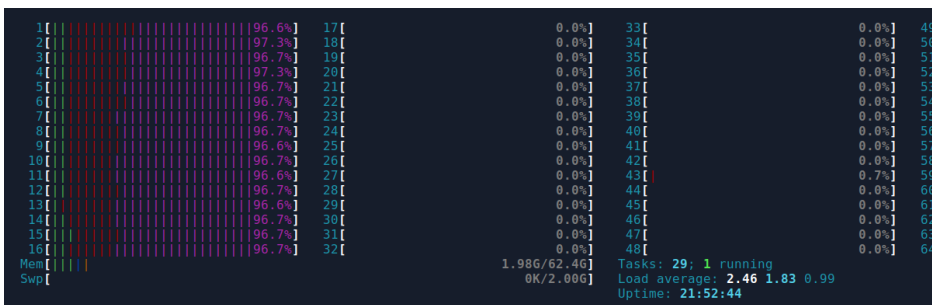


Výkon zpracování DNS dotazů po TCP

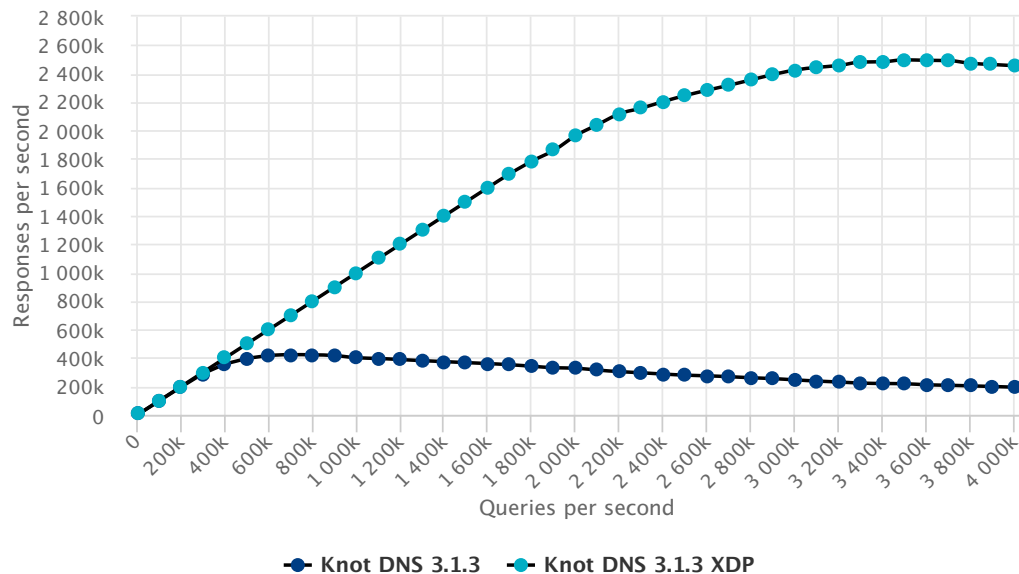
- htop: 16 cores, 4M qps (7Gbps), XDP



- htop: 16 cores, 4M qps



Response Rate
Linux 5.13.0, TCP TLD, (2021-11-02)



Měření výkonu zpracování DNS dotazů

- Dříve využívané nástroje již nevyhovují a špatně se používají
 - tcpreplay v režimu Netmap (UDP)
 - dpdk-tcp-generator (TCP)
- Absence vhodných alternativních nástrojů
- Potřeba vytvoření nového nástroje – **kxdpgun**
 - Součást Knot DNS – využití implementace režimu XDP
 - Vícevláknový generátor DNS provozu (UDP nebo TCP)
 - Vstupem je jednoduchý textový soubor s popisem dotazů (EDNS0, DO bit)
 - Na běžném serveru je možné generovat desítky miliónů dotazů za sekundu



Monitoring DNS provozu v režimu XDP

- Běžné nástroje (tcpdump, wireshark,...) nelze použít
- Modul statistik DNS provozu (**mod-stats**)
 - Rozšíření o nové čítače související s XDP
 - Přeprogramování modulu kvůli nežádoucí režii při vysokých rychlostech na procesorech s mnoha jádry
- Nový modul pro zjednodušený export DNS provozu (**mod-probe**)
 - Minimalistická a efektivní implementace bez zbytečných transformací
 - Předávání struktury parametrů (bez záznamů odpovědi) pro dvojici dotaz-odpověď
 - Možnost použití více paralelních instancí pro lepší rozložení zátěže
 - Možnost vzorkování od určité rychlosti
 - Rozhraní v C (integrace do dns-probe projektu ADAM) nebo Pythonu (jednoduché použití)



Režim XDP a routování

- Základní režim XDP je založen na symetrickém routování
 - Vstupní rozhraní dotazu = výstupní rozhraní odpovědi
 - Zdrojová MAC adresa dotazu = cílová MAC adresa dotazu
 - Pro správný chod vyžaduje předřazený router
- Knot DNS 3.1 přidává podporu „route-check“
 - Při zpracování dotazu v XDP se vyčte z OS routovací informace pro odpověď
 - Routovací tabulky může spravovat nezávislý démon (BIRD)
 - Nastavení správné cílové MAC adresy odpovědi
 - Pokud se výstupní rozhraní liší od vstupního, zpracuje se DNS dotaz konvenčně
 - Dotaz se zahodí pokud je routovací pravidlo blackhole/unreachable/prohibit



Plány do blízké budoucnosti

- Dokončení implementace DNS-over-TCP v režimu XDP
- Experimenty s DNS-over-QUIC v režimu XDP
- Testování se 100GbE síťovou kartou





Děkuji za pozornost

Daniel Salzman • daniel.salzman@nic.cz • knot-dns.cz