

Proiectarea și implementarea middleboxurilor software practice

Sumar

Universitatea Polithenica Bucureti



Teză de doctorat

Vladimir Olteanu

Coordonator: Valentin Cristea

Cuvinte-cheie

- rețele
- datacenter
- middlebox
- protocol

Sumar

Această teză abordează proiectarea middlebox-urilor software scalabile orizontal.

Middlebox-urile sunt mașini care procesează pachete de rețea, efectuând operațiuni care pot fi simple, ca filtrarea lor pe baza unor cămpuri, sau complexe, ca detectarea scurgerilor de informații secrete. Ele sunt un element de bază al internetului de astăzi.

În mod tradițional, middleboxurile erau dispozitive hardware monolitice. Aducerea lor la zi, fie din punct de vedere al performanței, fie al funcționalității, era de obicei un lucru costisitor, deoarece însmenă înlocuirea întregii mașini cu un model diferit. În plus, era necesar ca middleboxurile să facă față volumului de trafic de la orele de vârf; în afara acestor ore, capacitatea suplimentară era, în esență, irosită.

Pentru a rezolva aceste probleme, a apărut un nou domeniu de cercetare, numit Network Function Virtualization. Ideea de bază este că procesarea traficului nu se mai face în dispozitive dedicate, ci în software care rulează într-o manieră distribuită pe hardware generic. Astfel de middleboxuri pot fi scalate orizontal pe măsură ce volumul de trafic crește sau scade: capacitate suplimentară poate fi adăugată pornind câteva mașini în plus, și apoi eliminată oprind unele dintre ele.

Cel mai mare obstacol din calea scalabilității orizontale a middlebox-urilor este faptul că acestea pot ține stare, atât specifică fiecărui flow în parte, cât și partajată de mai multe flow-uri. Această stare este indispensabilă pentru procesarea corectă a traficului; dacă starea ar deveni indisponibilă, chiar și temporar, toate flow-urile care îi corespund ar putea suferi întreruperi. Prin urmare, scalarea orizontală nu implică doar redirecționarea unei părți a traficului către o altă mașină, ci și migrarea tuturor stărilor relevante către mașina menționată.

S-au făcut eforturi pentru a găsi un framework general care să permită evenimente de scalare orizontală aproape fără întreruperi, dar acestea tind să impună un cost destul de mare în ceea ce privește performanța. Subtextul acestei teze este că argumentăm împotriva unor astfel de soluții generalizate; în acest scop, proiectăm două middleboxuri separate și folosim soluții diferite pentru a le rezolva problemele de scalabilitate.

Primul middlebox pe care îl discutăm este un carrier-grade NAT. Acesta este un middlebox destul de simplu de implementat, care ține stare atât pentru fiecare flow, cât și pentru fiecare IP extern. Implementarea noastră are un algoritm de migrare a stării care impune o penalizare minimă de performanță. Aici arătăm că, pur și simplu riscând reordonarea unor pachete, dar păstrând toate celealte garanții de funcționare, perturbările de performanță impuse de evenimentele de scalare sunt minime. Reordonările occasionale ale pachetelor sunt, în general, un lucru căruia protocoalele de la nivelul 4 și mai sus îi pot face față cu ușurință.

Al doilea middlebox discutat este Beamer, un load-balancer stateless pentru cloud, care funcționează și cu MPTCP. Load-balancer-ele stateful își amintesc asociările dintre conexiuni și serverele care se ocupă de ele. Aceasta, în esență, încearcă să dubleze informațiile păstrate de serverele însăși în tabelele lor de conexiuni. Prin faptul că nu ține niciun fel de stare, Beamer este ușor de implementat atât în software, cât și în hardware, și evită multe dintre problemele cu care se confruntă load-balancer-ele în mod general. Beamer poate face față cu ușurință atacurilor de tip SYN flood, ceva cu care toate load-balancer-ele stateful au probleme.

În cele din urmă, vom discuta despre ceva care este un middlebox doar în sensul larg al cuvântului: am proiectat versiunea 6 a protocolului SOCKS și am implementat un proxy. Versiunea 6 tratează multe dintre problemele versiunii 5: penalizarea din punct de vedere al RTT-urilor este, în majoritatea cazurilor, zero (sau chiar negativ!), un serviciu DNS implicit este furnizat clienților, și protocolul profită din plin de noile îmbunătățiri introduse în TLS 1.3. Deoarece, strict vorbind, proxy-ul este un serviciu care rulează peste TCP, deployment-uri mari de SOCKSv6 pot fi făcute folosind Beamer.

Cuprins

1	Introducere	3
2	Fundal	5
2.1	Middleboxurile în arhitectura Internetului	5
2.2	OpenFlow și SDN	6
2.3	Multipath TCP	7
3	Principii pentru procesarea scalabilă de pachete	9
3.1	Starea de fapt	12
3.2	Contribuțiile tezei	13
3.3	Publicații	16
4	Carrier-grade NAT	17
4.1	Introducere și motivare	17
4.2	Proiectarea unui Carrier-Grade NAT	19
4.3	Implementare	22
4.4	Evaluare	24
4.5	Concluzii	29
5	Beamer	31
5.1	Introducere	31
5.2	Fundal	32
5.3	Limitările load balancingului stateful	34
5.4	Beamer: load-balancing stateless	36
5.4.1	Hashing stabil	37
5.4.2	Daisy chaining	39
5.4.3	Algoritmul din dataplane	40
5.5	Tratarea Multipath TCP	41
5.5.1	Abordarea MPLB	42
5.5.2	Abordarea Beamer	46
5.6	Control plane-ul	47
5.7	Defragmentare	49
5.8	Implementare	51

5.9	Evaluare	52
5.9.1	Micro-benchmarkuri	52
5.9.2	Scalabilitate și robustețe	55
5.9.3	Load balancing de HTTP peste MPTCP	59
5.9.4	Scalabilitatea controllerului	61
5.9.5	Hashing stabil	62
5.9.6	Defragmentare	62
5.10	Concluzii	63
6	Procoloul SOCKS, versiunea 6	65
6.1	Introducere	65
6.2	Modul de funcționare	66
6.3	Cereri	68
6.4	Răspunsuri pentru versiuni incompatibile	69
6.5	Răspunsuri pentru autentificare	70
6.6	Răspunsuri pentru operațiuni	71
6.6.1	Comanda CONNECT	72
6.6.2	Comanda BIND	72
6.6.3	Comanda UDP ASSOCIATE	73
6.7	Opțiuni SOCKS	77
6.7.1	Opțiuni pentru stiva de protocoale	77
6.7.2	Opțiuni care conțin date de autentificare	84
6.7.3	Opțiuni legate de metodele de autentificare	84
6.7.4	Opțiuni legate de sesiuni	86
6.7.5	Opțiuni pentru idempotență	88
6.8	Autentificare cu nume de utilizator și parolă	90
6.9	Autentificare folosind SSL	92
6.10	TCP Fast Open între client și proxy	92
6.11	Starturi false	92
6.12	DNS furnizat de SOCKS	92
6.13	Considerații de securitate	93
6.13.1	Cereri mari	93
6.13.2	Atacuri cu replay	94
6.13.3	Epuizarea resurselor	94
6.14	Considerații de confidențialitate	94
6.15	Timing-ul răspunsurilor	94
6.16	Implementare	96
6.17	Evaluare	96
6.18	Concluzii	97
7	Lucrări înrudite	99
8	Concluzii	101