

DETEKCE ANOMÁLIÍ V SÍŤOVÉM PROVOZU

Václav Bartoš

Výpočetní technika a informatika, 1. ročník, prezenční studium

Školitel: Lukáš Sekanina

Fakulta Informačních Technologií, Vysoké Učení Technické v Brně

Božetěchova 1/2, 612 66 Brno

ibartosv@fit.vutbr.cz

Abstrakt. Práce se zabývá detekcí anomálií v síťovém provozu. Jsou diskutovány některé praktické problémy výzkumu v této oblasti a je představena sada nástrojů ulehčujících řešení těchto problémů. Dále je představena myšlenka přesunu detekčního algoritmu co nejdříve zdroji dat, tedy do exportéru záznamů o síťových tocích, a jsou diskutovány možnosti dalšího směru autorova výzkumu a cíle jeho disertační práce.

Klíčová slova. síťová bezpečnost, detekce anomálií, framework, NetFlow

1 Úvod

Systémy detekce anomálií je kategorie IDS systémů, která pracuje na principu detekce statistických odchylek v naměřených datech o síťovém provozu. Narozdíl od systémů založených na vyhledávání vzorů tedy neprohledává obsah paketů. Výzkum v oblasti detekce anomálií je v posledních letech poměrně aktivní a bylo navrženo velké množství různých metod. Tato práce diskutuje některé problémy v této oblasti, kterým se chci v rámci své disertační práce věnovat.

Jedním z problémů současného stavu výzkumu metod detekce anomálií je právě množství různých metod spolu s nedostatkem jejich objektivního vyhodnocení či porovnání. Přestože existuje několik článků, které tyto metody shrnují a kategorizují podle různých vlastností (např. [7, 10]), nebyla publikována žádná práce, která by prezentovala přímé experimentální porovnání detekčních schopností více metod.

Takové porovnání lze v omezené míře najít jen v některých publikacích navrhuje nových metody, kde autoři srovnávají výsledky nově navržené metody s jednou nebo dvěma staršími (např. [9]). Takovéto srovnání však není ani úplné, ani objektivní, a navíc pro autory znamená mnoho práce navíc, protože musí implementovat i tyto starší metody. Obvykle totiž k těmto metodám neexistují veřejně dostupné implementace.

Provedení skutečně kvalitního porovnání metod detekce anomálií je však v současnosti velmi obtížné. Kromě nutnosti vytvořit vhodnou metodiku vyhodnocování je příčinou i nedostatek nástrojů pro práci se síťovými daty vhodných pro implementaci metod detekce anomálií a experimenty s nimi. Posledním problémem je nedostatek testovacích dat. Přestože nějaké veřejně dostupné sady existují (MAWI archive, CAIDA), pro mnoho případů jsou nedostatečné a především nejsou anotované, tzn. útoky v nich nejsou popsány. Právě znalost všech útoků a jiných událostí, které by měly být detekovány, je však pro vyhodnocení kvality detekčních metod nezbytná.

Mezi další problémy současných metod detekce anomálií patří často velké detekční zpoždění a také náchylnost měřících systémů k přetížení při velkém zvýšení provozu (např. při DDoS útoku).

V následující kapitole je krátké shrnutí výzkumu v této oblasti v České Republice. Dále je prezentován framework vytvořený s cílem pomoci vyřešit některé z výše popsanych problémů. V kapitole 4 je popsána myšlenka přesunu detekčního algoritmu přímo do sond měřících data o síťovém provozu, což je hlavní náplň mé aktuální práce. V kapitole 5 je pak popsán směr dalšího směřování mého výzkumu a zaměření mé disertační práce. Poslední kapitola shrnuje celou práci.

2 Výzkum síťové bezpečnosti v ČR

V České Republice se výzkumem v oblasti bezpečnosti počítačových sítí zabývá Fakulta informatiky Masarykovy univerzity a Fakulta informačních technologií Vysokého učení technického v Brně, do jisté míry i některé další univerzity (např. ČVUT, jejíž Agent Technology Center má bezpečnost sítí jako jednu z oblastí aplikace agentních systémů) a především správce české akademické sítě – CESNET. Všechny tyto organizace na výzkumu úzce spolupracují a já sám jsem členem výzkumných skupin na VUT i v CESNETu. Díky tomu mám přístup k velkému množství dat o síťovém provozu i k měřícím zařízením, což je nezbytný předpoklad pro provádění výzkumu v této oblasti a k dokončení mé disertační práce.

Dále se výzkumem v této oblasti v ČR zabývá několik málo firem, mezi nimi i AdvaICT a Cognitive-Security, které obě vznikly na základě výsledků výzkumu detekce anomálií na výše zmíněných univerzitách.

3 NADEX framework

Kvůli výše popsaným důvodům jsem se rozhodl vytvořit sadu nástrojů, jejímž cílem je usnadnit implementaci metod detekce anomálií, provádění různých experimentů s těmito metodami a umožnit snadné vyhodnocování a porovnání jejich schopností.

Tato sada nástrojů, či framework, je nazvána NADEX (*Network Anomaly Detection EXperiments*) a obsahuje implementace několika metod detekce anomálií navržených v literatuře a mnoho dalších modulů, skriptů a utilit, které výrazně usnadňují vytváření dalších algoritmů tím, že implementují nejčastěji používané funkce (např. načítání dat v různých formátech). Dále je k frameworku dodáváno několik datových sad využitelných pro testování metod.

Většina skriptů v tomto frameworku je napsána v Pythonu, několik utilit je, především z výkonnostních důvodů, implementováno v C a C++.

V první publikované verzi jsou zahrnuty tři metody detekce anomálií publikované v literatuře a prototyp vlastní metody určené pro rychlou detekci DDoS útoků na flow exportéru (viz kap. 4). Dvě z metod slouží k detekci odchylek v časových řadách. Lze použít např. časové řady objemu provozu v bytech, paketech či tocích, nebo řady entropie adres či portů vyskytujících se v síťovém provozu. První metoda je založena na exponenciálně váhovaném plovoucím průměru (EWMA), implementovaná podle popisu v [6], druhou metodou je algoritmus navržený v [1] využívající vlnkovou transformaci. Třetí metoda místo časových řad zpracovává přímo záznamy o tocích. Jedná se o metodu ASTUTE [9], která hledá korelované změny v objemu jednotlivých toků. Vychází z předpokladu, že v normálním provozu by vlastnosti všech toků měly být nezávislé, zatímco některé útoky (např. DDoS) způsobí nárůst či pokles objemu mnoha toků současně. V dalších verzích frameworku budou postupně přidávány implementace i dalších publikovaných metod.

Kromě těchto algoritmů obsahuje framework mnoho dalších modulů, např. moduly pro načítání různých typů dat – paketů, toků a časových řad. Modul pro načítání paketů z pcap souborů je navíc v současnosti pravděpodobně jediným modulem, který umožňuje v Pythonu načítat pcap soubory bez použití libpcap knihovny (která není multiplatformní). Podobně je tomu s načítáním záznamů o tocích ze souborů ve formátu programu nfdump. Protože nfdump neposkytuje žádný způsob načítání jeho souborů

dalšími aplikacemi, obsahuje framework knihovnu¹ poskytující jednoduché rozhraní pro čtení nfdump souborů. Tato knihovna je napsána v C a je poskytována i obálka v Pythonu, lze ji tedy snadno využít v obou těchto jazycích.

Framework obsahuje i dvě utility pro předzpracování dat. První z nich je program efektivně počítající entropii adres, portů a dalších položek i jejich kombinací ze záznamů o tocích. Tato entropie je využívána v některých metodách detekce anomálií, neboť mnoho typů útoků mění rozložení hodnot některých položek hlaviček paketů, což se projeví v entropii těchto hodnot. Kvůli snížení časové a paměťové náročnosti výpočtů je pro některé položky vypočítávána pouze aproximace a to pomocí metody založené na hashování, která byla prezentována v mé diplomové práci a na soutěži Student EEICT 2011 [2]. Druhým programem je softwarový generátor flow záznamů. Načítá pakety z pcap souboru, simuluje chování flow exportéru s nastavitelnými parametry flow cache a exportované záznamy o tocích ukládá do souboru s flexibilním formátem. Umožňuje navíc ukládat i část payloadu paketů.

Framework také integruje několik již existujících balíčků – *matplotlib* pro generování grafů, *pandas* pro práci s časovými řadami a *dpkt* pro parsování paketů. Dále jsou přiloženy různé skripty, např. pro generování DoS útoků nebo pro detekci skenování portů za základě jednoduchých pravidel.

Pro testování metod detekce anomálií je k frameworku přiloženo i několik datových sad. Jsou to časové řady objemu provozu (byty, pakety, toky) a jeho entropie, obojí navíc rozdělené podle protokolů (všechny, TCP, UDP, ICMP, ostatní). Tato data pocházejí ze dvou linek připojujících síť VUT v Brně do internetu z časového období od září 2011 do února 2012. Dále jsou poskytovány anonymizované NetFlow záznamy z těchto linek. Z důvodu velkého množství dat jsou volně ke stažení jen data z jednoho dne, na požádání je však možné poskytnout více.

Myšlenka vytvoření tohoto frameworku byla prezentována tento rok na konferenci AIMS [4], podrobný popis první zveřejněné verze byl publikován ve formě technické zprávy [3].

Části frameworku nyní sám intenzivně využívám při analýze síťových dat a experimentování s navrhováním nových detekčních algoritmů. Plánuji ho dále rozvíjet, a to především s pomocí bakalářských a magisterských studentů.

4 Detekce anomálií v exportéru

Většina současných metod detekce anomálií pracuje s informacemi získanými ze záznamů o síťových tocích (*flows*). Hlavním zdrojem takovýchto dat jsou tzv. *flow exportéry*, což mohou být buď dedikovaná zařízení zapojená na sledované lince, nebo může tuto činnost zastávat router. Exportér pomocí protokolu NetFlow nebo IPFIX odesílá záznamy o tocích do tzv. *kolektoru*, kde jsou ukládány do souborů nebo databáze. Jeden kolektor často sbírá data z více exportérů. Takto získaná data jsou pak předmětem další analýzy včetně detekce anomálií.

Záznamy o tocích jsou exportovány až po skončení daného toku nebo po vypršení timeoutu, který bývá nastaven až na několik minut. Navíc na kolektoru obvykle nejsou data analyzována ihned, ale po blocích v pravidelných intervalech (nejčastěji 5 minut). I bez započítání doby běhu detekčního algoritmu tak může trvat více než pět minut, než je nežádoucí provoz na síti detekován. Pro efektivní aplikaci protipatření (např. zablokování IP adresy útočnicka) a zmírnění dopadu útoku je však nutné útok detekovat co nejdříve.

Navrhujeme proto umístit detekční algoritmus co nejbližšímu zdroji dat, tedy přímo do měřících sond (flow exportérů). Odstraní se tak všechna výše popsaná zpoždění a útoky je možné detekovat v rámci sekund. Algoritmus běžící v exportéru navíc může využívat více informací, např. informace o jednotlivých paketech nebo různé statistiky o vyrovnávací paměti toků (např. aktuální zaplnění, počet nových toků za vteřinu apod.), která je součástí každého flow exportéru. Informace o probíhajícím útoku je poté odeslána do kolektoru.

¹Knihovna je založená na zdrojových kódech programu nfdump

Detekční algoritmus by měl být velmi jednoduchý a výpočetně nenáročný, aby neomezoval výkon exportéru. Některé typy útoků však lze detekovat poměrně snadno (např. velké DDoS útoky či rozsáhlé skenování sítě). Detekce sofistikovanějších útoků vyžadujících větší výpočetní výkon může být i nadále prováděna na kolektoru.

Na tomto návrhu i praktické realizaci spolupracuji s Rickem Hofstede z University of Twente, Nizozemsko. Vyvíjený systém je zaměřen především na rychlou detekci DDoS útoků, k čemuž je navržen vlastní, výpočetně nenáročný algoritmus. Na exportéru je měřen počet nových toků v pětisekundových intervalech. Časový průběh této hodnoty je analyzován a při detekci náhlého nárůstu je hlášena anomálie. Tato detekce je založena na predikci následující hodnoty, následném určení odchylky skutečné a predikované hodnoty a porovnání této odchylky s (dynamicky vypočtenou) mezí. Predikce bere v úvahu i periodické denní a týdenní odchylky v objemu provozu. Díky použití krátkého časového intervalu může být intenzivní útok detekován během 5 vteřin, méně intenzivní útok během několika desítek vteřin. Přesný popis algoritmu je mimo rámec tohoto textu, navíc se zatím jedná o prototyp, který bude pravděpodobně dále vylepšován.

Prototyp algoritmu byl nejprve implementován ve výše představeném frameworku a poté i jako plugin pro exportér FlowMon od firmy INVEA-Tech. Tento plugin navíc při detekci záplavy podobných paketů (např. SYN flood, což je nejčastější varianta DDoS útoku) filtruje záznamy o tocích této záplavy, čímž zabraňuje možnému přetížení kolektoru. Aby nebyly informace o tomto provozu zcela ztraceny, je nutno odeslat do kolektoru zprávu o vyfiltrovaných tocích. V současné době však používaný kolektor takové zprávy přijímat neumožňuje a proto jsou tyto informace zapisovány jen do souboru na exportéru.

Další možností je agregovat informace o tocích záplavy do jednoho meta-toku a ten odeslat na kolektor. Zde by bylo vhodné použít flexibilní formát IPFIX, kde lze definovat speciální šablonu pro reprezentaci takovýchto informací. Opět je ovšem nutné upravit kolektor, aby tato data uměl zpracovat.

Myšlenka detekčního algoritmu běžícího přímo v exportéru byla prezentována v mé předchozí práci [4] a trochu podrobněji i v práci Ricka Hofstede [8]. V době psaní tohoto textu je plugin experimentálně nasazen na sondě monitorující provoz na 10 Gb/s lince připojující kampus University of Twente do internetu. V nejbližší době je plánováno jeho nasazení i v české akademické síti CESNET. Po dokončení a otestování celého systému je plánována publikace výsledků na konferenci PAM 2013.

5 Další směřování výzkumu

V budoucnu se chci zaměřit na vývoj systému detekce anomálií, který bude založen na udržování profilu chování jednotlivých stanic.

O každé aktivní stanici (IP adrese) v síti lze měřit různé statistické informace, jako např. počet přijatých a odeslaných bajtů a paketů, počet stanic se kterými navázala spojení a kolik stanic navázalo spojení s ní, informace o používaných protokolech, průměrná délka přijatých a odeslaných paketů a mnoho dalších. Pokud budeme měřit či počítat N takových informací, z nichž každá je reprezentována jedním číslem, je chování jedné stanice reprezentováno N -rozměrným vektorem.

Na stanice pak lze pohlížet jako na body v N -rozměrném prostoru. Pomocí různých statistických metod pak lze například detekovat body, které v tomto prostoru leží daleko od ostatních, tedy stanice, které se chovají nestandardně (tj. jinak než ostatní stanice sítě). Pokud navíc budeme udržovat i historii všech hodnot, lze detekovat i změny v chování jednotlivých stanic. Detekce takových změn je velmi užitečná, protože může znamenat například nakažení stanice malwarem, který se snaží nakazit ostatní počítače nebo rozesílá spam.

Protože ne každá změna chování stanice znamená útok či napadení stanice malwarem, měl by systém obsahovat i možnost definovat pravidla, které změny mají být hlášeny a které ne. Podobně by v tomto systému mělo být možné definovat určité oblasti N -rozměrného prostoru, které odpovídají konkrétním typům útoku.

Hlavní výhodou přístupu, založeného na profilech chování jednotlivých stanic, je to, že narozdíl od mnoha jiných metod jsou v tomto případě při detekci nějaké anomálie automaticky známy i IP adresy, které jsou za anomálii odpovědné. Rychlá identifikace zdroje a/nebo cíle útoku je přitom nezbytná pro provedení jakýchkoliv protipatření.

Při vývoji tohoto systému bude nutné vyřešit několik výzkumných otázek. Mezi ně mimo jiné patří:

- Jaké údaje o chování stanic měřit či počítat?
Více údajů může umožnit přesnější detekci různých typů útoků, ale znamená také výpočetně a paměťově náročnější zpracování. Je třeba najít vhodnou množinu měřených statistik.
- Je možné rozšířit flow záznamy o další položky a tím usnadnit či zpřesnit detekci útoků?
Mezi potenciální rozšiřující údaje patří např. různé informace o L7 protokolu, URL extrahované z HTTP požadavků, informace o časových rozestupech paketů v toku a další. Je třeba prozkoumat užitečnost těchto informací a to i vzhledem k náročnosti jejich měření a uchovávání.
- Jaké algoritmy jsou pro detekci nestandardně se chovajících stanic a detekci změn nejvhodnější?
Nabízí se např. možnost použití hierarchického clusterování, jako v práci [5], tato metoda je však pro větší množství dat extrémně výpočetně náročná a v praxi proto nepoužitelná. Je tedy třeba najít jiné efektivnější algoritmy.
Při vývoji metod detekce změn chování budu vycházet ze zkušeností získaných při práci na detekci anomálií v časových řadách (viz předchozí kapitola).

Důležitou součástí systému by mělo být také udržování míry důvěryhodnosti jednotlivých stanic. Pro každou stanicu bude udržována historie incidentů či jen podezřelého chování, na základě čehož bude vypočítána hodnota důvěryhodnosti dané stanice. Pro stanice, které se v minulosti chovaly podezřele, pak například mohou platit přísnější limity detekčních algoritmů či mohou být tyto stanice sledovány podrobněji. Na hodnotu důvěryhodnosti mohou mít vliv i případné jiné detekční systémy.

Data pro vývoj a testování tohoto systému budou brána z kolektoru NetFlow dat, který sbírá data ze všech měřících bodů sítě CESNET. Budou tak k dispozici informace o většině provozu v této síti. Počet aktivních IP adres, jejichž statistiky bude nutno sledovat, se u této sítě pohybuje v řádu statisíců. Základní analýzu takového množství dat zvládne moderní server počítat v reálném čase, některé složitější algoritmy už však příliš pomalé. Z toho důvodu budou prozkoumány i možnosti akcelerace použitých algoritmů pomocí různých paralelních architektur (multi-core, GPU, FPGA).

Přestože vývoj celého systému bude prováděn především na reálných datech přímo na kolektoru, pro prvotní experimenty se statistikami jednotlivých stanic a pro vývoj dílčích součástí (detekčních algoritmů) bude opět využit framework NADEX prezentovaný v kapitole 3.

6 Závěr

V tomto textu byla popsána moje dokončená, aktuálně probíhající i do budoucna plánovaná práce. Nejprve byly popsány některé problémy současného výzkumu metod detekce anomálií a byl představen framework určený pro usnadnění implementace a testování těchto metod. Další část práce se zaměřuje na problém rychlosti detekce a odolnosti monitorovacích systémů proti přetížení při DDoS útoku. Byla zde prezentována myšlenka přesunu detekčního algoritmu z flow kolektoru do exportéru, tedy blíže ke zdroji dat. V současné době probíhá implementace a testování algoritmu určeného pro rychlou detekci DDoS útoků a podobných událostí, který běží jako plugin na NetFlow sondě. V závěru práce byl popsán plánovaný směr mého dalšího výzkumu, tedy cíle mé disertační práce.

Poděkování

Tato práce byla podpořena výzkumným záměrem MSM 0021630528, grantem BUT FIT-S-11-1, grantem VG20102015022 a projektem IT4Innovations Centre of Excellence CZ.1.05/1.1.00/02.0070.

Reference

- [1] Barford, P.; Kline, J.; Plonka, D.; aj.: A signal analysis of network traffic anomalies. In *Proceedings of the 2nd ACM SIGCOMM Workshop on Internet measurement*, IMW '02, New York, NY, USA: ACM, 2002, ISBN 1-58113-603-X, s. 71–82.
- [2] Bartoš, V.: Optimization of Entropy-based Method for Network Anomaly Detection. In *Proceedings of the 17th Conference STUDENT EEICT 2011*, Brno, CZ, 2011.
- [3] Bartoš, V.; Žádník, M.: Framework for comparison of network anomaly detection algorithms. Technická zpráva FIT-TR-2012-02, Fakulta Informačních Technologií, VUT v Brně, Brno, CZ, 2012.
- [4] Bartoš, V.; Žádník, M.: Network Anomaly Detection: Comparison and Real-Time Issues. In *Dependable Networks and Services*, Lecture Notes in Computer Science 7279, Springer, 2012, ISBN 978-3-642-30632-7, s. 118–121.
- [5] Carter, K. M.; Lippmann, R. P.; Boyer, S. W.: Temporally oblivious anomaly detection on large networks using functional peers. In *Proceedings of the 10th annual conference on Internet measurement*, IMC '10, New York, NY, USA: ACM, 2010, ISBN 978-1-4503-0483-2, s. 465–471.
- [6] Nong Ye, Y. Z., Connie Borrer: EWMA techniques for computer intrusion detection through anomalous changes in event intensity. *Qual. Reliab. Engng. Int.*, ročník 18, 2002: s. 443–451.
- [7] Patcha, A.; Park, J.-M.: An overview of anomaly detection techniques: Existing solutions and latest technological trends. *Computer Networks*, ročník 51, August 2007: s. 3448–3470, ISSN 1389-1286.
- [8] Rick Hofstede, A. P.: Real-Time and Resilient Intrusion Detection: A Flow-Based Approach. In *Dependable Networks and Services*, LNCS 7279, Springer Berlin / Heidelberg, 2012, s. 109–112.
- [9] Silveira, F.; Diot, C.; Taft, N.; aj.: ASTUTE: detecting a different class of traffic anomalies. In *Proceedings of the ACM SIGCOMM 2010 conference on SIGCOMM*, SIGCOMM '10, New York, NY, USA: ACM, 2010, ISBN 978-1-4503-0201-2, s. 267–278.
- [10] Zhang, W.; Yang, Q.; Geng, Y.: A Survey of Anomaly Detection Methods in Networks. In *International Symposium on Computer Network and Multimedia Technology*, CNMT 2009, January 2009.