

VILNIAUS UNIVERSITETAS
MATEMATIKOS IR INFORMATIKOS FAKULTETAS
PROGRAMŲ SISTEMŲ STUDIJŲ PROGRAMA

**Varžymosi principais pagrįstų atakų karkasų,
tinkamų kenkėjiško kodo obfuskacijai, analizė**

**Analysis of Adversarial Attack Frameworks for Malware
Obfuscation**

Kursinis darbas

Atliko: 4 kurso 3 grupės studentas

Liudas Kasperavičius

Darbo vadovas: prof. dr. Olga Kurasova

Vilnius – 2024

Turinys

SĄVOKŲ APIBRĖŽIMAI	3
SANTRUMPOS	4
ĮVADAS	5
1. LITERATŪROS APŽVALGA	6
1.1. AE perkeliamumas	6
1.2. Naudojami kenkėjiškų programų požymiai	6
1.2.1. PE formato programų požymiai	7
1.2.2. Baitų lygio požymiai (gali būti ištraukiami iš bet kokio formato failų).....	7
1.3. Perturbacijos	7
1.3.1. Baitų lygio perturbacijos.....	7
1.3.2. Semantinės perturbacijos	8
1.3.3. Kompleksinės perturbacijos	8
1.4. GAN tipo modelių karkasai	9
1.4.1. <i>MalGAN</i>	9
1.4.2. <i>N-gram MalGAN</i>	10
1.4.3. <i>MalFox</i>	11
1.5. Skatinamojo mokymosi tipo modelių karkasai	11
1.5.1. <i>DQEAF</i>	12
1.5.2. <i>MalInfo</i>	13
1.6. Genetinių algoritmų tipo modelių karkasai	13
1.6.1. <i>AIMED</i>	14
1.6.2. <i>GAMMA</i>	14
1.7. Nevalidaus PE formato problema	15
2. KARKASŲ VERTINIMAS	16
3. EKSPERIMENTINIS TYRIMAS	17
3.1. Tyrimo aprašymas	17
3.1.1. GAN implementacija	17
3.1.2. GAN mokymosi duomenys.....	17
3.1.3. Karkaso parametrai	17
3.1.4. Realų kenkėjiškų programų rinkinys	17
3.1.5. Komerciniai detektoriai	17
3.2. Mokymosi etapas	18
3.3. Varžymosi principais pagrįstos atakos prieš komercinius detektorius	18
REZULTATAI IR IŠVADOS	20
ŠALTINIAI	21
PRIEDAI	23
1 priedas. Kompleksinių perturbacijų iliustracijos	23

Sąvokų apibrėžimai

Karkasas (angl. *Framework*). Nurodo specifines technologijas, naudojamus požymius ir perturbacijas, siekiamus tikslus AE generacijai. Skirtas apibrėžti procesą ir įrankius, kuriuos naudojant būtų galima generuoti nurodytų tikslų siekiančius AE. 2, 6, 7, 9, 10, 11, 12, 13, 14, 15, 16, 17

Maišymo Funkcija (angl. *Hash Function*). Tai funkcija $f : \{0,1\}^* \rightarrow \{0,1\}^m$. Naudojama, kai iš begalinės įvesčių erdvės norima gauti fiksuoto dydžio (m) išvestį. 3, 7

Nulinės sumos žaidimas (angl. *Zero-Sum Game*). Dviejų žaidėjų žaidimas, kuriame galimas vienas laimėtojas. Laimėtojo laimėta suma yra lygi pralaimėtojo pralaimėtai sumai. 9

Pėdsakas (angl. *Signature*). Programos struktūros ir požymių santrauka, beveik unikalčiai identifikuojanči programą (pvz. maišymo funkcija). 5

Q-Funkcija (angl. *Q-Function*). $Q : S \times A \rightarrow \mathbb{R}$, čia S – galimų būsenų erdvė (angl. *State Space*), A – galimų veiksmų erdvė (angl. *Action Space*). 12

Sprendimų priėmimo riba (angl. *Decision Boundary*). Paprasčiausiems ML modeliams tai yra kreivė plokštumoje. Sudėtingesniems – daugiadimensiniams modeliams – daugdara (angl. *manifold*). 5

Strategija (angl. *Policy*). Tai funkcija $\pi : S \times A \rightarrow \{0,1\}$, čia S – galimų būsenų erdvė (angl. *State Space*), A – galimų veiksmų erdvė (angl. *Action Space*). Šią funkciją RL modelis „išmoksta“ mokymosi metu. 12

Surogatinis Modelis (angl. *Surrogate Model*). ML modelis, aproksimuojantis kitą ML modelį, kurio parametrai (svoriai) nėra žinomi. 6, 9, 10, 11, 12, 13, 14, 16

Varžymosi principais pagrįstos atakos (angl. *Adversarial Attacks*). Tai atakos, pritaikytos „apgauti“ ML klasifikatorius. 5, 6, 10, 12, 14, 16, 17, 19

Santrumpos

AE	varžymosi principais pagrįstomis atakomis obfuskuoti kenkėjiško kodo pavyzdžiai (<i>angl. adversarial examples</i>) 3, 5–7, 9, 11, 13, 14, 16–18
API	<i>angl. application programming interface</i> 7, 8, 10, 15
CNN	<i>angl. convolutional neural network</i> 11, 12
DI	dirbtinis intelektas (<i>angl. artificial Intelligence</i>) 5, 6
DLL	dinamiškai susieta biblioteka (<i>angl. dyanmic-link library</i>) 7–9, 11, 15
GA	genetiniai algoritmais pagrįstas ML modelis (<i>angl. genetic algorithms</i>) 13, 14
GAN	generatyviniai priešiški tinklai (<i>angl. generative adversarial networks</i>) (generatyviniai varžymosi principais pagrįsti tinklai) 9–11, 14, 17
GBDT	<i>angl. gradient boosted decision trees</i> 11, 12, 14
KNN	<i>angl. K-Nearest Neighbours</i> 11
ML	mašininis mokymasis (<i>angl. machine learning</i>) 3–7, 9–14, 17, 18
NLP	skaitmeninis natūraliosios kalbos apdorojimas (<i>angl. natural language processing</i>) 7
PE	<i>angl. portable executable</i> 6–8, 15, 17
RL	skatinamasis mokymasis (<i>angl. reinforcement learning</i>) 3, 11–14
SVM	<i>angl. support vector machine</i> 11

Įvadas

Pastaraisiais metais kenkėjiškas kodas ir programos kuriamos itin sparčiai (~ 450000 kenkėjiškų programų per dieną 2024 m. AV-TEST¹ duomenimis). Kenkėjiško kodo aptikimo programos, kurios tradiciškai remiasi programų pėdsakais (angl. *signature*), nespėja atnaujinti pėdsakų duomenų bazių pakankamai greitai. Dėl to dirbtinio intelekto (DI), tiksliau Mašininio Mokymosi (ML), naudojimas kenkėjiškų programų ar kenkėjiško kodo aptikimo srityje tapo itin populiarus [DCB⁺21]. Tačiau ML modeliai, nors ir geba aptikti kenkėjiškas programas iš naujų, dar nematytų, duomenų, yra pažeidžiami varžymosi principais pagrįstoms atakoms (angl. *adversarial attacks*) [CSD19; HT17; RSR⁺18; ZHZ⁺22]. Šių atakų principas yra ML modelio – klasifikatoriaus – sprendimų priėmimo ribos (angl. *decision boundary*) radimas – žinant šią ribą pakanka pakeisti kenkėjiškos programos veikimą taip, kad ML modelis priimtų sprendimą klasifikuoti ją kaip nekenksmingą [DCB⁺21]. Žinoma, rasti šią ribą nėra trivialus uždavinys. Mokslinėje literatūroje išskiriami 3 ribos paieškos atvejai [FWL⁺19]:

1. **Baltos dėžės** atvejis: kenkėjiško kodo kūrėjas turi visą informaciją apie ML modelį, t. y. modelio architektūrą, svorius, hiperparametrus.
2. **Juodos dėžės su pasitikėjimo įverčiu** atvejis: kenkėjiško kodo kūrėjas gali tik testuoti modelį – t. y. pateikti programą ir gauti atsakymą. Atsakymo forma – klasifikacija ir tikimybė, kad klasifikacija yra teisinga (pasitikėjimo įvertis).
3. **Juodos dėžės** atvejis: kenkėjiško kodo kūrėjas gali tik testuoti modelį. Atsakymo forma yra tik klasifikacija.

Akivaizdu, jog „juodos dėžės“ atvejis yra sudėtingiausias, bet ir labiausiai atitinka realias sąlygas [AKF⁺18]. Todėl šiame darbe nagrinėjami modeliai, gebantys generuoti varžymosi principais pagrįstų atakų obfusuotus kenkėjiško kodo pavyzdžius (AE) „juodos dėžės“ atvejais.

Tikslas – nustatyti labiausiai tinkantį modelį varžymosi principais pagrįstoms atakoms „juodos dėžės“ atvejais.

Uždaviniai:

1. Apžvelgti kenkėjiško kodo obfuskacijos metodus
2. Nustatyti kriterijus varžymosi principais grįstų atakų modeliams ir juos įvertinti
3. Atlikti eksperimentinį tyrimą su vienu iš įvertintų karkasų ir patikrinti vertinimo rezultatus

¹<https://www.av-test.org/en/statistics/malware>

1. Literatūros apžvalga

Mokslinėje literatūroje kenkėjiško kodo ir programų obfuskacijai daugiausia dėmesio skiriama AE generavimui, dėl paplitusio ML modelių naudojimo kenkėjiško kodo ir programų detekcijos srityje [AKF⁺18]. Dažniausia atakų platforma – operacinė sistema „Windows“ ir jos naudojami PE formato failai (*VirusTotal*² duomenimis).

Kenkėjiškų programų analizė gali būti išskiriama į **statinę analizę** ir **elgesio** (*angl. behavior*) **analizę**. Statinės analizės detektoriai yra labiau paplitę ir prieinami, kadangi požymių išskyrimas iš tam tikro iš anksto žinomo formato failų ir jų klasifikavimas yra sąlyginai nesudėtinga užduotis. Tuo tarpu elgesio analizės detektoriai dažniausiai neprieinami eiliniams vartotojams, o diegiami korporacijose [RSR⁺18]. Jų veikimas pagrįstas nuolatiniu sistemos stebėjimu ir anomalijų detekcija, tad yra kartu ir sudėtingesnė, ir reikalaujanti daugiau resursų, nei statinė analizė, užduotis.

Mokslinėje literatūroje nagrinėjant specifinį būdą generuoti AE, apibrėžiamas varžymosi principais pagrįstų atakų karkasas. Karkasą sudaro

- Surogatinio modelio tipas ir architektūra
- ML modelio (jei naudojamas) tipas ir architektūra
- ML modelio naudojami požymiai (žr. 1.2)
- ML modelio naudojamos perturbacijos (žr. 1.3)

1.1. AE perkeliamumas

Perkeliamumas DI modeliuose dažniausiai suprantamas kaip žinių perkeliamumas (*angl. knowledge transferability*). Tačiau tiriant varžymosi principais pagrįstas atakas buvo pastebėtas ir AE perkeliamumas (*angl. adversarial sample transferability*). Tai reiškia, jog gebant sukurti AE, kuriuos neteisingai klasifikuoja vienas modelis, tikėtina, jog kitas (tos pačios paskirties) modelis taip pat klasifikuos AE neteisingai. Be to, nustatyta, jog AE perkeliamumo savybė galioja skirtingų architektūrų modeliams [DCB⁺21]. Šia savybe remiasi visos „juodos dėžės“ atvejams pritaikytos atakos.

1.2. Naudojami kenkėjiškų programų požymiai

Varžymosi principais pagrįstos atakos taikosi į ML modeliais paremtus kenkėjiškų programų detektorius. Šie detektoriai yra klasifikatoriai – pateiktį (programą) klasifikuoja kaip kenkėjišką (*angl. malicious*) arba nekenkėjišką (*angl. benign*). Kadangi programos nėra fiksuoto dydžio, klasifikatoriai remiasi programų požymiais, kurie gaunami atliekant požymių ištraukimą (*angl. feature extraction*). Laikoma, jog „juodos dėžės“ atvejais sužinoti kokius tiksliai požymius vertina kenkėjiškų programų detektorius yra neįmanoma, tad karkasų apibrėžimuose, priklausomai nuo jų specifikos ir tikslų, neretai pateikiami jų vertinami programų požymiai. Šioje sekcijoje išskiriami ir klasifikuojami mokslinėje literatūroje minimi požymiai.

²<https://www.virustotal.com/gui/stats>

1.2.1. PE formato programų požymiai

- **DLL vardai (arba API vardai [HT17])** [ZCY⁺24]. PE faile turi būti nurodyti visi naudojami DLL ir jų API. Prieš pradedant treniruoti ML modelį, atliekama visų turimų programų analizė ir nustatoma visų naudojamų DLL ar jų API aibė D . Tarkime $|D| = n$. Tuomet, požymių vektorius programai, naudojančiai $X \subseteq D$ DLL, bus n -matis dvejetainis vektorius, kurio i -asis elementas yra
$$\begin{cases} 0, & \text{jei } D_i \notin X, \\ 1, & \text{jei } D_i \in X \end{cases} \quad \text{čia } D_i - i\text{-asis } D \text{ elementas.}$$
- **PE metaduomenys** [AKF⁺18]. Tai visi PE formato faile esantys metaduomenys, tokie, kaip sekcijų pavadinimai, sekcijų dydžiai, „ImportTable“ ir „ExportTable“ metaduomenys ir kt. Formuojant požymių vektorių skaičiuojama metaduomenų maišymo funkcija.

1.2.2. Baitų lygio požymiai (gali būti ištraukiami iš bet kokio formato failų)

- **Prasmingų žodžių (angl. Strings) kiekis** [AKF⁺18]. Prasmingus žodžius suprantame kaip turinčius prasmę žmogui (*angl. human readable*). Tai gali būti URL, failų keliai (*angl. file paths*) ar registro raktų pavadinimai. Kadangi prasmingų žodžių kiekis tėra vienas skaičius, požymių vektorius dažniausiai formuojamas prijungiant ir kitus požymius.
- **Baitų/entropijos histograma** [SB15]. Specifinis metodas, užkoduojantis dažniausiai pasikartojančias baitų ir entropijos poras n dimensijų vektoriumi.
- **n -gramos** [ZZY⁺22]. Dažniausiai sutinkamos skaitmeniniame natūraliosios kalbos apdorojime (NLP). Tai yra n žodžių junginiai, arba, sukompiliuotų programų apdorojimo kontekste, n baitų junginiai. Nustatant požymių vektorių, visos n -gramos surikiuojamos pagal pasikartojimą programoje mažėjimo tvarka („populiariausios“ viršuje). Iš pirmų m reikšmių sudaromas m -matis vektorius – tai ir yra požymių vektorius.

1.3. Perturbacijos

Perturbacijos – tai pagrindinis obfuskacijos metodas AE kūrimui. Perturbacijų tikslas yra pakeisti kenkėjiškos programos veikimą išsaugant originalų funkcionalumą. Perturbacijos gali būti sudėtingos ir apimti visą programą (pvz. visos programos užšifravimas ir pridėjimas prie kitos programos), semantinės (pvz. tam tikrų mašininio kodo instrukcijų keitimas į ekvivalentų rezultatą pasiekiančias) arba baitų lygio (pvz. nulinių baitų pridėjimas programos gale) [HT17]. Perturbacijų parinkimas įeina į karkaso apibrėžimą. Šioje sekcijoje aptariamos mokslinėje literatūroje minimos perturbacijos.

1.3.1. Baitų lygio perturbacijos

- **ARBE (Append Random Bytes at the End)** [FWL⁺19]. PE formato failo gale pridedami atsitiktiniai baitai.
- **ARI (Append Random Import)** [FWL⁺19]. PE formato failo „ImportAddressTable“ lentelėje pridedama atsitiktinai pavadinta biblioteka su atsitiktinai pavadinta funkcija.

- **ARS (Append Randomly named Section)** [FWL⁺19]. PE formato failo „SectionTable“ lentelėje pridedamos atsitiktinės sekcijos (sekcijos ir jų tipai yra apibrėžti PE formate).
- **RS (Remove Signature)** [FWL⁺19]. Sertifikato pašalinimas iš PE formato failo „CertificateTable“ lentelės.
- **Naujas įėjties taškas** [AKF⁺18]. Prasidėjus programai, iškart peršokama nuo naujo įėjties taško į originalųjį.
- **Header Fields** [DCB⁺21]. PE formato failo „PE Header“ ir „Optional Header“ dalių specifinių laukų keitimas (pvz. sekcijos pavadinimo keitimas [AKF⁺18]).
- **Partial DOS** [DCB⁺21]. PE formato failo „DOS Header“ dalies pirmi 58 baitai po „MZ“ skaičiaus yra nenaudojami moderniose operacinėse sistemose, tad juos galima keisti.
- **Slack Space** [DCB⁺21]. Dėl PE formato specifikos, kiekviena nauja sekcija turi prasidėti tam tikro skaičiaus, nurodyto „PE Header“ dalyje, kartotiniu nuo pradžios. Kompiliatoriai šį reikalavimą išpildo sekcijų gale pridėdami tiek nulinių baitų, kiek reikia teisingam sulygiavimui pasiekti. Būtent ši nulinių baitų erdvė gali būti keičiama be jokios įtakos originaliai programai.
- **Padding** [DCB⁺21]. Nulinių baitų pridėjimas failo gale.
- **Full DOS** [DCB⁺21]. Perturbacijos esmė tokia pat, kaip ir „Partial DOS“, tik naudojami visi „DOS“ dalies baitai, išskyrus „MZ“ ir „PE Offset“ („Partial DOS“ manipuliacijoms naudoja tik dalį tarp „MZ“ ir „PE Offset“).
- **Extend** [DCB⁺21]. Pakeičiama PE formato faile „DOS“ dalyje esanti „PE Offset“ reikšmė į didesnę³. Taip padidinama (išplečiama) visa „DOS“ dalis. Tolesnis perturbacijos principas yra toks pat, kaip ir „Full DOS“.
- **Shift** [DCB⁺21]. PE formato failuose kiekvienas sekcijos blokas prasideda su sekcijos vieta nuo pradžios (*angl. offset*). Tarkime ši reikšmė yra S . Sekcijos kodas pradedamas vykdyti tik nuo adreso $P + S$, kur P – programos pradžios adresas. Vadinasi, padidinus³ S per n , atsiranda n baitų laisvos vietos iki sekcijos pradžios, kurią galima keisti be jokios įtakos programos veikimui.

1.3.2. Semantinės perturbacijos

- **Nereikalingų DLL/API vardų požymių pridėjimas** [HT17]. PE formato faile „ImportTable“ lentelėje pridedami originalios programos nenaudojami DLL/API vardai.
- **Binary Rewriting** [DCB⁺21]. Semantinis instrukcijų perrašymas. Pavyzdžiui, $A + B$ instrukcijos pakeitimas į $A - (-B)$.

1.3.3. Kompleksinės perturbacijos

- **Obfusmal** [ZCY⁺24]. Užšifruojama originalios programos kodo sekcija. Sukuriama ir originalios programos gale pridedama programa *Shell.dll*, kurioje laikomas atšifravimo raktas, originalios programos kodo sekcijos adresas ir dydis. Be to, *Shell.dll* geba atšifruoti

³šios reikšmės padidinimas reiškia visos failo struktūros keitimą („DOS“ dalis yra failo pradžioje). Būtina pakeisti visų sekcijų vietas nuo pradžios (*angl. offset*) jų metaduomenyse.

originalios programos kodo sekciją ir jai perduoti kontrolę. *Shell.dll* pridedama prie naudojamų DLL, o programos pradžios taškas nustatomas į *Shell.dll* pradžios tašką. Iliustracija pateikiama 4-ame pav.

- **Stealmal** [ZCY⁺24]. Visa originali programa užšifruojama ir pridedama prie programos *Shell.exe* galo. *Shell.exe* geba atšifruoti originalią programą ir perduoti jai kontrolę. Iliustracija pateikiama 4-ame pav.
- **Hollowmal** [ZCY⁺24]. Užšifruojama visa originali programa. Ji pridedama prie kurios nors nekenksmingos programos galo. Prie šio junginio galo pridedama *Hollow.dll* programa, kurios veikiamas panašus į *Shell.exe* iš „Stealmal“. Viso junginio pradžios taškas nustatomas į *Hollowmal.dll* pradžios tašką. Iliustracija pateikiama 4-ame pav.

1.4. GAN tipo modelių karkasai

GAN modelių karkasai paremti Generatyviniais Priešiškais Tinklais (*angl. generative adversarial networks*), kurių veikimo principas yra du neuroniniai tinklai (generatorius ir diskriminatorius), žaidžiantys nulinės sumos žaidimą (*angl. zero-sum game*) [CDH⁺16]. Kenkėjiško kodo obfuskacijos kontekste ir ypač „juodos dėžės“ atvejais, diskriminatorius atlieka surogatinio modelio vaidmenį. Bendras GAN modelių mokymosi etapas yra tokia seka:

1. Generatorius, naudodamas požymių vektorių ir tokios pačios dimensijos „triukšmo“ (*angl. noise*) vektorių, sugeneruoja perturbacijas.
2. Originali kenkėjiška programa modifikuojama pagal perturbacijas (sukuriamas AE).
3. Diskriminatorius klasifikuoja sugeneruotą AE (kenkėjiškas / nekenkėjiškas). Pagal klasifikacijos rezultatą skaičiuojamos diskriminatoriaus ir generatoriaus nuostolių funkcijos reikšmės (diskriminatoriaus nuostolių funkcijos reikšmė priklauso nuo tikro detektoriaus klasifikacijos).
4. Visa seka kartojama nustatytą kiekį kartų.

[HT17; ZCY⁺24; ZZY⁺22].

1.4.1. MalGAN

Tai vienas iš pirmųjų ir populiariausių GAN tipo modelių karkasų. *MalGAN* karkaso [HT17] tikslas ir apibrėžimas pateikiami 1-oje lentelėje.

Tikslas	Efektyviai išvengti AE aptikimo, kai ML kenkėjiškų programų detektoriaus implementacija nežinoma („juodos dėžės“ atvejis)
Surogatinis modelis	Daugiasluoksnis tiesioginio sklaidimo neuroninis tinklas – klasifikatorius. Įvestis – programos požymių vektorius. Išvestis – klasifikacija į kenksmingą arba nekenksmingą. Šis tinklas taip pat naudojamas kaip diskriminatorius GAN architektūroje.

ML modelis	Daugiasluoksnis tiesioginio sklaidimo neuroninis tinklas. Įvestis – programos požymių vektorius ir tokios pačios dimensijos „triukšmo“ vektorius. Išvestis – modifikuotas požymių vektorius. Šis tinklas naudojamas kaip generatorius GAN architektūroje.
Požymiai	<ul style="list-style-type: none"> • <i>MalGan</i> straipsnyje naudojami tik API vardų požymiai, patenkantys į PE formato programų požymių kategoriją (žr. 1.2.1), tačiau autoriai nurodo, jog gali būti naudojami bet kokie požymiai⁴.
Perturbacijos	<ul style="list-style-type: none"> • Semantinės perturbacijos (1.3.2) – nereikalingų API vardų požymių pridėjimas

1 lentelė. *MalGAN* karkasas

1.4.2. *N-gram MalGAN*

Šis karkasas remiasi *MalGAN* (1.4.1) karkasu ir siekia jį pagerinti. *N-gram MalGAN* karkaso [ZZY⁺22] tikslas ir apibrėžimas pateikiami 2-oje lentelėje.

Tikslas	Supaprastinti, pagreitinti ir pagerinti varžymosi principais pagrįstas atakas. Pašalinti prielaidas ⁴ apie detektorių „juodos dėžės“ atvejais.
Surogatinis modelis	Surogatinio modelio veikimas ir architektūra tokia pati, kaip ir <i>MalGAN</i> (1.4.1)
ML modelis	Pagrindinio modelio veikimas ir architektūra labai panašūs į <i>MalGAN</i> (1.4.1), tačiau norėdami stabilizuoti mokymosi procesą, autoriai siūlo nenaudoti „triukšmo“ vektoriaus. Vietoje to, generatoriaus išvestis (n -matis vektorius) modifikuojama nekeičiant pirmų m dimensijų, o kitas $n - m$ pakeičiant nekenksmingų programų požymiais.
Požymiai	<ul style="list-style-type: none"> • Baitų lygio požymiai (1.2.2) – n-gramos.

⁴ autoriai nagrinėja „juodos dėžės“ atvejį su prielaida, jog detektoriaus naudojami požymiai yra žinomi.

Perturbacijos	<p>Autoriai neatliko eksperimentų su perturbuotomis programomis, tačiau pažymi, jog norint gauti sugeneruotus požymių vektorius užtenka pridėti reikiamus baitus programos gale. Tai atitinka 1.3.1 apibrėžtą</p> <ul style="list-style-type: none"> baitų lygio perturbaciją <i>ARBE</i>, tik šiuo atveju pridedami baitai nebūtų atsitiktiniai, o norimos n-gramos.
----------------------	---

2 lentelė. *N-gram MalGAN* karkasas

1.4.3. *MalFox*

MalFox taip pat remiasi *MalGAN* (1.4.1), tačiau siekia kurti AE realiomis sąlygomis, dėl to atlieka esminius pakeitimus. *MalFox* karkaso [ZCY⁺24] tikslas ir apibrėžimas pateikiami 3-oje lentelėje.

Tikslas	Generuoti AE, kurių neaptiktų komerciniai detektoriai (prieš tai aptarti karkasai eksperimentams kaip nepriklausomą detektorių naudojo tokius ML modelius, kaip SVM, KNN, GBDT ir kt., bet ne komercinius detektorius). Šio karkaso detektorius yra <i>VirusTotal</i> (viešai prieinama paslauga, agreguojanti virš 70 komercinių kenkėjiškų programų detektorių).
Surogatinis modelis	Surogatinis modelis, kaip ir kituose GAN tipo modelių karkasuose, naudojamas kaip diskriminatorius. Įvestis – perturbuota programa. Išvestis – klasifikacija į kenksmingą arba nekenksmingą. Implementacija – konvoliucinis neuroninis tinklas (CNN).
ML modelis	Standartinis GAN generatorius, požymių vektorių sujungiantis su „triukšmo“ vektoriumi. Implementacija – konvoliucinis neuroninis tinklas (CNN).
Požymiai	<ul style="list-style-type: none"> PE formato programų požymiai (1.2.1) – DLL vardai.
Perturbacijos	<ul style="list-style-type: none"> Visos kompleksinės perturbacijos (1.3.3)

3 lentelė. *MalFox* karkasas

1.5. Skatinamojo mokymosi tipo modelių karkasai

Skatinamojo mokymosi (RL) modeliai susideda iš agento ir aplinkos. Aplinka susideda iš informatyvių požymių ištraukimo metodo (*angl. feature extraction*) ir kenkėjiškų programų detektoriaus. Šiuo atveju aplinkos būsenų erdvė S yra požymių vektorių erdvė. Agentas – tai

algoritmas ar neuroninis tinklas, kurio tikslas yra surasti optimalią strategiją (angl. *policy*). Šiuo atveju strategijos veiksmų erdvė A susideda iš perturbacijų (žr. 1.3) [FWL⁺19]. Bendras RL modelių mokymosi etapas yra tokia seka:

1. Agentas, naudodamas dabartinę aplinkos būseną ir praeito veiksmo atlygį (angl. *reward*), parenka sekantį veiksmą iš galimų veiksmų aibės ir taiko mokymosi algoritmą (algoritmas priklauso nuo agento implementacijos).
2. Atliekamas veiksmas – perturbuojama programa arba požymių vektorius (priklauso nuo karkaso).
3. Gaunami aplinkos kitimo įverčiai – nauja būsena ir atlygis, skaičiuojamas pagal detektoriaus klasifikacijos rezultatą.
4. Seka kartojama tol, kol agentas nelaiko strategijos optimalia arba nustatytą kiekį kartų.

[FWL⁺19; ZHZ⁺22]

1.5.1. DQEAF

Šis karkasas taiko gilųjį skatinamąjį mokymąsi, kai agentas implementuojamas kaip gilusis neuroninis tinklas. DQEAF karkaso [FWL⁺19] tikslas ir apibrėžimas pateikiami 4-oje lentelėje.

Tikslas	Parodyti, jog ML kenkėjiškų programų detektoriai, ypač modeliai, treniruoti prižiūrimu mokymusi, yra pažeidžiami varžymosi principais pagrįstoms atakoms
Surogatinis modelis	RL karkasuose nenaudojami surogatiniai modeliai. Kaip „juodos dėžės“ detektorius pasirinktas GBDT modelis.
ML modelis	Agentas implementuotas kaip gilusis Q -tinklas (CNN praplėtimas, kai tinklas naudojamas kaip Q -funkcijos aproksimacija). Taip pat taikomas prioritetizuotas patirčių pakartojimo metodas (angl. <i>prioritized experience replay</i>), kuomet agentas treniruojamas tik su aukštą atlygį gavusiais perėjimais ($S \times A$).
Požymiai	Požymių vektorius taip pat apibrėžia visų būsenų erdvę S . Šiuo atveju $S = \mathbb{R}^{513}$. <ul style="list-style-type: none"> • Baitų lygio požymiai (1.2.2) – baitų/entropijos histograma.
Perturbacijos	Perturbacijos apibrėžia visų galimų agento veiksmų erdvę A . Šiuo atveju $A = \{0,1\}^4$. <ul style="list-style-type: none"> • Baitų lygio perturbacijos (1.3.1) <ul style="list-style-type: none"> – ARBE – ARI – ARS – RS

4 lentelė. DQEAF karkasas

1.5.2. *MalInfo*

MalInfo remiasi *MalFox* (1.4.3). *MalInfo* karkaso [ZHZ⁺22] tikslas ir apibrėžimas pateikiami 5-oje lentelėje.

Tikslas	Surasti optimalią obfuskacijos strategiją konkrečiai programai, pagal kurią sukurtas AE nebūtų aptiktas komercinių kenkėjiškų programų detektorių.
Surogatinis modelis	RL surogatinis modelis nenaudojamas. „Juodos dėžės“ detektoriumi pasirinkti komerciniai detektoriai („VirusTotal“)
ML modelis	Agentas implementuotas kaip klasikiniai ML algoritmai (konkrečiai dinaminis programavimas ir skirtumų laike (<i>angl. temporal difference</i>) algoritmas).
Požymiai	Agentas nėra neuroninis tinklas ir požymių iš programos netraukia. Agentas mokosi tik iš perėjimų, o būsenų erdvė S yra originali programa ir perturbuoti jos variantai. Teoriškai perturbuotų programos variantų galėtų būti be galo daug, tuomet $S = A^\infty$, $ S = \aleph_0$, tačiau autoriai nurodo, jog daugiau nei 3 sluoksniai kompleksinių perturbacijų reikšmingai paveikia programos veikimo laiką, o tai gali „sukelti įtarimų“ komerciniams detektoriams. Todėl pasirinkta $S = A^3$
Perturbacijos	$A = \{0,1,2,3\}$ <ul style="list-style-type: none"> • „Null“ perturbacija – naudinga tik formaliam pilnumui (atitinka nulinį A veiksmą) • Visos kompleksinės perturbacijos (1.3.3), t. y. tokios pačios, kaip ir <i>MalFox</i> (1.4.3) karkaso.

5 lentelė. *MalInfo* karkasas

1.6. Genetinių algoritmų tipo modelių karkasai

Genetiniai algoritmai (GA) yra viena seniausių mašininio mokymosi ML apraiškų; jų veikimas paremtas evoliucija [CSD19]. Kenkėjiškų programų obfuskacijai AE generavimas taikant GA yra tokia seka:

1. Sukuriama pradinė populiacija (perturbacijos metodai pradinei populiacijai priklauso nuo karkaso).
2. Atliekamas tinkamumo (*angl. fitness*) vertinimas.
3. Atliekama selekcija – dažniausiai pasirenkami geriausiai įvertinti populiacijos AE, tačiau galimos ir kitos selekcijos strategijos.
4. Atliekamas selekcijos atrinktų AE kryžminimas (po 2) taip sukuriant naują AE, turintį po dalį genų iš abiejų kryžmintų AE.
5. Tam tikrai daliai AE atliekama dalies genų mutacija.

6. Vertinama, ar sugeneruoti AE atitinka kriterijus (vertina detektorius)
7. Jei kriterijai nėra tenkinami, seka kartojama nuo 2-o žingsnio.

[YPT22]

1.6.1. AIMED

AIMED karkaso [CSD19] tikslas ir apibrėžimas pateikiami 6-oje lentelėje.

Tikslas	AE generavimo greičio padidinimas ir modelių kompleksiskumo sumažinimas, lyginant su GAN ir RL tipo modelių karkasais.
Surogatinis modelis	Surogatinis modelis nenaudojamas. Naudojami „juodos dėžės“ detektoriai yra 3 komerciniai (<i>Kaspersky, ESET, Sophos</i>) ir vienas ML modelis – GBDT.
ML modelis	Klasikinis GA modelis – veikimas visiškai atitinką bendrą seką. Tinkamumo (<i>angl. fitness</i>) vertinamas remiasi AE požymių vektoriaus panašumu į originalios programos požymių vektorius (kuo mažiau panašūs, tuo tinkamumo įvertinimas didesnis).
Požymiai	<ul style="list-style-type: none"> Baitų lygio požymiai (1.2.2) – atskiras n-gramų atvejis, kai $n = 1$
Perturbacijos	<ul style="list-style-type: none"> Baitų lygio perturbacijos⁵ (1.3.1).

6 lentelė. AIMED karkasas

1.6.2. GAMMA

GAMMA karkaso [DBL⁺21] tikslas ir apibrėžimas pateikiami 7-oje lentelėje.

Tikslas	Efektyvus (neaptikimo šansų didinimas naudojant perturbacijas, paremtas nekenksmingomis programomis) varžymosi principais pagrįstų atakų kūrimas.
Surogatinis modelis	Surogatinis modelis nenaudojamas. GBDT ir <i>MalConv</i> pasirinkti kaip „juodos dėžės“ detektoriai.
ML modelis	Pagrindinė modelio idėja yra požymių ištraukimas iš nekenksmingų programų ir jų pridėjimas, naudojant tam pritaikytas perturbacijas, į kenksmingas programas kiekvienos populiacijos generavimo metu. Tinkamumo (<i>angl. fitness</i>) ir kriterijų vertinimas atliekamas naudojant detektorius ir pridėtų požymių dydį baitais (norima pridėti kuo mažiau požymių).

⁵ autoriai rėmėsi perturbacijomis, aprašytomis [AKF⁺18]

Požymiai	<ul style="list-style-type: none"> • PE formato programų požymiai (1.2.1). • Kodas sekcijose (nestandartinis požymis).
Perturbacijos	<ul style="list-style-type: none"> • Visos baitų lygio perturbacijos (1.3.1), gebančios pridėti baitus. • Autoriai pažymi, jog gali būti naudojama ir DLL / API vardų pridėjimo semantinė perturbacija (1.3.2).

7 lentelė. GAMMA karkasas

1.7. Nevalidaus PE formato problema

Anderson et al., atlikdami eksperimentus su funkcionalumą išlaikančiomis perturbacijomis PE formato failams, pastebėjo, jog ne visais atvejais perturbuotos programos veikia teisingai. Dėl *Windows* operacinės sistemos PE formato failų interpretavimo ir paleidimo specifikos, programas įmanoma parašyti tokiu būdu, jog pakeitus kodo ar kitų sekcijų turinį nekeičiant originalių mašininio kodo instrukcijų, programa neveiktų. Techniškai, programų rašymas tokiu būdu pažeidžia patį PE formato standartą, tačiau šią praktiką neretai naudoja kenkėjiškų programų autoriai [AKF⁺18].

Norint visiškai išvengti nevalidaus PE formato problemos tenka taikyti perturbacijas, nekeičiančias originalių programų, o taikančias kitokius obfuskacijos metodus. Iš 1.3 sekcijoje aptartų perturbacijų, tokias sąlygas atitinka tik 2 kompleksinės perturbacijos (1.3.3) – *Stealmal* ir *Hollowmal*.

2. Karkasų vertinimas

Karkasų vertinimui apibrėšime kriterijus, kurie padės objektyviai išrinkti realioms varžymosi principais pagrįstoms atakoms tinkančius karkasus. Kriterijus laikysime esant dviejų tipų: **kokybiniai** (karkasas atitinka kriterijų arba ne) ir **kiekybiniai** (karkasas atitinka kriterijų dalinai: 0% – visiškai neatitinka, 100% – visiškai atitinka).

Kriterijai (surikiuoti pagal svarbą mažėjančia tvarka):

K-1. Karkasas pritaikytas „juodos dėžės“ atvejams (**kokybinis**).

K-2. Karkasas pritaikytas komerciniams detektoriams (**kokybinis**).

K-3. Karkasas užtikrina realistišką atakos efektyvumo įvertinimą – naudoja surogatinį modelį (**kokybinis**).

K-4. Karkasas užtikrina atakos efektyvumą (**kiekybinis**).

Karkasas	K-1	K-2	K-3	K-4
<i>MalFox</i> (1.4.3)	✓	✓	✓	56,00%
<i>MalInfo</i> (1.5.2)	✓	✓	✗	59,40%
<i>AIMED</i> (1.6.1)	✓	✓	✗	47,98%
<i>N-gram MalGAN</i> (1.4.2)	✓	✗	✓	88,58%
<i>MalGAN</i> (1.4.1)	✓	✗	✓	81,00%
<i>GAMMA</i> (1.6.2)	✓	✗	✗	53,00%
<i>DQEAF</i> (1.5.1)	✓	✗	✗	46,56%

8 lentelė. Karkasų vertinimas pagal kriterijus. Karkasai surikiuoti pagal įvertinimą mažėjančia tvarka.

Kriterijų vertinimo lentelėje pastebime, jog dviejų aukščiausiai pagal K-4 kriterijų įvertintų karkasų (*N-gram MalGAN* (1.4.2) ir *MalGAN* (1.4.1)) įvertinimai žymiai aukštesni, nei vidurkis (nuokrypis per 26,79% ir 19,21% atitinkamai), tačiau bendras jų įvertinimas pakankamai žemas (4 ir 5 vietos). Tai lemia K-2 kriterijaus įvertinimas. Kadangi vertinimo tikslas yra nustatyti realioms varžymosi principais pagrįstoms atakoms tinkančius karkasus, negalime teigti, jog karkasai, puikiai konstruojantys AE akademiniams modeliams, gebės konstruoti tokios pačios kokybės AE komerciniams modeliams. Tai ir parodysime 3-ame skyriuje.

3. Eksperimentinis tyrimas

3.1. Tyrimo aprašymas

Naudodami *MalGAN* (1.4.1) karkasą paruošime varžymosi principais pagrįstas atakas ir apskaičiuosime jų efektyvumą prieš komercinius detektorius.

3.1.1. GAN implementacija

Naudosime viešai *Github* platformoje paskelbtą⁶ *MalGAN* karkaso ML modelio implementaciją.

3.1.2. GAN mokymosi duomenys

Mokymuisi naudosime *EMBER* [AR18] duomenų rinkinį. Treniravimo aibė susideda iš 2000 programų požymių (*API* vardų iš *ImportTable* PE formato failo sekcijos). Aibė subalansuota – joje yra po 1000 kenkėjiškų ir nekenkėjiškų programų požymių. Mokymosi etapui taip pat taikysime treniravimo / validacijos skaidinį santykiu 4 : 1.

3.1.3. Karkaso parametrai

- Dvejjetainio požymių vektoriaus M dimensija $d_M = 10000$
- „Triukšmo“ vektoriaus Z dimensija $d_Z = 100$
- Neuronų skaičius generatoriaus „paslėptuose“ sluoksniuose: 100,256,512,1024
- Neuronų skaičius diskriminatoriaus „paslėptuose“ sluoksniuose: 256,256
- Diskriminatoriaus mokymosi greitis $\alpha_D = 4 \cdot 10^{-4}$
- Generatoriaus mokymosi greitis $\alpha_G = 10^{-6}$

3.1.4. Realų kenkėjiškų programų rinkinys

Atakoms naudosime realias kenkėjiškas PE formato programas iš *VirusShare*⁷. Duomenų rinkinį sudaro 100 kenkėjiškų programų.

3.1.5. Komerciniai detektoriai

Tiek originalias, tiek modifikuotas (AE) kenkėjiškas programas patikrinsime su *VirusTotal*⁸ – paslauga, kuri agreguoja daugiau nei 70 komercinių kenkėjiškų programų detektorių.

⁶<https://github.com/ZaydH/MalwareGAN>

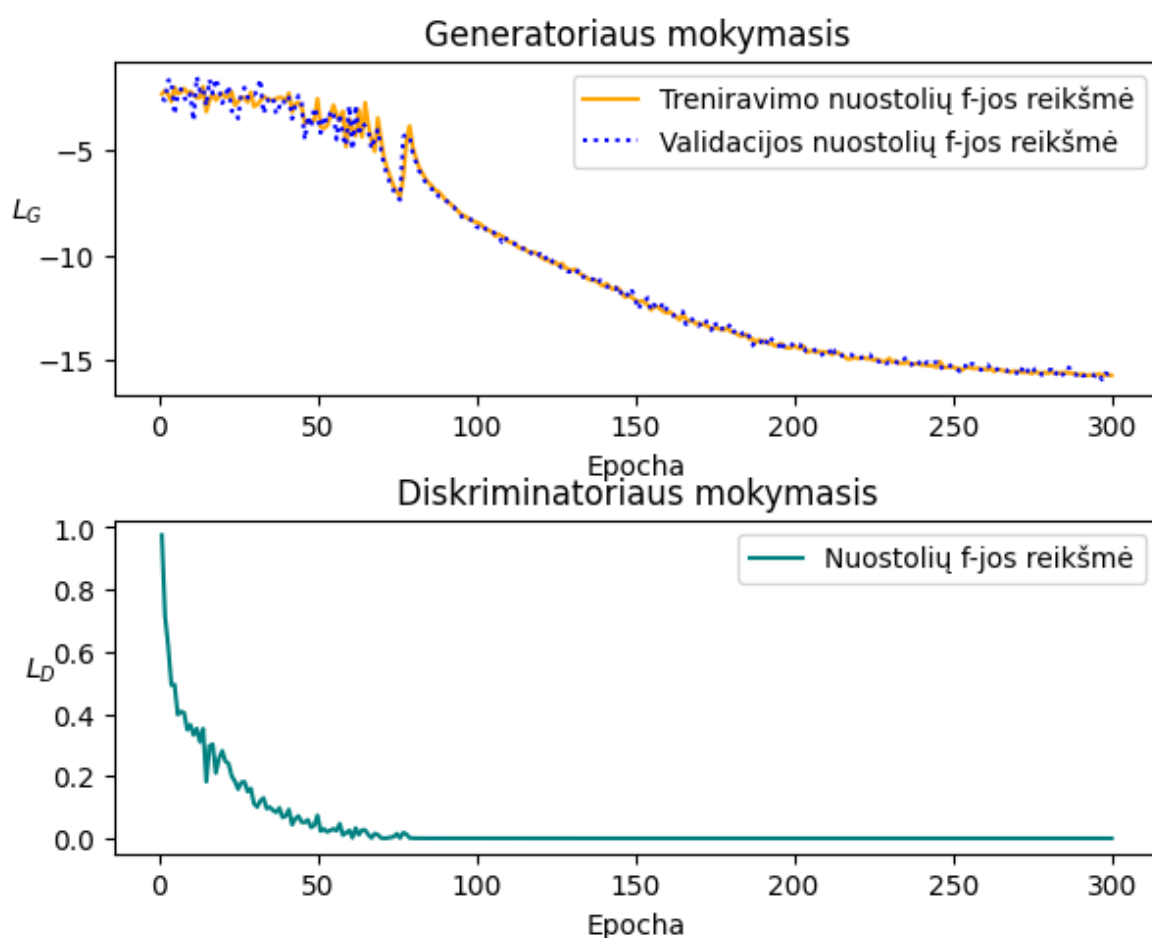
⁷<https://virusshare.com/>

⁸<https://www.virustotal.com>

3.2. Mokymosi etapas

Specifinis *MalGAN* (1.4.1) mokymosi etapas yra tokia seka:

1. Išmokomas „juodos dėžės“ modelis (šiuo atveju pasirinktas „atsitiktinio miško“ (*angl. random forest*) modelis).
2. Mokoma pagal bendrą seką (aptartą 1.4), kurioje:
 - Praleidžiamas originalios programos perturbacijos žingsnis (diskriminatoriui iškart perduodamas AE požymių vektorius).
 - Generatorius ir diskriminatorius mokomi kartu, kadangi diskriminatorius turi naudoti generatoriaus sugeneruotus AE, o generatoriaus nuostolių funkcija priklauso nuo diskriminatoriaus išvesties. Šis žingsnis pavaizduotas 1-ame pav.

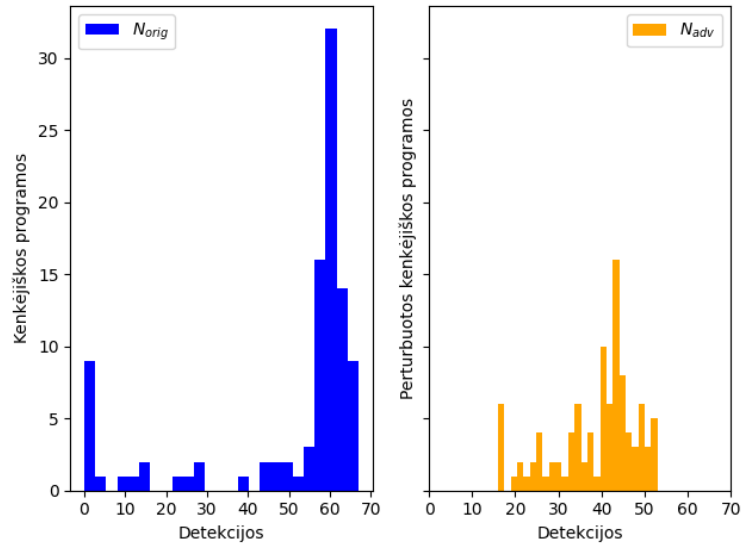


1 pav. *MalGAN* generatoriaus ir diskriminatoriaus nuostolių funkcijų reikšmės mokymosi etape

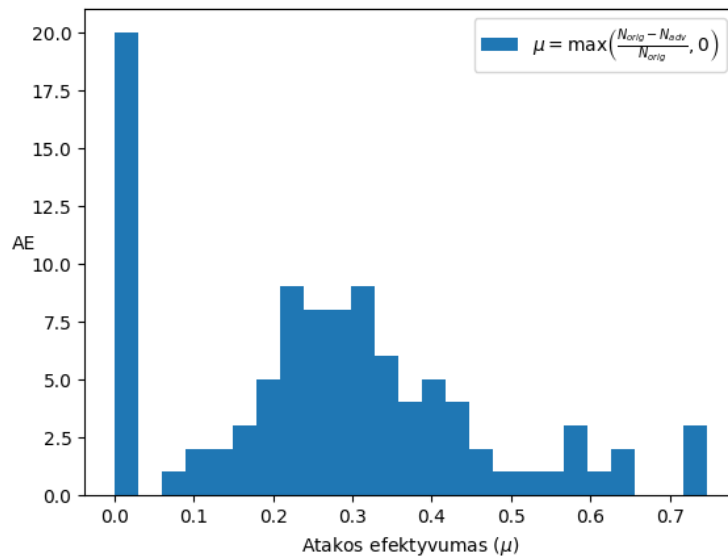
3.3. Varžymosi principais pagrįstos atakos prieš komercinius detektorius

Naudodami mokymosi etape ištreniruotą ML modelį ir 3.1.4 aptartą duomenų rinkinį, sugeneruosime AE ir palyginsime, kiek komercinių detektorių aptinka originalias programas (N_{orig}) ir kiek obfuskuotas pagal *MalGAN* karkasą (N_{adv}).

Atakų efektyvumą μ skaičiuosime remiantis Zhong et al. pasiūlyta formule $\mu' = \frac{N_{orig} - N_{adv}}{N_{orig}}$ [ZCY⁺24], laikydami, jog atakos efektyvumas neneigiamas $\mu = \max(\mu', 0)$.



2 pav. Originalių (N_{orig}) ir perturbuotų (N_{adv}) kenkėjiškų programų detekcijų pasiskirstymas



3 pav. Varžymosi principais pagrįstų atakų efektyvumo pasiskirstymas

Efektyvumo vidurkis: 26,85%

Standartinis nuokrypis: 18,97%

Rezultatai ir išvados

Rezultatų ir išvadų dalyje turi būti aiškiai išdėstomi pagrindiniai darbo rezultatai (kažkas išanalizuota, kažkas sukurta, kažkas įdiegta) ir pateikiamos išvados (daromi nagrinėtų problemų sprendimo metodų palyginimai, teikiamos rekomendacijos, akcentuojamos naujovės).

Šaltiniai

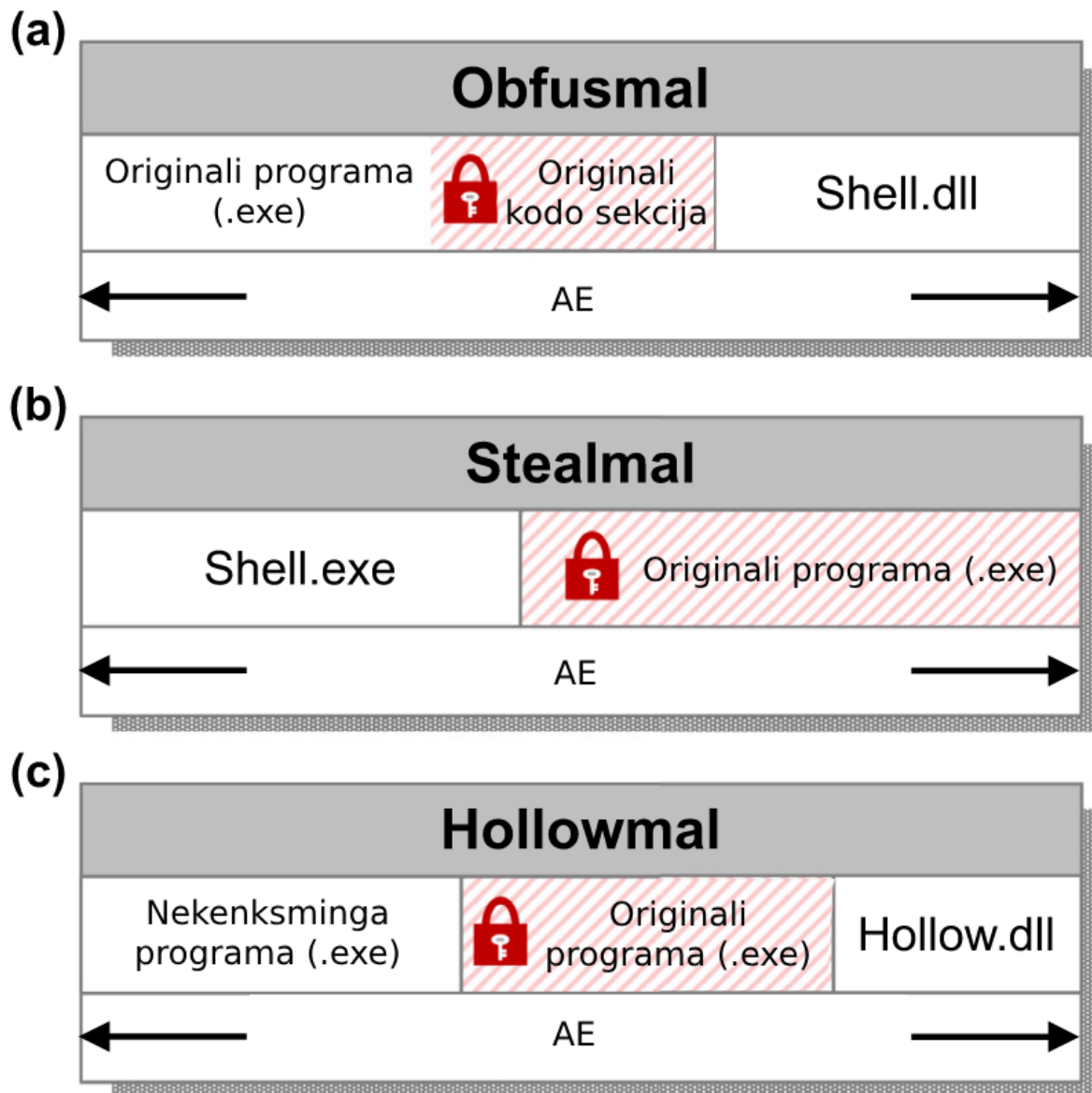
- [AKF⁺18] H. S. Anderson, A. Kharkar, B. Filar, D. Evans, P. Roth. *Learning to Evade Static PE Machine Learning Malware Models via Reinforcement Learning*. 2018-01-30. [žiūrėta 2024-09-30]. Prieiga per internetą: <https://doi.org/10.48550/arXiv.1801.08917>.
- [AR18] H. S. Anderson, P. Roth. *EMBER: An Open Dataset for Training Static PE Malware Machine Learning Models*. arXiv e-prints, 2018-04-01. [žiūrėta 2024-11-17]. Prieiga per internetą: <https://doi.org/10.48550/arXiv.1804.04637>.
- [CDH⁺16] X. Chen, Y. Duan, R. Houthooft, J. Schulman, I. Sutskever, P. Abbeel. *InfoGAN: Interpretable Representation Learning by Information Maximizing Generative Adversarial Nets*. 2016-06-11. [žiūrėta 2024-10-07]. Prieiga per internetą: <https://doi.org/10.48550/arXiv.1606.03657>.
- [CSD19] R. L. Castro, C. Schmitt, G. Dreo. AIMED: Evolving Malware with Genetic Programming to Evade Detection. Iš: *2019 18th IEEE International Conference On Trust, Security And Privacy In Computing And Communications/13th IEEE International Conference On Big Data Science And Engineering (TrustCom/BigDataSE)*. 2019, p. 240–247 [žiūrėta 2024-09-23]. ISSN 2324-9013. Prieiga per internetą: <https://doi.org/10.1109/TrustCom/BigDataSE.2019.00040>.
- [DBL⁺21] L. Demetrio, B. Biggio, G. Lagorio, F. Roli, A. Armando. Functionality-Preserving Black-Box Optimization of Adversarial Windows Malware. *IEEE Transactions on Information Forensics and Security*. 2021, tomas 16, p. 3469–3478 [žiūrėta 2024-10-14]. ISSN 1556-6021. Prieiga per internetą: <https://doi.org/10.1109/TIFS.2021.3082330>.
- [DCB⁺21] L. Demetrio, S. E. Coull, B. Biggio, G. Lagorio, A. Armando, F. Roli. Adversarial EXamples: A Survey and Experimental Evaluation of Practical Attacks on Machine Learning for Windows Malware Detection. *ACM Trans. Priv. Secur.* 2021, tomas 24, numeris 4, 27:1–27:31 [žiūrėta 2024-09-30]. ISSN 2471-2566. Prieiga per internetą: <https://doi.org/10.1145/3473039>.
- [FWL⁺19] Z. Fang, J. Wang, B. Li, S. Wu, Y. Zhou, H. Huang. Evading Anti-Malware Engines With Deep Reinforcement Learning. *IEEE Access*. 2019, tomas 7, p. 48867–48879 [žiūrėta 2024-09-18]. ISSN 2169-3536. Prieiga per internetą: <https://doi.org/10.1109/ACCESS.2019.2908033>.
- [HT17] W. Hu, Y. Tan. *Generating Adversarial Malware Examples for Black-Box Attacks Based on GAN*. 2017-02-20. [žiūrėta 2024-09-18]. Prieiga per internetą: <http://arxiv.org/abs/1702.05983>.

- [YPT22] J. Yuste, E. G. Pardo, J. Tapiador. Optimization of Code Caves in Malware Binaries to Evade Machine Learning Detectors. *Computers & Security*. 2022, tomas 116, p. 102643 [žiūrėta 2024-10-07]. ISSN 0167-4048. Prieiga per internetą: <https://doi.org/10.1016/j.cose.2022.102643>.
- [RSR⁺18] I. Rosenberg, A. Shabtai, L. Rokach, Y. Elovici. Generic Black-Box End-to-End Attack Against State of the Art API Call Based Malware Classifiers. Iš: M. Bailey, T. Holz, M. Stamatogiannakis, S. Ioannidis (sudarytojai). *Research in Attacks, Intrusions, and Defenses*. Cham: Springer International Publishing, 2018, p. 490–510. ISBN 978-3-030-00470-5. Prieiga per internetą: https://doi.org/10.1007/978-3-030-00470-5_23.
- [SB15] J. Saxe, K. Berlin. Deep Neural Network Based Malware Detection Using Two Dimensional Binary Program Features. Iš: *2015 10th International Conference on Malicious and Unwanted Software (MALWARE)*. 2015, p. 11–20 [žiūrėta 2024-11-04]. Prieiga per internetą: <https://doi.org/10.1109/MALWARE.2015.7413680>.
- [ZCY⁺24] F. Zhong, X. Cheng, D. Yu, B. Gong, S. Song, J. Yu. MalFox: Camouflaged Adversarial Malware Example Generation Based on Conv-GANs Against Black-Box Detectors. *IEEE Transactions on Computers*. 2024, tomas 73, numeris 4, p. 980–993 [žiūrėta 2024-09-15]. ISSN 1557-9956. Prieiga per internetą: <https://doi.org/10.1109/TC.2023.3236901>.
- [ZHZ⁺22] F. Zhong, P. Hu, G. Zhang, H. Li, X. Cheng. Reinforcement Learning Based Adversarial Malware Example Generation against Black-Box Detectors. *Computers & Security*. 2022, tomas 121, p. 102869 [žiūrėta 2024-09-14]. ISSN 0167-4048. Prieiga per internetą: <https://doi.org/10.1016/j.cose.2022.102869>.
- [ZZY⁺22] E. Zhu, J. Zhang, J. Yan, K. Chen, C. Gao. N-Gram MalGAN: Evading Machine Learning Detection via Feature n-Gram. *Digital Communications and Networks*. 2022, tomas 8, numeris 4, p. 485–491 [žiūrėta 2024-09-23]. ISSN 2352-8648. Prieiga per internetą: <https://doi.org/10.1016/j.dcan.2021.11.007>.

Priedai

Priedas nr. 1

Kompleksinių perturbacijų iliustracijos



4 pav. Obfusmal (a), Stealmal (b) ir Hollowmal (c) perturbacijų veikimo principų iliustracijos. Adaptuota iš [ZHZ⁺22]