

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 .....	7
1.3. Perturbacijos .....	7
1.3.1. Baitų lygio perturbacijos .....	8
1.3.2. Semantinės perturbacijos .....	9
1.3.3. Kompleksinės perturbacijos .....	9
1.4. GAN tipo modelių karkasai .....	9
1.4.1. <i>MalGAN</i> .....	10
1.4.2. <i>N-gram MalGAN</i> .....	11
1.4.3. <i>MalFox</i> .....	11
1.5. Skatinamojo mokymosi tipo modelių karkasai .....	12
1.5.1. <i>DQEAF</i> .....	12
1.5.2. <i>MalInfo</i> .....	13
1.6. Genetinių algoritmų tipo modelių karkasai .....	14
1.6.1. <i>AIMED</i> .....	14
1.6.2. <i>GAMMA</i> .....	15
1.7. Nevalidaus PE formato problema .....	15
2. KARKASŲ VERTINIMAS .....	17
3. EKSPERIMENTINIS TYRIMAS .....	18
3.1. Tyrimo aprašymas .....	18
3.1.1. GAN įgyvendinimas .....	18
3.1.2. GAN mokymosi duomenys .....	18
3.1.3. Karkaso parametrai .....	19
3.1.4. Realų kenkėjiškų programų rinkinys .....	19
3.1.5. Komerciniai detektoriai .....	19
3.2. Mokymosi etapas .....	20
3.3. Varžymosi principais pagrįstos atakos prieš komercinius detektorius .....	21
REZULTATAI IR IŠVADOS .....	22
ŠALTINIAI .....	23
PRIEDAI .....	25
1 priedas. Tarpiniai tyrimo rezultatai .....	25

## 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

**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

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

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

**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*)

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

**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

**Surogatinis Modelis (angl. *Surrogate Model*).** ML modelis, aproksimuojantis kitą ML modelį, kurio parametrai (svoriai) nėra žinomi

**Varžymosi principais pagrįstos atakos (angl. *Adversarial Attacks*).** Tai atakos, pritaikytos „apgauti“ ML klasifikatorius

## Santrumpos

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

## Įvadas

Pastaraisiais metais kenkėjiškas kodas ir programos kuriamos itin sparčiai ( $\sim 450000$  kenkėjiškų programų per dieną 2024 m. AV-TEST<sup>1</sup> 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<sup>+</sup>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<sup>+</sup>18; ZHZ<sup>+</sup>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<sup>+</sup>21]. Žinoma, rasti šią ribą nėra trivialus uždavinys. Mokslinėje literatūroje išskiriami 3 ribos paieškos atvejai [FWL<sup>+</sup>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<sup>+</sup>18]. Todėl šiame darbe nagrinėjami modeliai, gebantys generuoti varžymosi principais pagrįstų atakų obfusuotus kenkėjiško kodo pavyzdžius (AE) „juodos dėžės“ atvejais.

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

### Uždaviniai:

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

---

<sup>1</sup><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<sup>+</sup>18]. Dažniausia atakų platforma – operacinė sistema *Windows* ir jos naudojami PE formato failai (*VirusTotal*<sup>2</sup> 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<sup>+</sup>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<sup>+</sup>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. Šiame poskyryje išskiriami ir klasifikuojami mokslinėje literatūroje minimi požymiai.

---

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

### 1.2.1. PE formato programų požymiai

Išskiriami šie pagrindiniai PE formato programų požymiai:

- **DLL vardai (arba API vardai [HT17])** [ZCY<sup>+</sup>24]. PE faile turi būti nurodyti visi naudojami DLL ir jų API. Prieš pradedant mokytį 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<sup>+</sup>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

Baitų lygio požymiai gali būti ištraukiami iš bet kokio formato failų. Mokslinėje literatūroje minimi šie pagrindiniai baitų lygio požymiai:

- **Prasmingų žodžių (angl. Strings) kiekis** [AKF<sup>+</sup>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<sup>+</sup>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. Šiame poskyryje aptariamos mokslinėje literatūroje minimos perturbacijos.

### 1.3.1. Baitų lygio perturbacijos

Pačias paprasčiausias baitų lygio perturbacijas galima taikyti bet kokio formato failams, tačiau labiau prasmingos perturbacijos taikomos PE formato failams. Išskiriamos šios pagrindinės baitų lygio perturbacijos:

- **ARBE (Append Random Bytes at the End)** [FWL<sup>+</sup>19]. PE formato failo gale pridedami atsitiktiniai baitai.
- **ARI (Append Random Import)** [FWL<sup>+</sup>19]. PE formato failo *ImportAddressTable* lentelėje pridedama atsitiktinai pavadinta biblioteka su atsitiktinai pavadinta funkcija.
- **ARS (Append Randomly named Section)** [FWL<sup>+</sup>19]. PE formato failo *SectionTable* lentelėje pridedamos atsitiktinės sekcijos (sekcijos ir jų tipai yra apibrėžti PE formate).
- **RS (Remove Signature)** [FWL<sup>+</sup>19]. Sertifikato pašalinimas iš PE formato failo *CertificateTable* lentelės.
- **Naujas įeities taškas** [AKF<sup>+</sup>18]. Prasidėjus programai, iškart peršokama nuo naujo įeities taško į originalųjį.
- **Header Fields** [DCB<sup>+</sup>21]. PE formato failo *PE Header* ir *Optional Header* dalių specifinių laukų keitimas (pvz., sekcijos pavadinimo keitimas [AKF<sup>+</sup>18]).
- **Partial DOS** [DCB<sup>+</sup>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<sup>+</sup>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<sup>+</sup>21]. Nulinių baitų pridėjimas failo gale.
- **Full DOS** [DCB<sup>+</sup>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<sup>+</sup>21]. Pakeičiama PE formato faile *DOS* dalyje esanti *PE Offset* reikšmė į didesnę<sup>3</sup>. Taip padidinama (išplečiama) visa *DOS* dalis. Tolesnis perturbacijos principas yra toks pat, kaip ir *Full DOS*.
- **Shift** [DCB<sup>+</sup>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, padidinus<sup>3</sup>  $S$  per  $n$ , atsiranda  $n$  baitų laisvos vietos iki sekcijos pradžios, kurią galima keisti be jokios įtakos programos veikimui.

---

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



### 1.3.2. Semantinės perturbacijos

Semantinių perturbacijų įgyvendinimas taip pat atliekamas baitų lygyje, tačiau šie pokyčiai turi aukštesnio lygio prasmę. Išskiriamos šios semantinės perturbacijos:

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

### 1.3.3. Kompleksinės perturbacijos

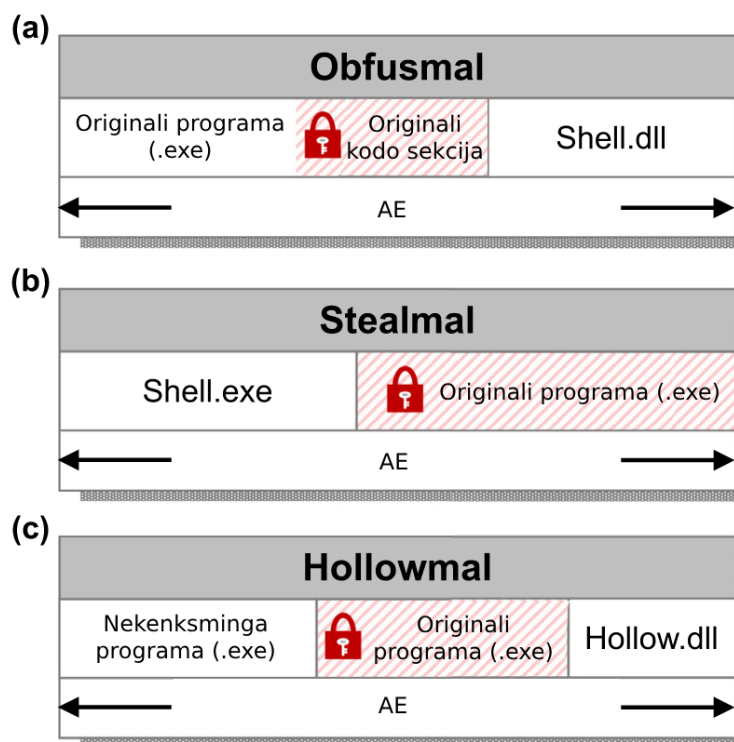
Kompleksinės perturbacijos yra pritaikomos tam tikriems tikslams. Obfuskacijos ir varžymosi principais pagrįstų atakų tikslams literatūroje minimos šios kompleksinės perturbacijos:

- **Obfusmal** [ZCY<sup>+</sup>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 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 1-ame pav.
- **Stealmal** [ZCY<sup>+</sup>24]. Visa originali programa užšifruojama ir pridedama prie programos *Shell.exe* galo. *Shell.exe* geba atšifruoti originalią programą ir perduoti jai kontrolę. Iliustracija pateikiama 1-ame pav.
- **Hollowmal** [ZCY<sup>+</sup>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 1-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<sup>+</sup>16]. Kenkėjiško kodo obfuskacijos kontekste ir ypač „juodos dėžės“ atvejais, diskriminatorius atlieka surogatinio modelio vaidmenį. Bendras GAN modelių mokymosi etapas yra tokia seka [HT17; ZCY<sup>+</sup>24; ZZY<sup>+</sup>22]:

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).



1 pav. Obfusmal (a), Stealmal (b) ir Hollowmal (c) perturbacijų veikimo principų iliustracijos.  
Adaptuota iš [ZHZ<sup>+</sup>22]

4. Visa seka kartojama nustatytą kiekį kartų.

#### 1.4.1. *MalGAN*

Tai vienas iš pirmųjų ir populiariausių GAN tipo modelių karkasų. „Juodos dėžės“ detektoriai čia apibrėžiami kaip populiarūs ML klasifikatoriai, tokie, kaip MLP (*angl. Multilayer Perceptron*), RF (*angl. Random Forest*), DT (*angl. Decision Tree*), SVM (*angl. Support Vector Machine*). *MalGAN* karkaso [HT17] tikslas ir apibrėžimas pateikiami 1-oje lentelėje.

1 lentelė. *MalGAN* karkasas

<b>Tikslas</b>	Efektyviai išvengti AE aptikimo, kai ML kenkėjiškų programų detektoriaus įgyvendinimas nežinomas („juodos dėžės“ atvejis).
<b>Surogatinis modelis</b>	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.
<b>ML modelis</b>	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.

<b>Požymiai</b>	<i>MalGAN</i> straipsnyje [HT17] 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 požymiai <sup>4</sup> .
<b>Perturbacijos</b>	Semantinės perturbacijos (1.3.2) – nereikalingų API vardų požymių pridėjimas.

#### 1.4.2. *N-gram MalGAN*

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

2 lentelė. *N-gram MalGAN* karkasas

<b>Tikslas</b>	Supaprastinti, pagreitinti ir pagerinti varžymosi principais pagrįstas atakas. Pašalinti prielaidas <sup>4</sup> apie detektorių „juodos dėžės“ atvejais.
<b>Surogatinis modelis</b>	Surogatinio modelio veikimas ir architektūra tokia pati, kaip ir <i>MalGAN</i> (1.4.1).
<b>ML modelis</b>	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.
<b>Požymiai</b>	Baitų lygio požymiai (1.2.2) – $n$ -gramos.
<b>Perturbacijos</b>	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 baitų lygio perturbaciją <i>ARBE</i> , tik šiuo atveju pridėdami baitai nebūtų atsitiktiniai, o norimos $n$ -gramos.

#### 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<sup>+</sup>24] tikslas ir apibrėžimas pateikiami 3-oje lentelėje.

3 lentelė. *MalFox* karkasas

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

<b>Tikslas</b>	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).
<b>Surogatinis modelis</b>	Surogatinis modelis, kaip ir kituose GAN tipo modelių karkasuose, naudojamas kaip diskriminatorius. Įvestis – perturbuota programa. Išvestis – klasifikacija į kenksmingą arba nekenksmingą. Įgyvendinimas – konvoliucinis neuroninis tinklas (CNN).
<b>ML modelis</b>	Standartinis GAN generatorius, požymių vektorių sujungiantis su „triukšmo“ vektoriumi. Įgyvendinimas – konvoliucinis neuroninis tinklas (CNN).
<b>Požymiai</b>	PE formato programų požymiai (1.2.1) – DLL vardai.
<b>Perturbacijos</b>	Visos kompleksinės perturbacijos (1.3.3).

## 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<sup>+</sup>19]. Bendras RL modelių mokymosi etapas yra tokia seka [FWL<sup>+</sup>19; ZHZ<sup>+</sup>22]:

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 įgyvendinimo).
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.

### 1.5.1. DQEAf

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

4 lentelė. DQEAf karkasas

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

### 1.5.2. *MalInfo*

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

5 lentelė. *MalInfo* karkasas

<b>Tikslas</b>	Surasti optimalią obfuskacijos strategiją konkrečiai programai, pagal kurią sukurtas AE nebūtų aptiktas komercinių kenkėjiškų programų detektorius.
<b>Surogatinis modelis</b>	RL surogatinis modelis nenaudojamas. „Juodos dėžės“ detektoriumi pasirinkti komerciniai detektoriai ( <i>VirusTotal</i> ).
<b>ML modelis</b>	Agentas implementuotas kaip klasikiniai ML algoritmai (konkrečiai dinaminis programavimas ir skirtumų laike ( <i>angl. temporal difference</i> ) algoritmas).

<b>Požymiai</b>	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 $S = A^3$ .
<b>Perturbacijos</b>	$A = \{0,1,2,3\}$ <ul style="list-style-type: none"> <li>• <i>Null</i> perturbacija – naudinga tik formaliam pilnumui (atitinka nulinę <math>A</math> veiksmą).</li> <li>• Visos kompleksinės perturbacijos (1.3.3), t. y. tokios pačios, kaip ir <i>MalFox</i> (1.4.3) karkaso.</li> </ul>

## 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 [YPT22]:

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.

### 1.6.1. AIMED

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

6 lentelė. *AIMED* karkasas

<b>Tikslas</b>	AE generavimo greičio padidinimas ir modelių kompleksiskumo sumažinimas, lyginant su GAN ir RL tipo modelių karkasais.
<b>Surogatinis modelis</b>	Surogatinis modelis nenaudojamas. Naudojami „juodos dėžės“ detektoriai yra <b>3 komerciniai</b> ( <i>Kaspersky</i> , <i>ESET</i> , <i>Sophos</i> ) ir vienas ML modelis – GBDT.

<b>ML modelis</b>	Klasikinis GA modelis – veikimas visiškai atitinką bendrą seką. Tinkamumo ( <i>angl. fitness</i> ) vertinamas remiasi AE požymių vektoriaus panašumu į originalios programos požymių vektorių (kuo mažiau panašūs, tuo tinkamumo įvertinimas didesnis).
<b>Požymiai</b>	Baitų lygio požymiai (1.2.2) – atskiras $n$ -gramų atvejis, kai $n = 1$ .
<b>Perturbacijos</b>	Baitų lygio perturbacijos <sup>5</sup> (1.3.1).

### 1.6.2. GAMMA

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

7 lentelė. GAMMA karkasas

<b>Tikslas</b>	Efektyvus (neaptikimo šansų didinimas naudojant perturbacijas, paremtas nekenksmingomis programomis) varžymosi principais pagrįstų atakų kūrimas.
<b>Surogatinis modelis</b>	Surogatinis modelis nenaudojamas. GBDT ir <i>MalConv</i> pasirinkti kaip „juodos dėžės“ detektoriai.
<b>ML modelis</b>	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 detektorių ir pridėtų požymių dydį baitais (norima pridėti kuo mažiau požymių).
<b>Požymiai</b>	<ul style="list-style-type: none"> <li>• PE formato programų požymiai (1.2.1).</li> <li>• Kodas sekcijose (nestandartinis požymis).</li> </ul>
<b>Perturbacijos</b>	<ul style="list-style-type: none"> <li>• Visos baitų lygio perturbacijos (1.3.1), gebančios pridėti baitus.</li> <li>• Autoriai pažymi, jog gali būti naudojama ir DLL / API vardų pridėjimo semantinė perturbacija (1.3.2).</li> </ul>

## 1.7. Nevalidaus PE formato problema

Anderson et al. [AKF<sup>+</sup>18], atlikdami eksperimentus su funkcionalumą išlaikančiomis perturbacijomis PE formato failams, pastebėjo, jog ne visais atvejais perturbuotos programos veikia

<sup>5</sup> autoriai rėmėsi perturbacijomis, aprašytomis [AKF<sup>+</sup>18]

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.

Norint visiškai išvengti nevalidaus PE formato problemos tenka taikyti perturbacijas, nekeičiančias originalių programų, o taikančias kitokius obfuskacijos metodus. Iš 1.3 poskyryje 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**).

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

Karkasas	K-1	K-2	K-3	K-4
<i>MalFox</i> (1.4.3)	✓	✓	✓	56,00%
<i>MallInfo</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%

Kriterijų vertinimo (8-oje) 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 įgyvendinimas

Naudosime viešai *Github* platformoje paskelbtą<sup>6</sup> *MalGAN* karkaso ML modelio įgyvendinimą.

$$L_D = -\mathbb{E}_{x \in BB_{Benign}} \log(1 - D_{\theta_d}(x)) - \mathbb{E}_{x \in BB_{Malware}} \log D_{\theta_d}(x). \quad (1)$$

$$L_G = \mathbb{E}_{m \in S_{Malware}, z \sim p_{\text{uniform}}[0,1]} \log D_{\theta_d}(G_{\theta_g}(m, z)) \quad (2)$$

1-oje ir 2-oje formulėse atitinkamai pateikiamos *MalGAN* diskriminatoriaus ir generatoriaus nuostolių funkcijos. Šiose formulėse naudojamų simbolių prasmės paaiškinimas pateikiamas 9-oje lentelėje.

9 lentelė. *MalGAN* nuostolių funkcijų formulėse naudojamų simbolių paaiškinimas

Simbolis	Prasmė
$BB_{Benign}$	Požymių vektorių, kuriuos „juodos dėžės“ detektorius klasifikuoja kaip <b>nekenksmingus</b> , aibė
$BB_{Malware}$	Požymių vektorių, kuriuos „juodos dėžės“ detektorius klasifikuoja kaip <b>kenksmingus</b> , aibė
$S_{Malware}$	Tikrų kenkėjiškų programų požymių vektorių aibė
$D_{\theta_d}$	Diskriminatorius su $\theta_d$ svoriais ( $D_{\theta_d} : \{0,1\}^{n_m} \rightarrow [0,1]$ , $n_m \in \mathbb{N}$ )
$G_{\theta_g}$	Generatorius su $\theta_g$ svoriais ( $G_{\theta_g} : \{0,1\}^{n_m} \times [0,1]^{n_z} \rightarrow \{0,1\}^{n_m}$ , $n_m, n_z \in \mathbb{N}$ )
$z$	Triukšmo vektorius

##### 3.1.2. GAN mokymosi duomenys

Mokymuisi naudosime *EMBER* [AR18] duomenų rinkinį. Mokymosi 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 mokymosi / validacijos skaidinį santykiu 4 : 1.

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

### 3.1.3. Karkaso parametrai

Tyrimo metu atlikti 2 eksperimentai (**E-1** ir **E-2**), kuriems parinkti 10-oje lentelėje aprašyti parametrai. **E-1** atitinka *MalGAN* straipsnyje [HT17] aprašytą eksperimentą, **E-2** – eksperimentas su autoriaus parinktais parametrais siekiant padidinti atakų efektyvumą.

10 lentelė. Eksperimentuose naudoti ML modelio parametrai

Parametras	E-1	E-2
Dvejetaus požymių vektoriaus $M$ dimensija	160	10000
„Triukšmo“ vektoriaus $Z$ dimensija	10	100
Neuronų skaičius generatoriaus „paslėptuose“ sluoksniuose	256	100, 256, 512, 1024
Neuronų skaičius diskriminatoriaus „paslėptuose“ sluoksniuose	256	256, 256
Diskriminatoriaus mokymosi greitis	$10^{-3}$	$4 \cdot 10^{-4}$
Generatoriaus mokymosi greitis	$10^{-3}$	$10^{-6}$

### 3.1.4. Realių kenkėjiškų programų rinkinys

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

### 3.1.5. Komerciniai detektoriai

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

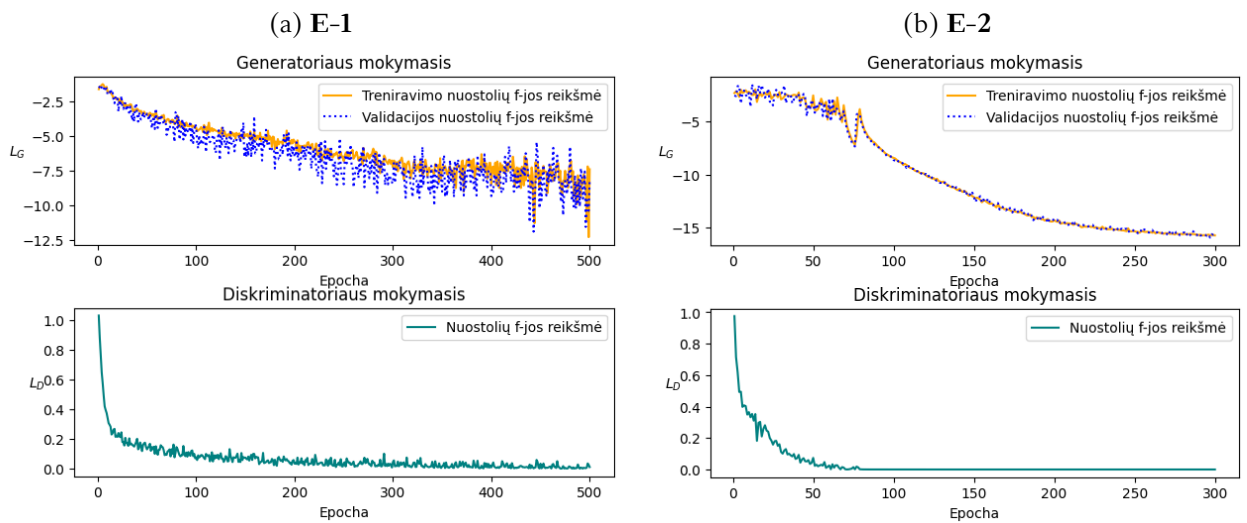
<sup>7</sup><https://virusshare.com/>

<sup>8</sup><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 2-ame pav.



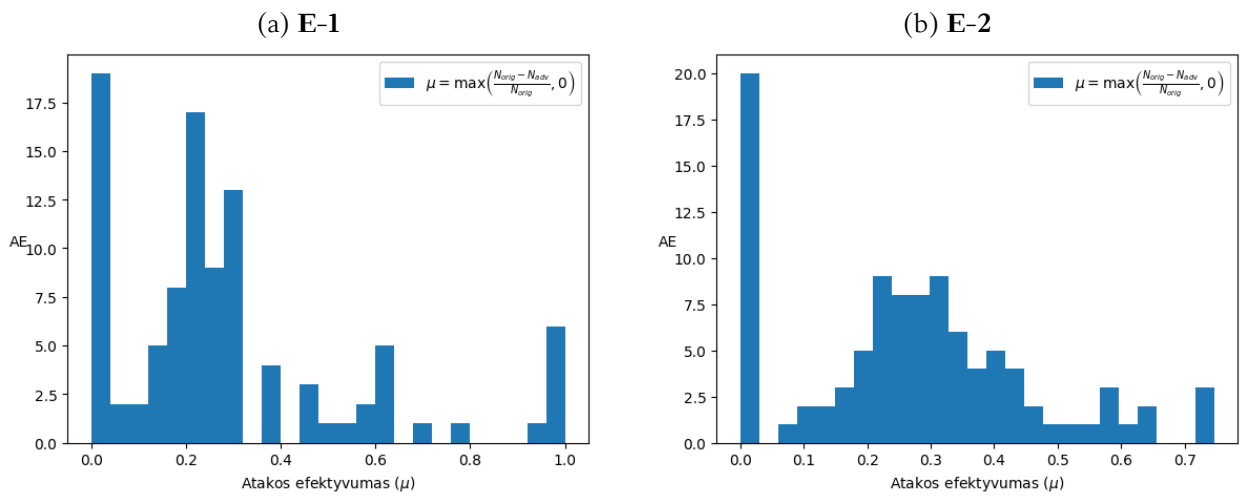
2 pav. *MalGAN* generatoriaus ir diskriminatoriaus nuostolių funkcijų reikšmės mokymosi etape

### 3.3. Varžymosi principais pagrįstos atakos prieš komercinius detektorius

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

Atakų efektyvumą  $\mu$  skaičiuosime remiantis Zhong et al. pasiūlyta formule  $\mu' = \frac{N_{orig} - N_{adv}}{N_{orig}}$  [ZCY<sup>+</sup>24]. Eksperimentuose galimi atvejai, kai originalią programą aptinka mažesnis detektorių kiekis, nei obfusuotą ( $N_{orig} < N_{adv}$ ). Tokiais atvejais  $\mu' < 0$ . Kadangi neigiamas atakos efektyvumas neturi prasmės, laikome, jog atvejais, kai  $\mu' < 0$  ataka nėra efektyvi (t. y. jos efektyvumas lygus 0)  $\Rightarrow \mu = \max(\mu', 0)$ .

Atakų efektyvumų pasiskirstymas pavaizduotas 3-ame pav., tarpiniai rezultatai ( $N_{orig}$  ir  $N_{adv}$ ), naudojami pasiskirstymų skaičiavimui pateikiami 4-ame pav. (1-ame priede).



3 pav. Varžymosi principais pagrįstų atakų efektyvumo prieš komercinius detektorius pasiskirstymas

Tyrimo apibendrinimui 11-oje lentelėje pateikiamos abiejų eksperimentų statistikos, apskaičiuotos iš 3-ame pav. vaizduojamų duomenų įtraukiant ir neefektyvias atakas (kai  $\mu = 0$ ).

11 lentelė. Eksperimentų statistikos

	E-1	E-2
<b>Efektyvumo vidurkis (<math>\bar{\mu}</math>)</b>	28,89 %	26,85 %
<b>Standartinis nuokrypis (<math>\sigma</math>)</b>	26,31 %	18,97 %

# Rezultatai ir išvados

## Rezultatai

1. Išanalizuoti 8 trims kategorijoms (GAN, RL, GA) priklausančios varžymosi principais pagrįstų atakų karkasai.
2. Jiems apibrėžti vertinimo kriterijai, išskiriantys geriausiai realioms sąlygoms pritaikytus karkasus.
3. Atliktas tyrimas, kurio metu K-2 kriterijaus neįgyvendinantis varžymosi principais pagrįstų atakų karkasas *MalGAN* (1.4.1) naudojamas AE generavimui prieš komercinius detektorius ir renkami K-4 kriterijų atitinkantys duomenys (t. y. nustatomas karkaso K-4 kriterijus tokiomis sąlygomis, kai tenkinamas K-2).

## Išvados

1. Tyrimo metu atlikti 2 eksperimentai parodo, jog net ir padidinus *MalGAN* ML modelio galimybes (suteikus daugiau kompiuterijos resursų), atakų efektyvumas prieš komercinius detektorius nepadidėja ( $\bar{\mu}_{E-2} < \bar{\mu}_{E-1}$ ), tačiau yra stabilesnis ( $\sigma_{E-2} < \sigma_{E-1}$ ).
2. Abiem atvejais tyrimo metu gautas K-4 įvertis ( $\bar{\mu}$ ) yra žymiai mažesnis už kriterijų vertinimo metu naudojamą 81,00 % *MalGAN* karkaso įvertinimą (kai K-2 netenkinamas). Tai paaiškina 8-oje lentelėje matomas K-4 kriterijaus anomalijas ir parodo, jog kriterijai ir karkasų vertinimas yra adekvatūs.
3. Laikant, jog parinkti kriterijai ir jų vertinimas yra patikimi, geriausiai iš analizuotų karkasų „juodos dėžės“ atvejams pritaikytas varžymosi principais pagrįstų atakų karkasas yra *MalFox* (1.4.3), tenkinantis visus 3 kokybinius kriterijus ir siekiantis 56,00 % kiekybinį įvertinimą (atakų efektyvumą).

## Šaltiniai

- [AKF<sup>+</sup>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<sup>+</sup>16] X. Chen, Y. Duan, R. Houthoof, 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<sup>+</sup>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<sup>+</sup>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<sup>+</sup>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>.

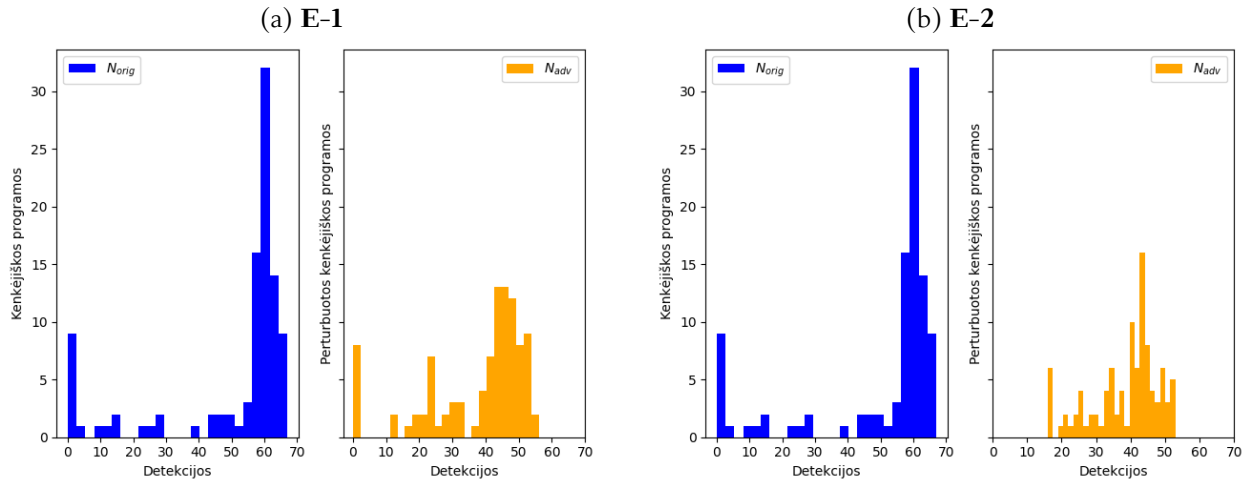
- [YPT22] J. Yuste, E. G. Pardo, J. Tapiador. Optimization of Code Caves in Malware Binaries to Evade Machine Learning Detectors. *Computers & Security*. 2022, tomas 116, p. 102643 [žiūrėta 2024-10-07]. ISSN 0167-4048. Prieiga per internetą: <https://doi.org/10.1016/j.cose.2022.102643>.
- [RSR<sup>+</sup>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](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<sup>+</sup>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<sup>+</sup>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<sup>+</sup>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

### Tarpiniai tyrimo rezultatai



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