

# The Robust Malware Detection Challenge and Greedy Random Accelerated Multi-Bit Search

AdvML & GRAMS



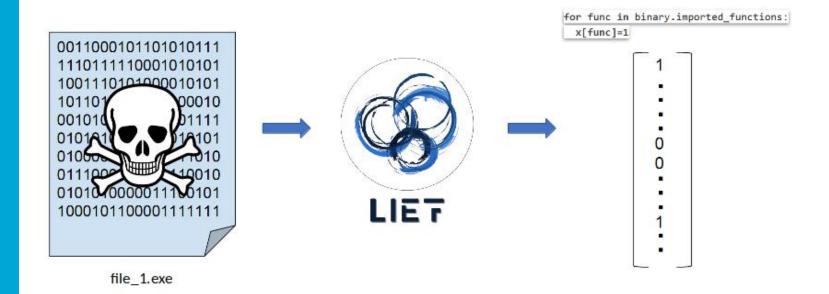
<u>Sicco Verwer</u>, Azqa Nadeem, Christian Hammerschmidt, Laurens Bliek, Abdullah Al-Dujaili, Una-May O'Reilly

# The challenge



ALFA:

Thanks to Abdullah Al-Dujaili and Una-May O'Reilly



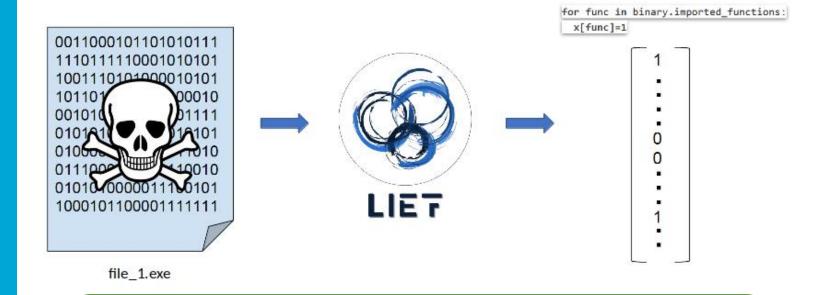


# The challenge



ALFA

Thanks to Abdullah Al-Dujaili and Una-May O'Reilly





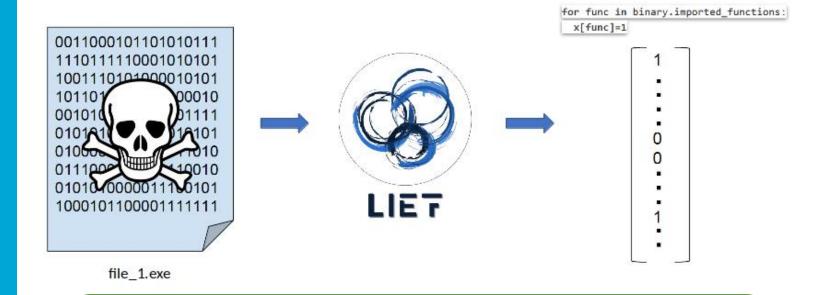
22.761 bits per file, 34.200 files each bit indicates the presence of one sys call

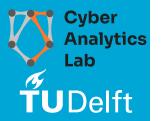
# The challenge



ALFA

Thanks to Abdullah Al-Dujaili and Una-May O'Reilly



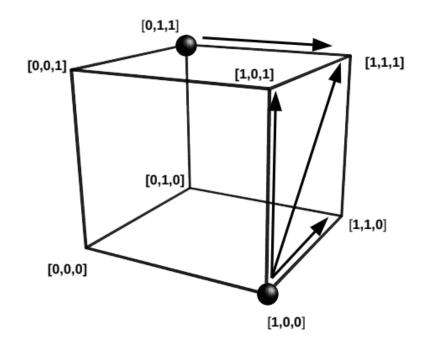


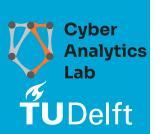
Classifiers obtain over 90% accuracy Malware authors will evade, but how?

#### The threat model



- Problems in existing threat models:
  - An ε radius does not capture adding sys calls
  - Removing sys calls can nullify the malware
- Solution, or a step in the right direction:





#### The threat model



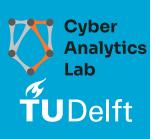
- Any number of 0 valued bits can be flipped to 1 values
- Use adversarial training to attack and defend

$$\theta^* \in \arg\min_{\theta \in \mathbb{R}^p} \mathbb{E}_{(\mathbf{x},y) \sim \mathcal{D}} \left[ \underbrace{\max_{\bar{\mathbf{x}} \in \mathcal{S}(\mathbf{x})} L(\theta, \bar{\mathbf{x}}, y)}_{\text{adversarial learning}} \right]$$

Implemented in the SLEIPNIR framework

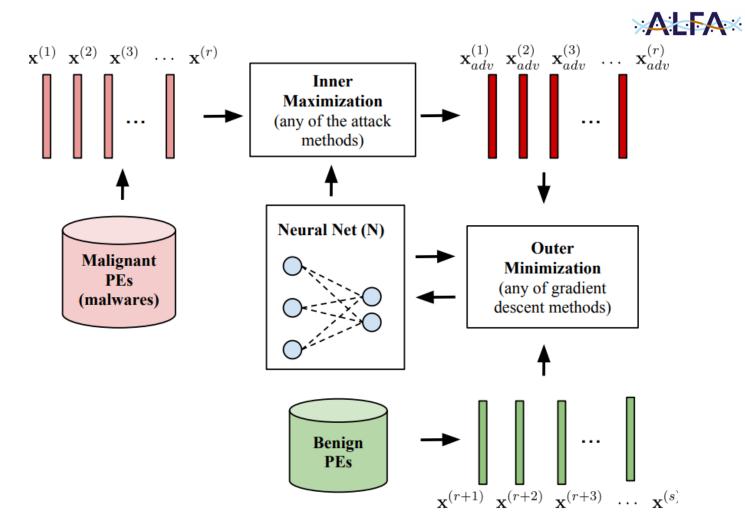
Al-Dujaili, A. Huang, E. Hemberg and U. O'Reilly, *Adversarial Deep Learning for Robust Detection of Binary Encoded Malware 2018 IEEE Security and Privacy Workshops (SPW)* 

https://github.com/ALFA-group/robust-adv-malware-detection



#### **SLEIPNIR**

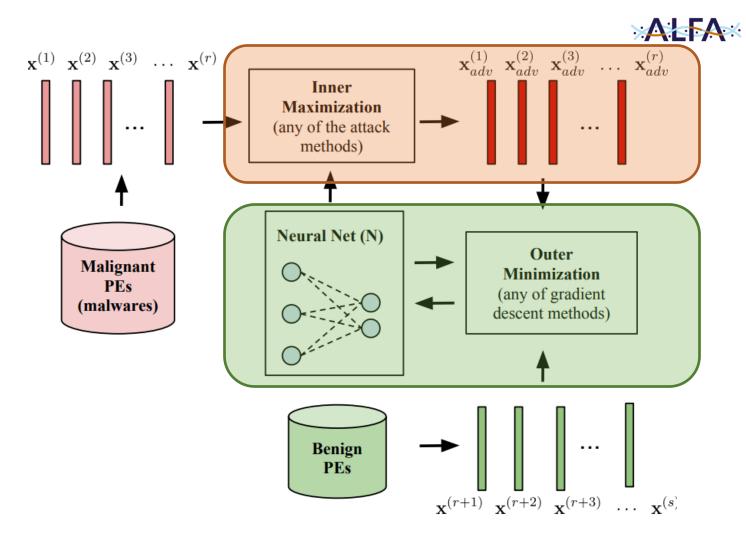






#### Attack – inner maximizer







Defense – outer maximizer

#### **Evasion rates**





Model		Adversary							
	Natural	${\tt dFGSM}^k$	$\mathtt{rFGSM}^k$	$\mathbf{BGA}^k$	$\mathtt{BCA}^k$				
Natural	8.1	99.7	99.7	99.7	41.7				
$\mathbf{dFGSM}^k$	6.4	6.4	21.1	7.3	27.4				
$\mathtt{rFGSM}^k$	5.7	7.0	5.9	5.9	6.8				
$\mathbf{BGA}^k$	7.6	39.6	17.8	7.6	10.9				
$\mathtt{BCA}^k$	7.6	99.5	99.5	91.8	7.9				



dFGSM<sup>k</sup> = deterministic Fast Gradient Sign Method BGA<sup>k</sup> = Bitwise Gradient Ascend





Model			Adve	rsary	
	Natural	$\mathbf{dFGSM}^k$	$\mathtt{rFGSM}^k$	$\mathtt{BGA}^k$	?
Natural	8.1	99.7	99.7	99.7	?
$\mathbf{dFGSM}^k$	6.4	6.4	21.1	7.3	
${\tt rFGSM}^k$	5.7	7.0	5.9	5.9	?
$\mathbf{BGA}^k$	7.6	39.6	17.8	7.6	
$\mathtt{BCA}^k$	7.6	99.5	99.5	91.8	7.9

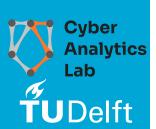
Is there a better attacker?







Model			Adve	rsary	
	Natural	$\mathtt{dFGSM}^k$	$\mathtt{rFGSM}^k$	$\mathbf{BGA}^k$	?
Natural	8.1	99.7	99.7	99.7	?
$\mathbf{dFGSM}^k$	6.4	6.4	21.1	7.3	
$\mathtt{rFGSM}^k$	5.7	7.0	5.9	5.9	?
$\mathtt{BGA}^k$	7.6	39.6	17.8	7.6	
?		?			7.9



Is there a better attacker? Is there a better defender?





Model	_		Adve	rsary	
	Natural	$\mathtt{dFGSM}^k$	$\mathtt{rFGSM}^k$	$\mathtt{BGA}^k$	?
Natural	8.1	99.7	99.7	99.7	?
$\mathbf{dFGSM}^k$	6.4	6.4	21.1	7.3	
$\mathtt{rFGSM}^k$	5.7	7.0	5.9	5.9	?
$\mathbf{BGA}^k$	7.6	39.6	17.8	7.6	
?		?			?



Is there a better attacker?
Is there a better defender?
What if they do not know each other?



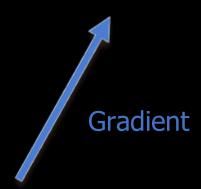


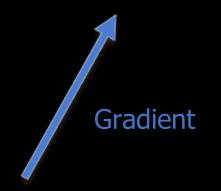
Model			Adve	rsary	
	Natural	$\mathbf{dFGSM}^k$	$\mathtt{rFGSM}^k$	$\mathtt{BGA}^k$	?
Natural	8.1	99.7	99.7	99.7	?
$\mathbf{dFGSM}^k$	6.4	6.4	21.1	7.3	
$\mathtt{rFGSM}^k$	5.7	7.0	5.9	5.9	?
$\mathbf{BGA}^k$	7.6	39.6	17.8	7.6	
?		?			?

Is there a better attacker?

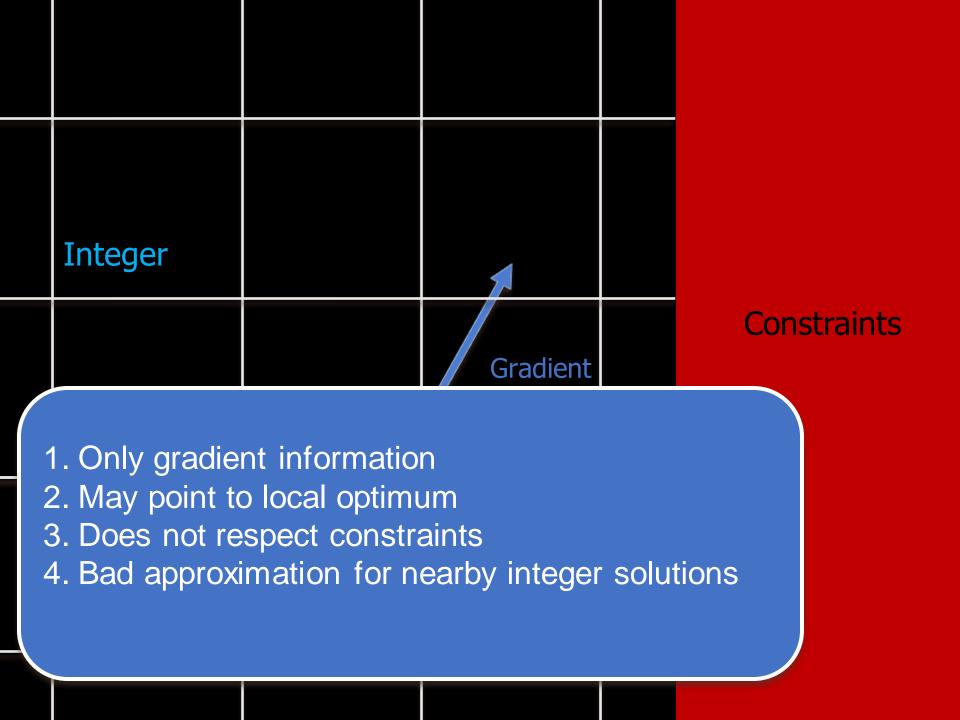


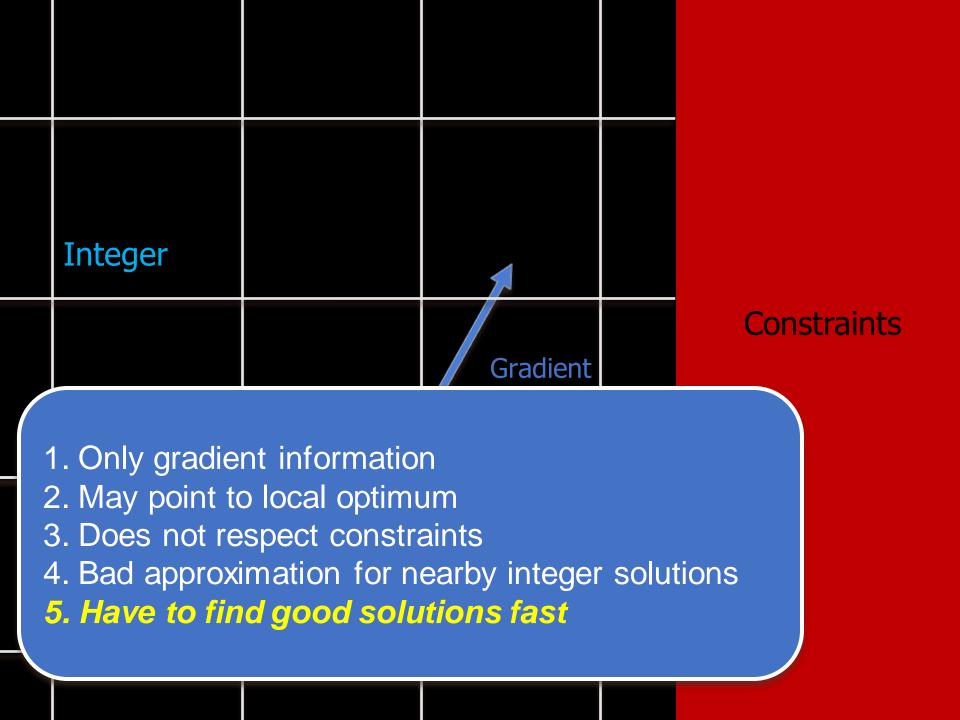
Our goal: build a new inner maximizer





#### Constraints





# GRAMS – an algorithm that



- Greedy follows the gradient
- Random is stochastic
- Adaptive slow when difficult, fast when easy
- Multibit flips multiple bits
- Search …, well, it does hill-climbing





Given a batch of data points

0	0	1	1	0	1	0	0
0	1	0	0	1	0	1	1
0	0	1	0	1	1	1	1
0	1	0	1	1	0	0	0

• We perform bitflips, which ones to flip?







Given a batch of data points

0	0	1	1	0	1	0	0
0	1	0	0	1	0	1	1
0	0	1	0	1	1	1	1
0	1	0	1	1	0	0	0

- We perform bitflips, which ones to flip?
  - Let's compute the gradient:

0.5	-0.5	0.0	-1.3	1.5	-0.6	0.3	0.4
0.2	-0.9	0.1	0.1	1.2	0.4	-0.4	-0.2
-1.2	0.5	-0.3	0.7	-1.4	-0.2	0.4	0.2
0.8	-0.2	0.8	-0.3	0.2	0.5	-0.1	0.0







Given a batch of data points

0	0	1	1	0	1	0	0
0	1	0	0	1	0	1	1
0	0	1	0	1	1	1	1
0	1	0	1	1	0	0	0

- We perform bitflips, which ones to flip?
  - Let's compute the gradient:

0.5		0.0	-1.0	1.0	-0.6	0.3	0.4
0.2	-0.9	0.1	0.1		0.4	-0.4	-0.2
	0.5	-0.3	0.7	-1.0	-0.2		
0.8	-0.2	0.8	-0.3		0.5		0.0







Given a batch of data points

0	0	1	1	0	1	0	0
0	1	0	0	1	0	1	1
0	0	1	0	1	1	1	1
0	1	0	1	1	0	0	0

- We perform bitflips, which ones to flip?
  - Let's compute the gradient:

0.5		0.0		1.0		0.3	0.4
0.2		0.1	0.1		0.4		
	0.5		0.7				
0.8		0.8			0.5		0.0



only change 0 to 1 (wrt original point)





0	0	1	1	0	1	0	0
0	1	0	0	1	0	1	1
0	0	1	0	1	1	1	1
0	1	0	1	1	0	0	0

- We perform bitflips, which ones to flip?
  - Let's compute the gradient:

0.5				1.0		
0.2					0.4	
	0.5		0.7			
0.8		8.0				



keep the top *k* in every row





Given a batch of data points

1	0	1	1	1	1	0	0
1	1	0	0	1	1	1	1
0	1	1	1	1	1	1	1
1	1	1	1	1	0	0	0

- We perform bitflips, which ones to flip?
  - Let's compute the gradient:

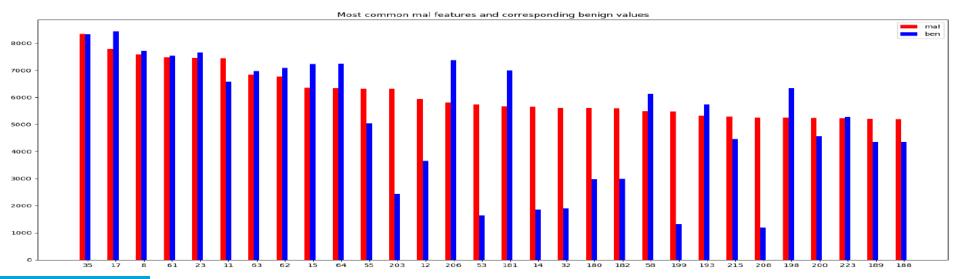
0.5				1.0		
0.2					0.4	
	0.5		0.7			
0.8		8.0				



flip these rows

#### Random







- But instead of fully random, fit a multivariate Bernoulli distribution D to the Benign datapoints
  - Sample bits from D, i.e., sample from Benign space
  - Restart 10 times
  - Follow gradients until convergence

# Adaptive



- What value to use for k?
- Set *k* to 8
- If there is no improvement:
  - Divide *k* by 2
  - Undo the bit update
- If there in an improvement:
  - Multiply k by 2
- If *k* < 0.5
  - Stop



#### Results – attack – evasion rates



	Nat.	rFGSM	BGA	Grosse	GRAMS	AGA	GwT	ENA
Natural	6.9	99.9	99.9	99.8	93.9	94.2	100.0	84.7
rFGSM	5.8	5.8	5.8	6.3	44.3	0.0	0.0	27.6
AME-AT	5.8	5.8	5.8	6.3				
GRAMS	8.0	9.0	8.1	9.7				
NNWC	10.8	10.8	10.8	10.8		?		
NNNN	<b>5.4</b>	5.4	5.4	5.4				
RC	6.8	7.6	7.0	8.5				



GRAMS is the most effective attacker Even knowing the defender, NN cannot be evaded

#### Results – attack – evasion rates



	Nat.	rFGSM	BGA	Grosse	GRAMS	AGA	GwT	ENA
Natural	6.9	99.9	99.9	99.8	93.9	94.2	100.0	84.7
rFGSM	5.8	5.8	5.8	6.3	44.3	0.0	0.0	27.6
AME-AT	5.8	5.8	5.8	6.3	44.3	0.0	0.0	29.6
GRAMS	8.0	9.0	8.1	9.7	4.7	0.0	0.0	2.6
NNWC	10.8	10.8	10.8	10.8	3.7	0.0	0.0	2.8
NNNN	<b>5.4</b>	5.4	5.4	5.4	2.0	0.0	0.0	1.6
RC	6.8	7.6	7.0	8.5	27.2	0.0	0.0	6.2



GRAMS is the most effective attacker against unknown defenses

#### Results – attack – evasion rates



	Nat.	rFGSM	BGA	Grosse	GRAMS	AGA	GwT	ENA
Natural	6.9	99.9	99.9	99.8	93.9	94.2	100.0	84.7
rFGSM	5.8	5.8	5.8	6.3	44.3	0.0	0.0	27.6
AME-AT	5.8	5.8	5.8	6.3	44.3	0.0	0.0	29.6
GRAMS	8.0	9.0	8.1	9.7	4.7	0.0	0.0	2.6
NNWC	10.8	10.8	10.8	10.8	3.7	0.0	0.0	2.8
NNNN	5.4	5.4	5.4	5.4	2.0	0.0	0.0	1.6
RC	6.8	7.6	7.0	8.5	27.2	0.0	0.0	6.2



GRAMS is the most effective attacker even against unknown defenses

but natural has a higher evasion rate!

#### Results – defense – F1 scores





	Nat.	rFGSM	BGA	Grosse	GRAMS	AGA	GwT	ENA
Natural rFGSM	0.913 <b>0.921</b>	0.001 0.918	0.001 0.892	0.004 0.604	0.104 0.519	0.099 0.948	0.000 0.948	0.243 0.790
AME-AT	0.919	0.919	0.919	0.917				
GRAMS	0.905	0.899	0.904	0.895				
NNWC	0.880	0.880	0.880	0.880		?		
NNNN	0.883	0.883	0.883	0.883				
RC	0.918	0.914	0.917	0.909				

AME-AT obtains the best F1-scores against known attackers



#### Results – defense – F1 scores





	Nat.	rFGSM	BGA	Grosse	GRAMS	AGA	GwT	ENA
Natural rFGSM	0.913 <b>0.921</b>	0.001 0.918	0.001 0.892	0.004 0.604	0.104 0.519	0.099 0.948	0.000 0.948	0.243 0.790
AME-AT	0.919	0.919	0.919	0.917	0.670	0.949	0.949	0.778
GRAMS	0.905	0.899	0.904	0.895	0.922	0.946	0.946	0.933
NNWC	0.880	0.880	0.880	0.880	0.917	0.936	0.936	0.922
NNNN	0.883	0.883	0.883	0.883	0.901	0.910	0.910	0.903
RC	0.918	0.914	0.917	0.909	0.797	0.953	0.953	0.921



AME-AT obtains the best F1-scores against known attackers But GRAMS against unknown defenses!

#### Results – defense – F1 scores



	Nat.	rFGSM	BGA	Grosse	GRAMS	AGA	GwT	ENA
Natural rFGSM	0.913 <b>0.921</b>	0.001 0.918	0.001 0.892	0.004 0.604	0.104 0.519	0.099 0.948	0.000 0.948	0.243 0.790
AME-AT	0.919	0.919	0.919	0.917	0.670	0.949	0.949	0.778
GRAMS	0.905	0.899	0.904	0.895	0.922	0.946	0.946	0.933
NNWC	0.880	0.880	0.880	0.880	0.917	0.936	0.936	0.922
NNNN	0.883	0.883	0.883	0.883	0.901	0.910	0.910	0.903
RC	0.918	0.914	0.917	0.909	0.797	0.953	0.953	0.921

AME-AT obtains the best F1-scores against known attackers But GRAMS against unknown defenses!



but NNNN has the smallest evasion rates therefore ties with GRAMS in the defense track

#### Conclusion



- SLEIPNIR AdvML framework
  - Highlighy recommended for teaching AdvML
- AdvML challenge:
  - Constrained and integer adversarial attack and defense
  - Attackers do not know defenders
  - Defenders do not know attackers
- GRAMS
  - Effective and efficient, a simple multiple bit-flip search
  - Open source and available at:
    - https://github.com/tudelft-cda-lab/GRAMS
- NNNN (not our work)
  - Cannot be attacked by flipping bits to 1

