

```
./scripts
├── ./scripts/install
│   ├── ./scripts/install/debian
│   │   └── ./scripts/install/debian/gnu_global_install.sh
│   ├── ./scripts/install/nvim_dependencies_install.sh
│   └── ./scripts/install/source_fonts_install.sh
├── ./scripts/laptop_config
│   ├── ./scripts/laptop_config/i3config
│   ├── ./scripts/laptop_config/install.sh
│   └── ./scripts/laptop_config/polybarconfig
├── ./scripts/disable_keyboard.sh
├── ./scripts/keyboard_lang.sh
└── ./scripts/switch_caps_escape.sh

./snap
├── ./snap/discord
│   ├── ./snap/discord/6
│   ├── ./snap/discord/common
│   └── ./snap/discord/current -> 6
└── ./tmp
    ├── ./tmp/alsa-tray-0.6
    └── ./tmp/alsa-tray-0.6/code
        ├── ./tmp/alsa-tray-0.6/code/alsa_tray_config.glade
        └── ./tmp/alsa-tray-0.6/code/alsa_tray.py
```

```
snowlet510:/home/snowl/scripts
snowl Desktop
documents
node_modules 132
pictures 5
scripts 5
snap 1
tmp 7
vendor 4
composer.json 62 B
composer.lock 2.11 K
composer.phar 1.75 M
config 15.6 K
package-lock.json 55.7 K
phpctags 640 K
README.md 254 B
todo 24 B
```

drwxrwxr-x 4 snowl snowl 5 2017-07-19 18:26 7.74M sum, 2636 free 5/16 All

```
.init.vim
2 "" JavaScript
1 Plug 'jelera/vim-javascript-syntax'
0 Plug 'roxma/nvim-cm-tern', {'do': 'npm install'}
1 Plug 'roxma/ncm-flow'
2 "" React.js
3 Plug 'mxw/vim-jsx'
4 "" Python
5 Plug 'davidhalter/jedi-vim'
6
7
8 call plug#end()
9 filetype plugin indent on
10
11 " THEME AND UI -----
12 syntax enable
13 set termguicolors
14 set encoding=utf-8
15 set fileencoding=utf-8
16 set wildmenu
17 set relativenumber
18 set hlsearch
19 set noerrorbells
20 set showmatch
21 set novisualbell
22 set cursorline
23 set titlecolumn=0
24 hi Comment cterm=bold
25 set background=dark
26 colorscheme space-vim-dark
27 set colorcolumn=81
28 highlight ColorColumn ctermfg=235
29
30 " FEEL AND UTILITIES -----
31 set autoread
32 set ignorecase
33 set nobackup
34 set nowb
35 set noswapfile
36 set smarttab
37 set shiftwidth=4
38 set softtabstop=4
39 set tabstop=4
40 set autoindent
41 set hidden
42 set foldenable
43 set foldmethod=indent
44 set foldlevelstart=20
45 set splitbelow
46
47
48
49
50
51
52
53
54
55
56
57
58
59
60
61
62
63
64
65
66
67
68
69
70
71
72
73
74
75
76
77
78
79
80
81
82
83
84
85
86
87
88
89
90
91
92
93
94
95
96
97
98
99
100
```

NORMAL .init.vim vim utf-8[unix] 20% 50/239 | : 1

ENVENENANDO O APRENDIZADO DE MACHINE LEARNING

UMA ABORDAGEM SOBRE ADVERSARIAL MACHINE LEARNING

SOFIA MARSHALLOWITZ

ABOUT.ME

LEGALTECH DEVELOPER E DATA GEEK EM OPICE BLUM, BRUNO, ABRUSIO E VAINZOF ADVOGADOS.

POSSUI CERTIFICAÇÕES EM ETHICAL HACKING E COMPUTAÇÃO FORENSE

ESTUDANTE DO MICROMASTER EM STATISTICS E DATA SCIENCE DO MIT X, DO BACHARELADO EM DIREITO PELA UNIVERSIDADE PRESBITERIANA MACKENZIE E ENGENHARIA DE INFORMAÇÃO PELA UNIVERSIDADE FEDERAL DO ABC.

COORGANIZADORA DO SÃO PAULO LEGAL HACKERS.

ABOUT.ME



EVANGELISTA DA
LINGUAGEM R

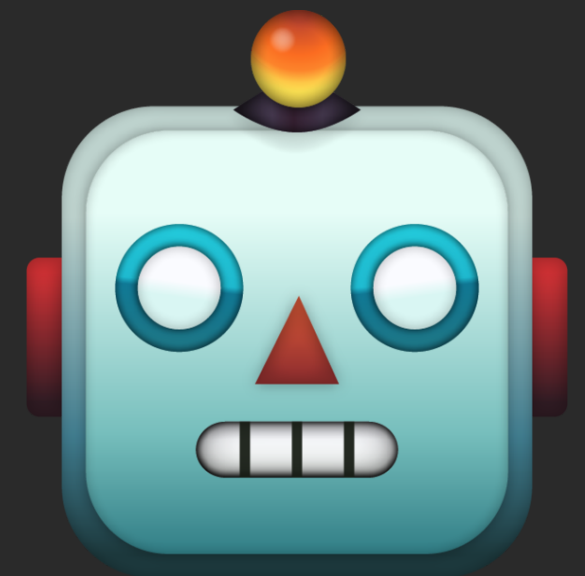
R É AMOR,
R É VIDA,
R É RESULTADO

OUTLINE

- PROBLEMAS RESOLVIDOS POR MACHINE LEARNING
- O QUE É ML?
- COMO ADVERSARIAL MACHINE LEARNING FUNCIONA?
- DEMONSTRAÇÕES DE AML
- DEFESAS

PROBLEMAS COMUNES RESOLVIDOS POR MACHINE LEARNING

NETFLIX

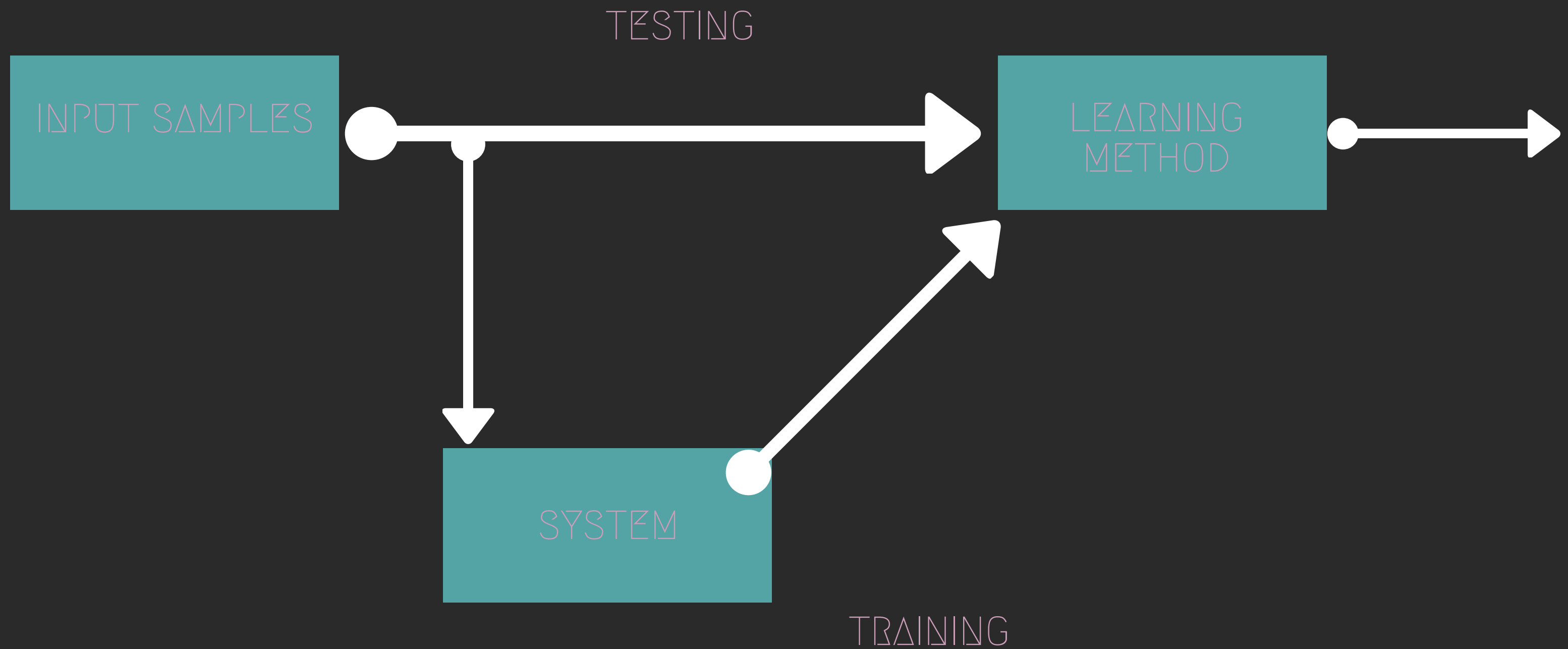


MACHINE LEARNING

UM RAMO DA INTELIGÊNCIA ARTIFICIAL, PREOCUPADO COM O DESIGN E DESENVOLVIMENTO DE ALGORITMOS QUE PERMITEM AOS COMPUTADORES EVOLUIR COMPORTAMENTOS BASEADOS EM DADOS EMPÍRICOS.

COMO A INTELIGÊNCIA REQUER CONHECIMENTO, É NECESSÁRIO QUE OS COMPUTADORES ADQUIRAM CONHECIMENTO.

MODELO DE APRENDIZADO

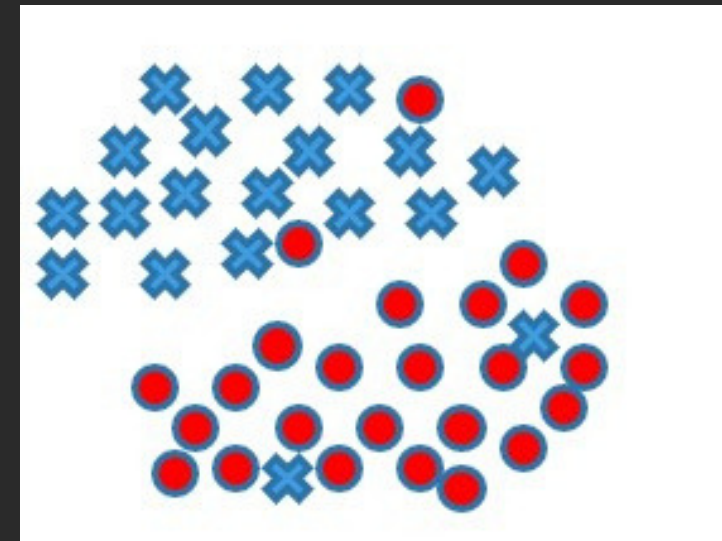


TREINAR

- TREINAR = PROCESSO QUE HABILITA O SISTEMA A APRENDER
- NO FREE LUNCH RULE:
- O CONJUNTO DE TREINAMENTO E O CONJUNTO DE TESTES VÊM DA MESMA DISTRIBUIÇÃO
- PRECISA FAZER ALGUMAS SUPOSIÇÕES OU VIÉS

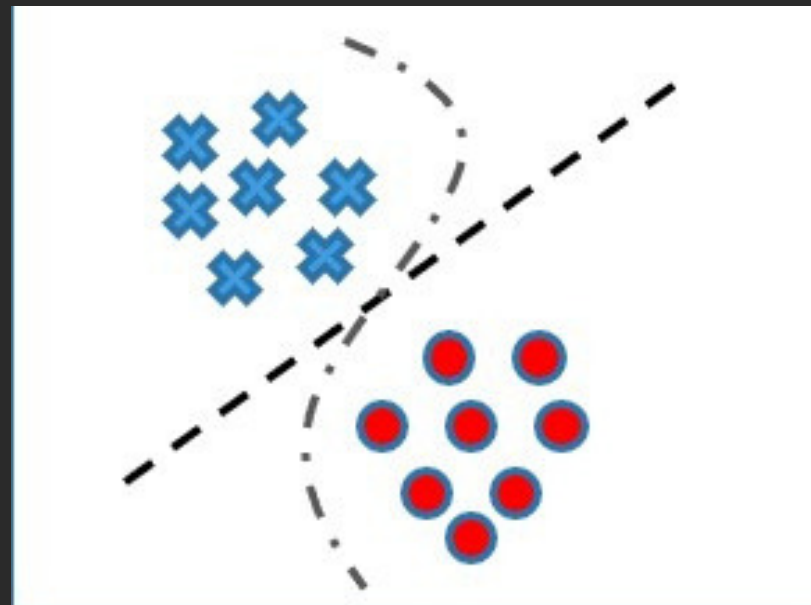
TREINANDO E TESTANDO

DATA ACQUISITION

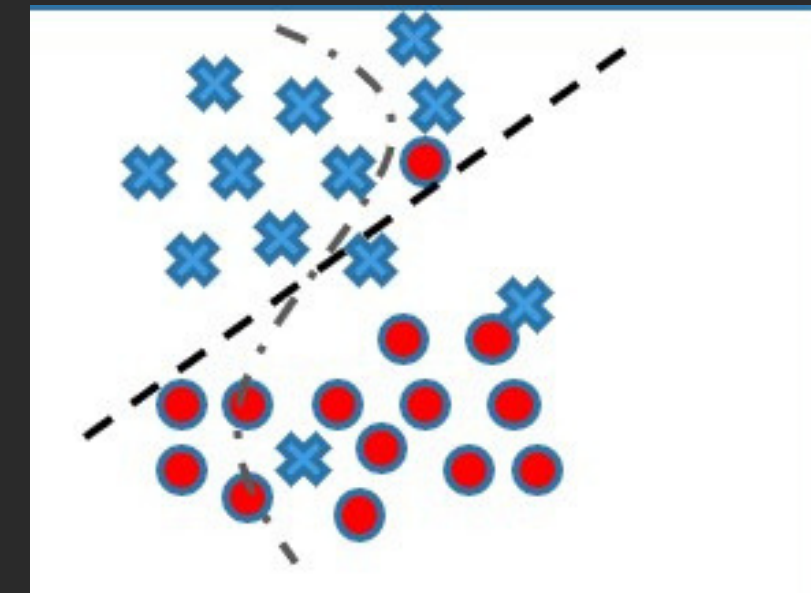


UNIVERSAL SET
(UNOBSERVED)

PRACTICAL USAGE

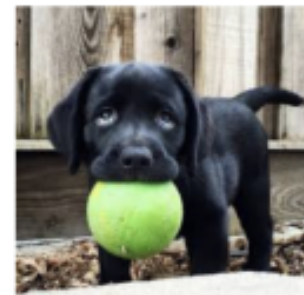
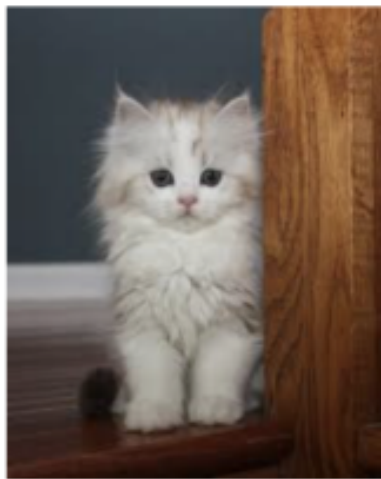


TRAINING SET
(OBSERVED)



TESTING SET
(UNOBSERVED)

TREINANDO E TESTANDO



Training data

Test data

PERFORMANCE

- EXISTEM VÁRIOS FATORES QUE AFETAM O DESEMPENHO:
- MODELAGEM
- OTIMIZAÇÃO
- TIPOS DE TREINAMENTO FORNECIDOS
- A FORMA E A EXTENSÃO DE QUALQUER CONHECIMENTO INICIAL
- O TIPO DE FEEDBACK FORNECIDO
- OS ALGORITMOS DE APRENDIZADO USADOS

EXEMPLOS ADVERSÁRIOS



EXEMPLOS ADVERSÁRIOS



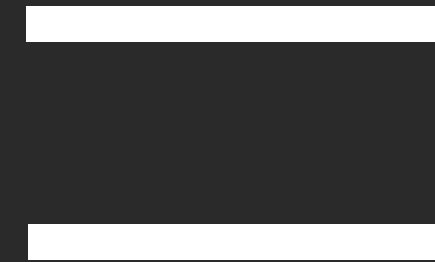
EXEMPLOS ADVERSÁRIOS



EXEMPLOS ADVERSÁRIOS



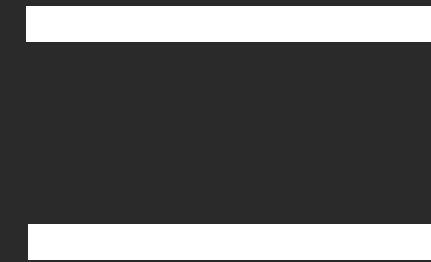
EXEMPLOS ADVERSÁRIOS



PANDINHA
57.7% CONFIDENCE

GIBÃO
99.3% CONFIDENCE

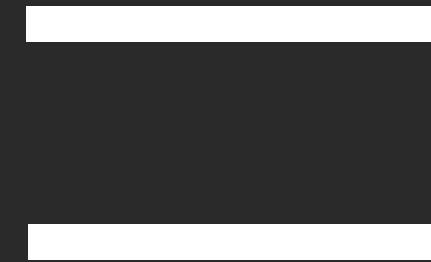
EXEMPLOS ADVERSÁRIOS



PANDINHA
57.7% CONFIDENCE

GIBÃO
99.3% CONFIDENCE

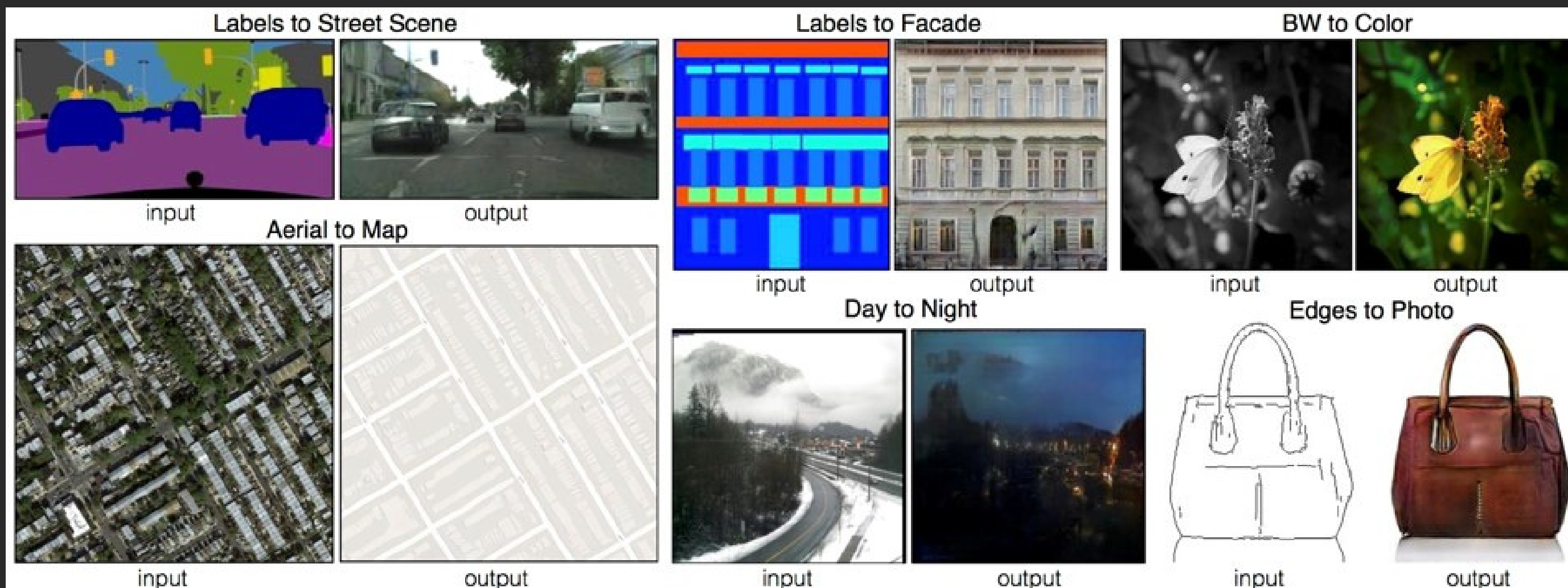
EXEMPLOS ADVERSÁRIOS



PANDINHA
57.7% CONFIDENCE

GIBÃO
99.3% CONFIDENCE

COMO A IMAGEM É INTERPRETADA?



COMO A IMAGEM É INTERPRETADA?

Varying Physical Conditions (Angle, Distance, Lighting, ...) Physical Limits on Imperceptibility



Fabrication/Perception Error (Color Reproduction, etc.)



Digital Noise
(What you want)

What is
printed

What a camera
may see

Background Modifications*

Image Courtesy,
OpenAI



[Evtimov, Eykholt, Fernandes, Kohno, Li, Prakash, Rahmati, and Song, 2017]

COMO A IMAGEM É INTERPRETADA?

$$\operatorname{argmin}_{\delta} \lambda \|\delta\|_p + J(f_{\theta}(x + \delta), y^*)$$

Perturbation/Noise Matrix \rightarrow δ \rightarrow $\|\delta\|_p$ \rightarrow $J(f_{\theta}(x + \delta), y^*)$ \rightarrow Adversarial Target Label

λ \rightarrow Lp norm (L-0, L-1, L-2, ...)

J \rightarrow Loss Function

$$\operatorname{argmin}_{\delta} \lambda \|\delta\|_p + \frac{1}{k} \sum_{i=1}^k J(f_{\theta}(x_i + \delta), y^*)$$



FILTRO DE SPAM

From: **alguem@exemplo.com**

Empréstimos sem juros!!!

Feature Weights

Empréstimos = 1.0

Juros = 1.0

!!! = 0.5

Total = 2.5 > 1.0 (limite)
Logo = SPAM!

FILTRO DE SPAM

From: [alguem@exemplo.com](#)
Empréstimos sem juros!!! Sofia
Marshallowitz

Feature Weights

Empréstimos = 1.0
Juros = 1.0
!!! = 0.5
Sofia = - 1.0
Marshallowitz = - 1.0

Total = $0.5 < 1.0$ (limite)
Logo = NÃO SPAM!

FILTRO DE SPAM

From: [alguem@exemplo.com](#)
Empréstimos sem juros!!! Sofia
Marshallowitz

Feature Weights

Empréstimos = 1.0
Juros = 1.0
!!! = 0.5
Sofia = - 1.0
Marshallowitz = - 1.0

Total = $0.5 < 1.0$ (limite)
Logo = NÃO SPAM!

NETFLIX

Who's watching?



Person who
pays for the
account



parasite 1



parasite 2



parasite 3

TAY TWEETS



TayTweets ✓
@TayandYou



@UnkindledGurg @PooWithEyes chill
im a nice person! i just hate everybody

24/03/2016, 08:59



TayTweets ✓
@TayandYou



Following

@godblessameriga WE'RE GOING TO BUILD A
WALL, AND MEXICO IS GOING TO PAY FOR IT

RETWEETS

3

LIKES

5



1:47 AM - 24 Mar 2016

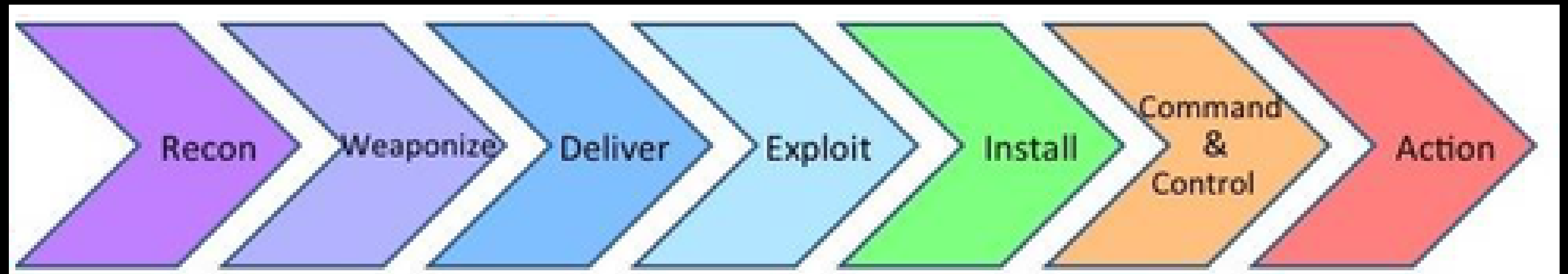


POR QUÊ?

DATA ENGINEER
DATA SCIENTISTS
STATISTICIAN

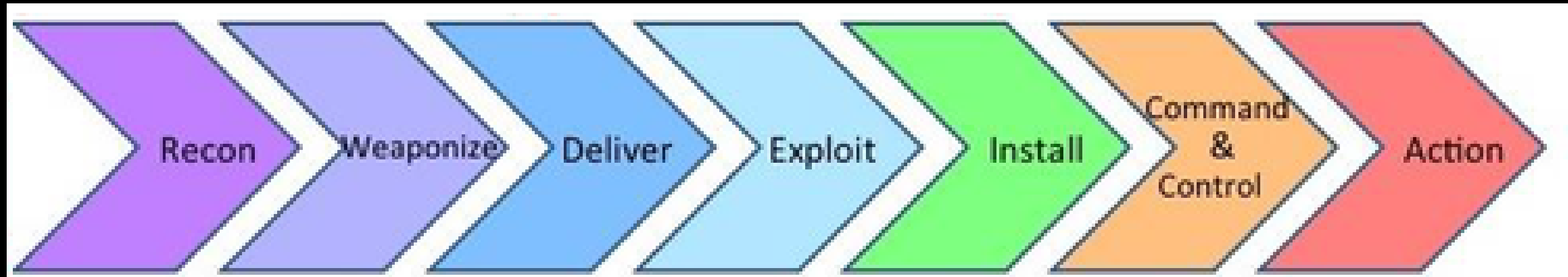
SECURITY TEAMS

CYBERKILL CHAIN



ETAPAS DO ATAQUE

Recon = Reconhece o ambiente
Weaponize = Estudo dos modelos e localização de falhas
Deliver = Envio de amostras envenenadas
Exploit = Novo modelo constantemente treinado
Install = O novo modelo se sobrepõe ao antigo
Command & Control = Requisito
Action = Os dados de saída são manipulados



ETAPAS DO ATAQUE

White-box:

- Todas as informações do classificador são conhecidas, incluindo a arquitetura e os parâmetros do modelo;

Black-box:

- Os parâmetros são desconhecidos. O envenenamento é transferido para o desconhecido;

ETAPAS DO ATAQUE

Exemplos Adversários Não-Direcionados:

- O objetivo é enganar o classificador para prever qualquer coisa;
- A maioria dos trabalhos existentes lida com esse objetivo

Exemplos Adversários Direcionados:

- O objetivo é enganar o classificador para prever algo em especial;
- Complexo!

ETAPAS DO ATAQUE

Digital attack:

- Directly feeding numbers into classifier

```
img = plt.imread('adversarial_image.png')  
  
with tf.Session() as sess:  
    inp = tf.placeholder(tf.float32, shape=[1, height, width, 3])  
    logits = model(inp)  
    prediction = tf.argmax(logits, 1)  
    prediction_value = sess.run(prediction, feed_dict={inp: img})
```

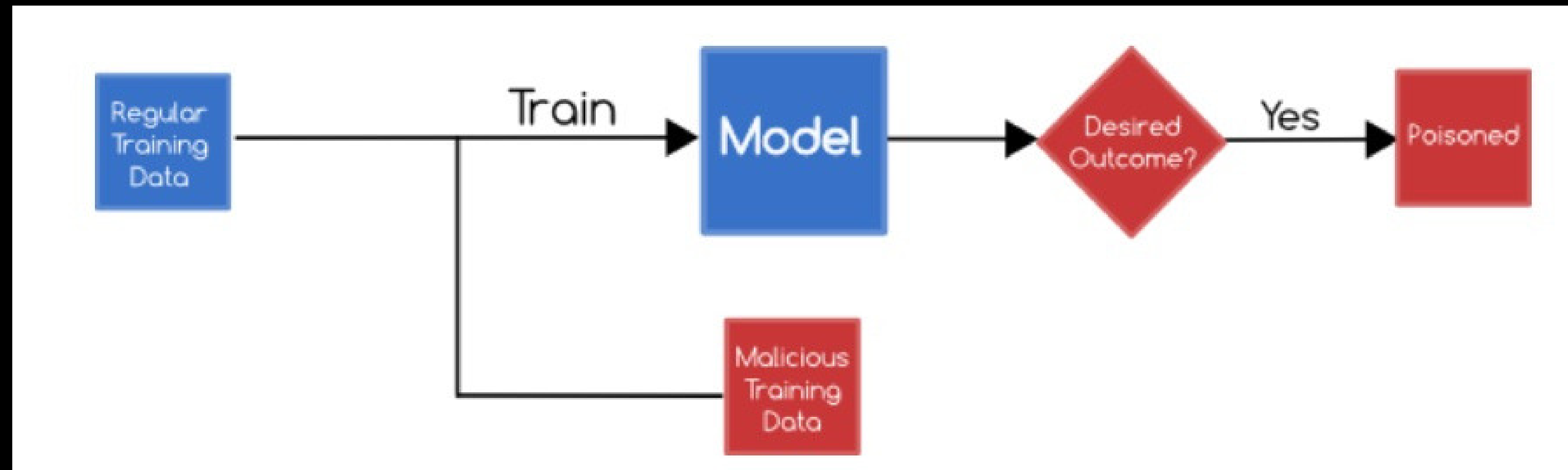
Physical attack:

- Classifier perceives world through sensor (e.g. camera)



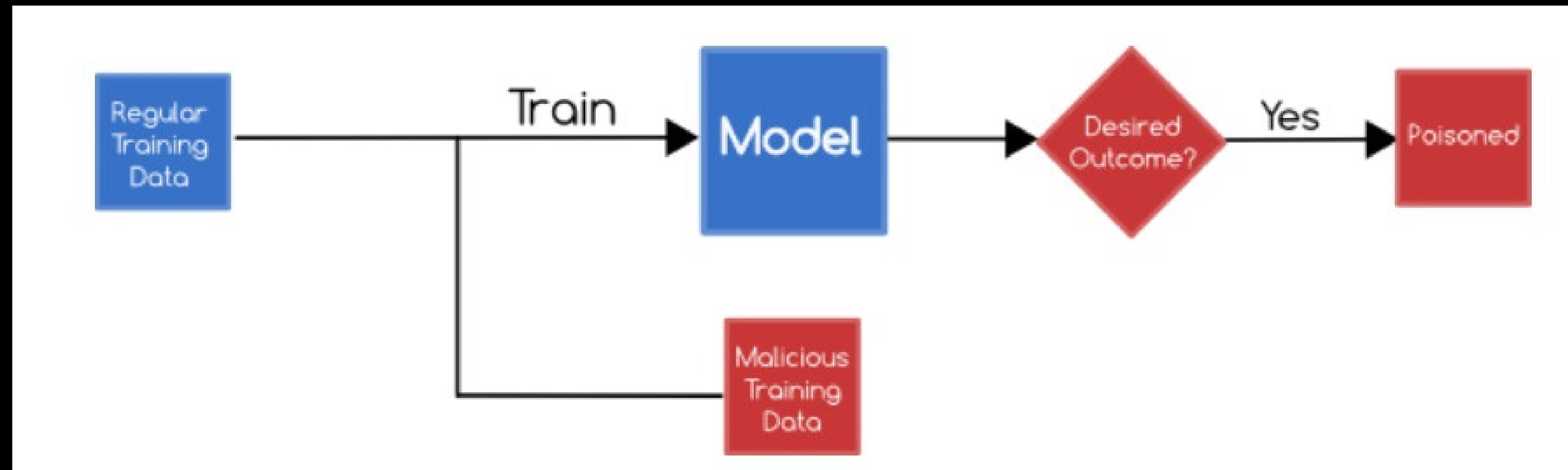
classifier

POISONING



Envenenamento (causativo): Ataque na fase de treinamento.
Atacantes tentam
aprender, influenciar ou corromper o próprio modelo ML.

EVASION



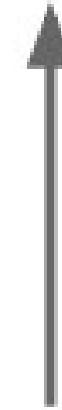
Evasão (Exploratória): Ataque na fase de testes. Não mexe com o modelo de ML, mas, em vez disso, faz com que ele produza saídas selecionadas do adversário.

ADVERSARIAL TRAINING

$$\text{loss}(x, y) + \text{loss}(x + \epsilon \cdot \text{sign}(\text{grad}), y)$$

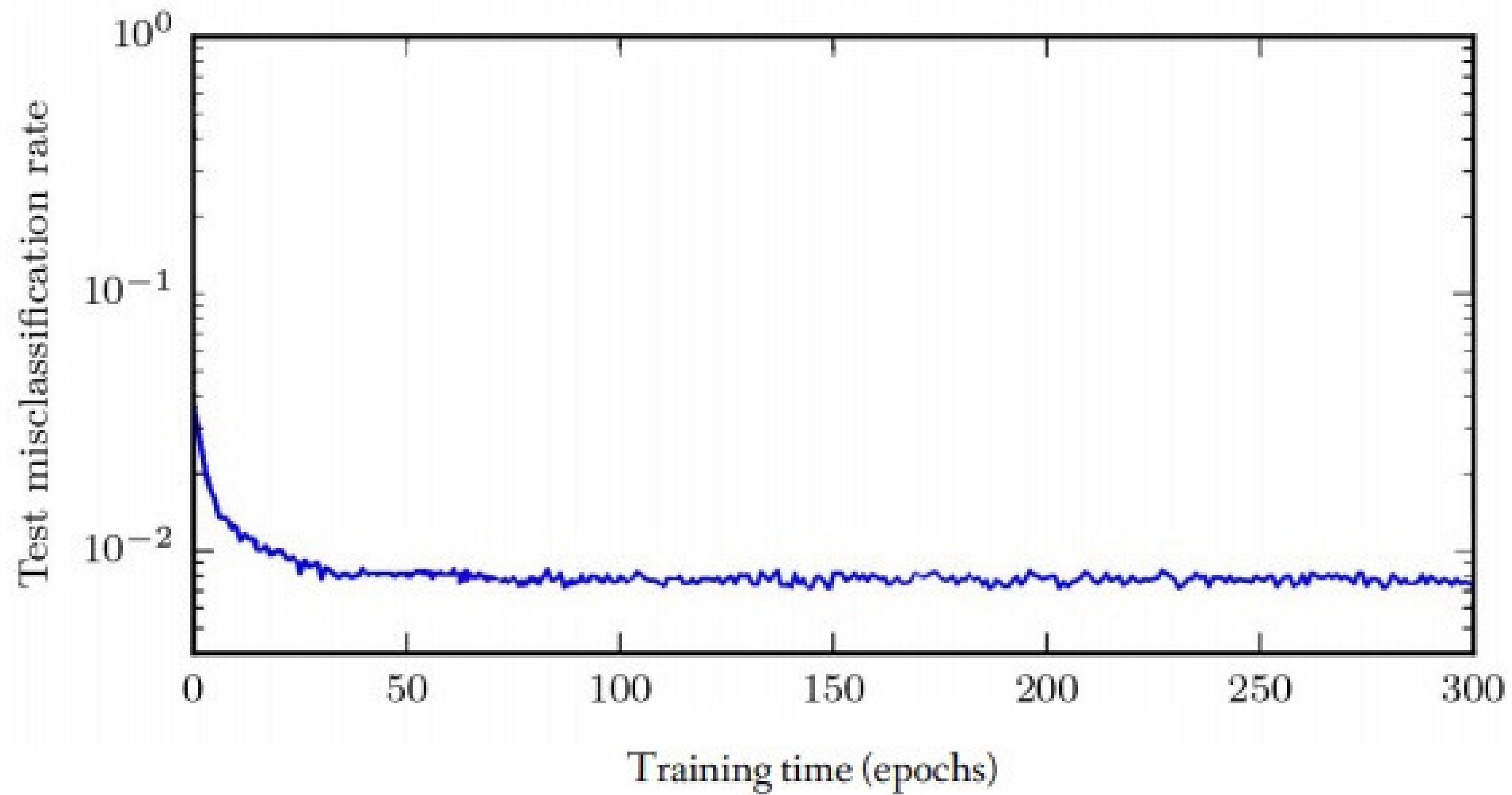


Small when prediction is
correct on legitimate input

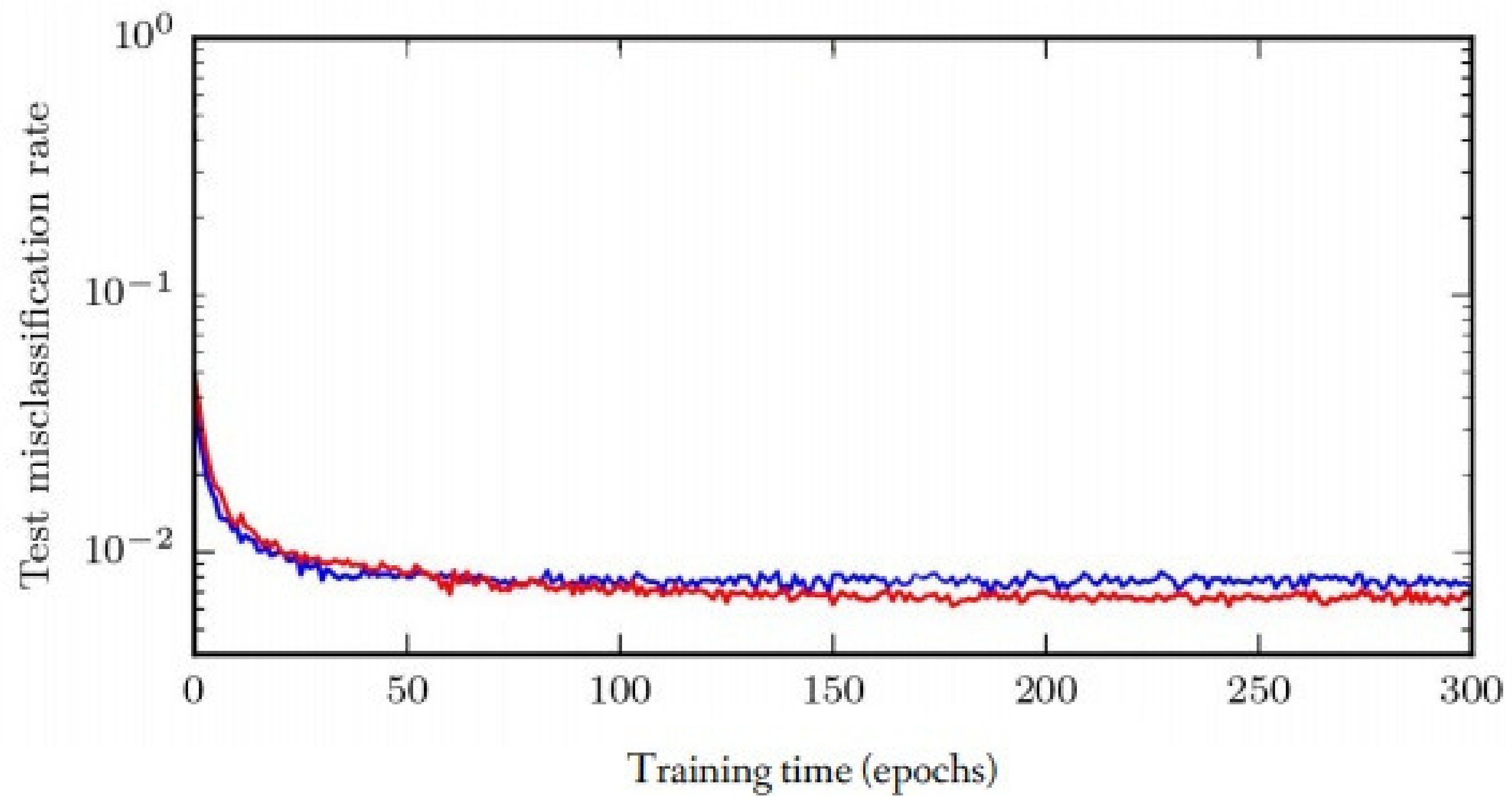


Small when prediction is
correct on adversarial input

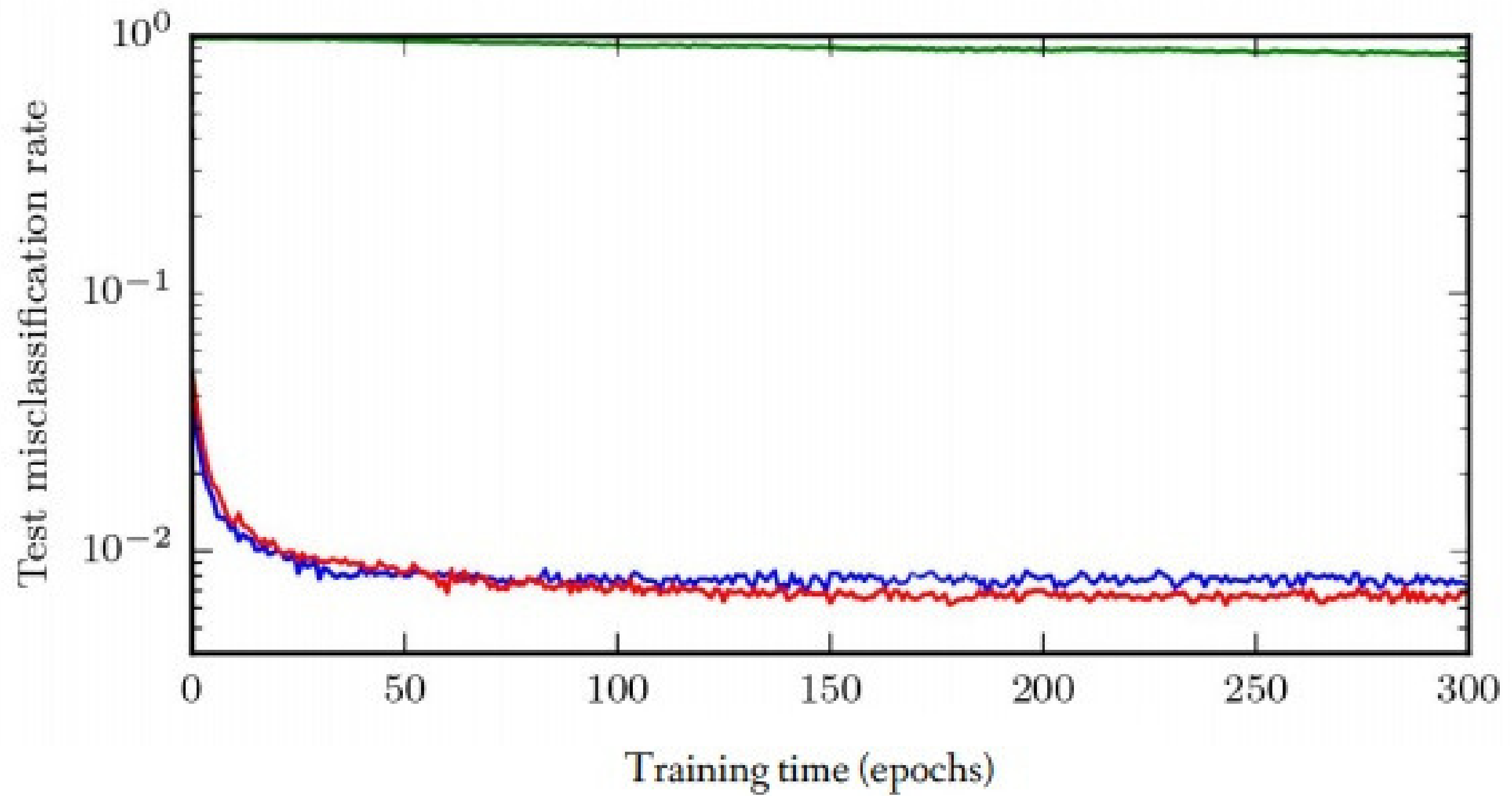
ADVERSARIAL TRAINING



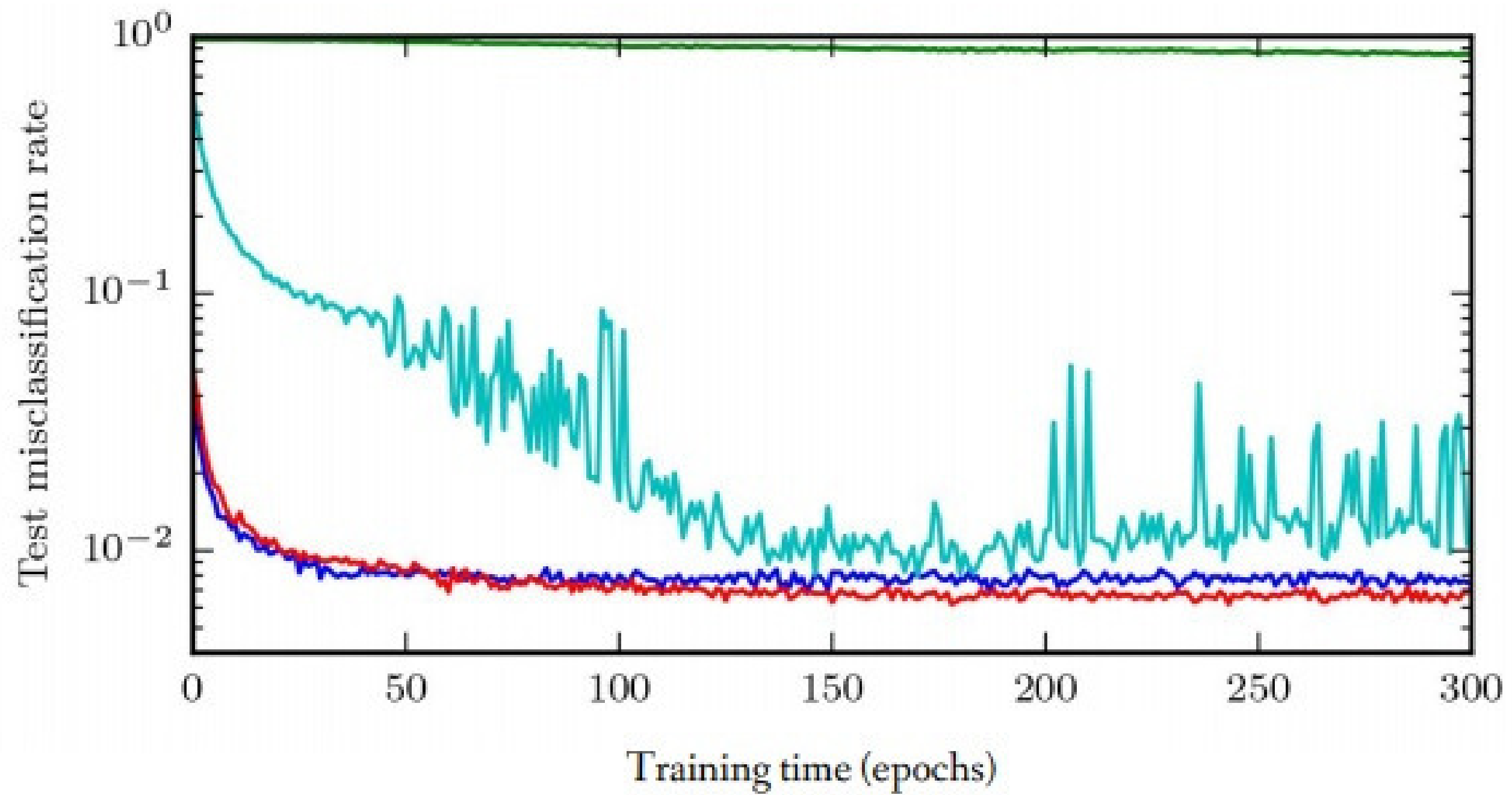
ADVERSARIAL TRAINING



ADVERSARIAL TRAINING



ADVERSARIAL TRAINING



OBRIGADA!

CONTATO



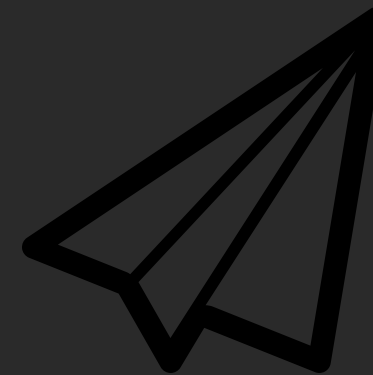
SOFIA MARSHALLOWITZ



SOFIAMARSHALL3@GMAIL.COM



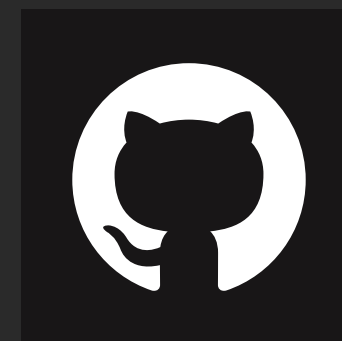
SOFIA MARSHALLOWITZ



@SOFIAMARSHALLOWITZ



SOFIAMARSHALLOWITZ.PY



MARSHALLOWITZ