# Robust Models for Source Code: Techniques and Applications

**Pavol Bielik**
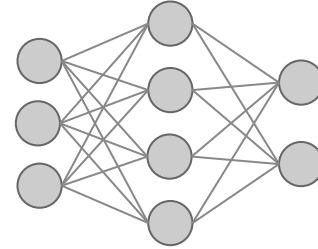pavol.bielik@inf.ethz.ch

Department of Computer Science
**ETH** *zürich*

SRILAB

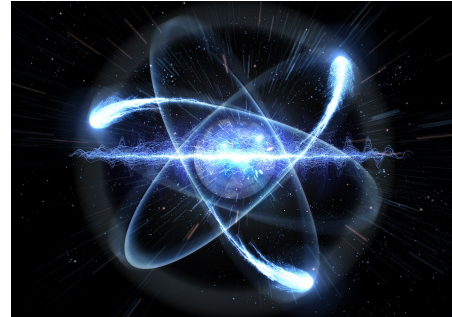# Secure, Reliable, and Intelligent Systems Lab @ ETH



**Robust ML**



**Neural network verification**



**ML for programming**



**Probabilistic programming**

# Secure, Reliable, and Intelligent Systems Lab @ ETH



**Robust ML**



Neural network verification



**ML for programming**



Probabilistic programming
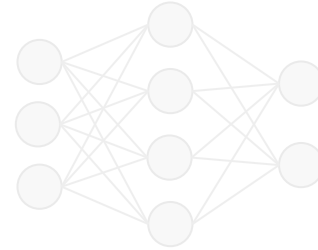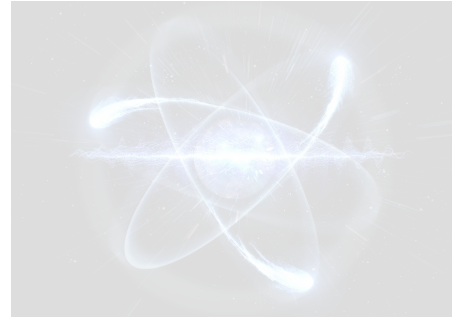
# Statistical Programming Tools

## Code Completion

```
Camera camera = Camera.open();
camera.SetDisplayOrientation(90);
              ?
```
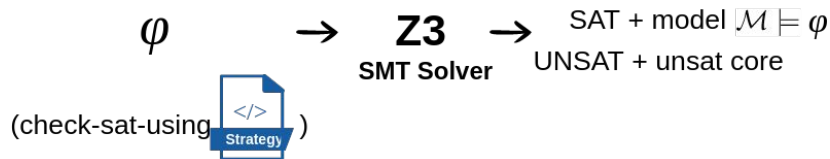
## Port Code

| C# | | Java | Translate |
|---|---|---|---|
| Console.WriteLine("Hi"); | | System.out.println("Hi"); | |
| ... | | ... | |

## Program Synthesis

## Learning to Solve Formulas

$\varphi \quad \rightarrow \quad$ **Z3** $\rightarrow$ SAT + model $\mathcal{M} \models \varphi$

**SMT Solver**    UNSAT + unsat core

(check-sat-using   Strategy   )

**Up to 100x speed-up over Z3**

## Bug Detection

```
...                    likely error
for x in range(a):
    print a[x]
```

# Deep Learning + Code

Malware Detection

Bug Detection     Loop Invariants     Bug Repair

Code Classification     Code Search     Type Inference     Neural Decompilation

Code Captioning     Code Completion     Variable Naming     Program Translation

| 2016 | 2017 | 2018 | 2019 |

# Deep Learning + Code

| 2016 | 2017 | 2018 | 2019 |
|------|------|------|------|
| | Bug Detection | Loop Invariants | Malware Detection |
| Code Classification | Code Search | Type Inference | Bug Repair |
| Code Captioning | Code Completion | Variable Naming | Neural Decompilation |
| | | | Program Translation |

**Techniques are general and apply to other tasks**

# Deep Learning + Code

Malware Detection

Bug Detection     Loop Invariants     Bug Repair

Code Classification     Code Search     Type Inference     Neural Decompilation

Code Captioning     Code Completion     Variable Naming     Program Translation

| 2016 | 2017 | 2018 | 2019 |

↑

## Majority is based on deep learning models

## Techniques are general and apply to other tasks
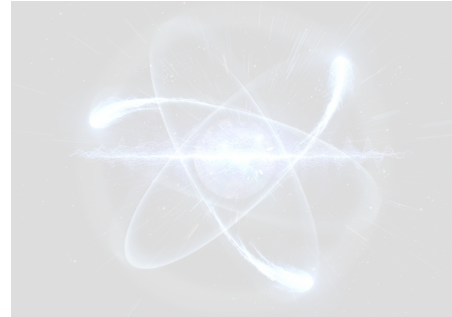
# Secure, Reliable, and Intelligent Systems Lab @ ETH
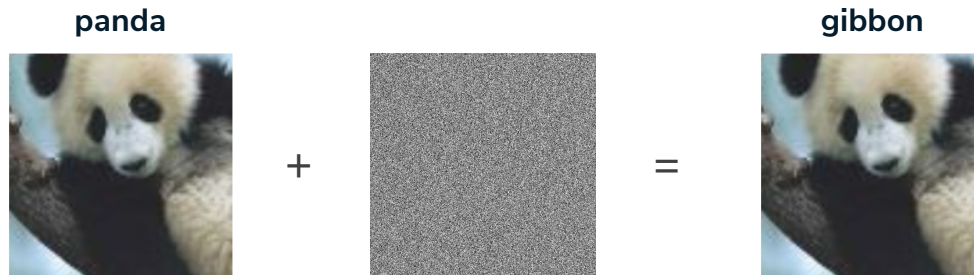


**Robust ML**



**Neural Network Verification**



**ML for programming**



**Probabilistic programming**

# Adversarial Robustness

Vision

**panda**                    **gibbon**



+                    =

Explaining and Harnessing Adversarial Examples. Goodfellow et. al. ICLR'15

# Adversarial Robustness

Vision

**stop sign**
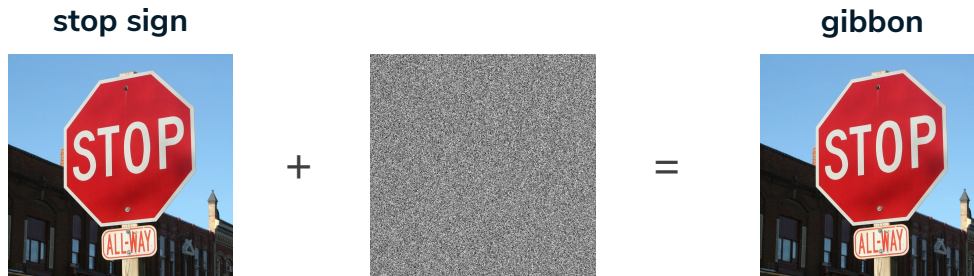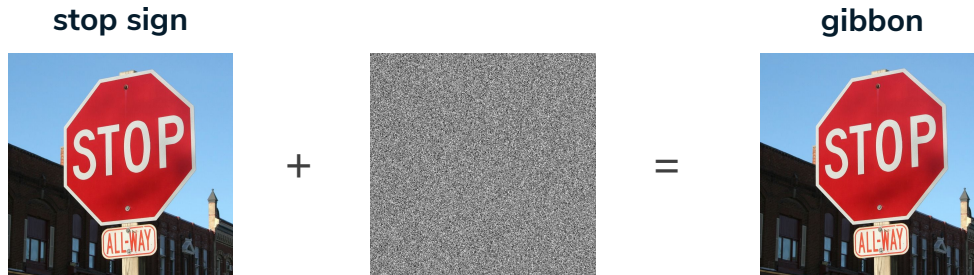
 +  = 

**gibbon**

# Adversarial Robustness



**Vision**

stop sign + = gibbon

**Sound**

+ **noise** =

Audio Adversarial Examples: Targeted Attacks on Speech-to-Text. Carlini et. al. ICML'18 workshop

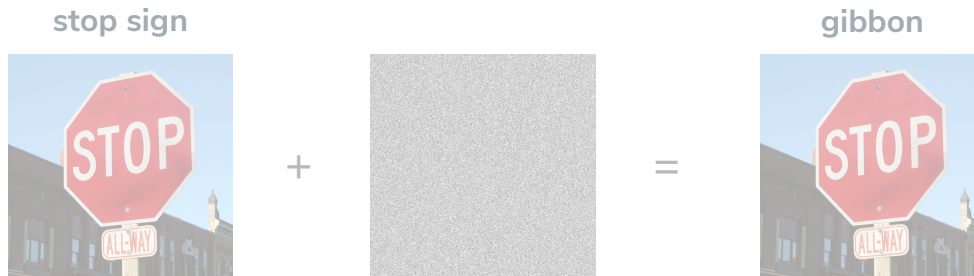# Adversarial Robustness

**stop sign**                                **gibbon**

👁 Vision

**Sound**

Audio Adversarial Examples: Targeted Attacks on Speech-to-Text. Carlini et. al. ICML'18 workshop

</> Code       +    **code refactoring**    =

# Adversarial Robustness



stop sign + = gibbon

◉ Vision

〰 Sound + noise =

Audio Adversarial Examples: Targeted Attacks on Speech-to-Text. Carlini et. al. ICML'18 workshop

</> Code + **code refactoring** =

# Adversarial Robustness for Code

( 1 )   **How robust are existing models?**

</> Code    | **90%** accuracy |  +  | **code refactoring** |  =  | **??** robustness |

# Adversarial Robustness for Code

( 1 )  **How robust are existing models?**

</> Code

$$90\% \text{ accuracy} + \textbf{code refactoring} = 0\% \text{ robustness}$$

# Adversarial Robustness for Code

**(1)** **How robust are existing models?**

**(2)** **How to find adversarial examples?**

</> Code

| **90%** accuracy | + | **code refactoring** | = | **0%** robustness |

# Adversarial Robustness for Code

**(1)** **How robust are existing models?**

**(2)** **How to find adversarial examples?**

**(3)** **How to improve robustness?**

</> Code

**90%**
accuracy

+

**code refactoring**

=

**0%**
robustness

# Adversarial Robustness Example



Input Program
**x**

```
...
v = parseInt(
  hex.substr(1),
  radix
)
...
```

Model
f(**x**) → **y**

(Type Inference)
Program Properties
**y**

```
...
v[num] = parseInt[num](
  hex[str].substr[str](1),
  radix[num]
)
...
```
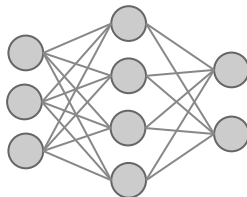
# Adversarial Robustness Example

**Input Program**
**x**

```
...
v = parseInt(
  hex.substr(1),
  radix
)
...
```

**Model**
f(**x**) → **y**

**(Type Inference)**
**Program Properties**
**y**

```
...
v[num] = parseInt[num](
  hex[str].substr[str](1),
  radix[num]
)
...
```

**Goal (Adversarially Robustness):**
**Model is correct for *all* label preserving program transformations**

```
...
v = parseInt(
  color.substr(1),
  radix
)
...
```
variable renaming

```
...
v = parseInt(
  hex.substr(42),
  radix
)
...
```
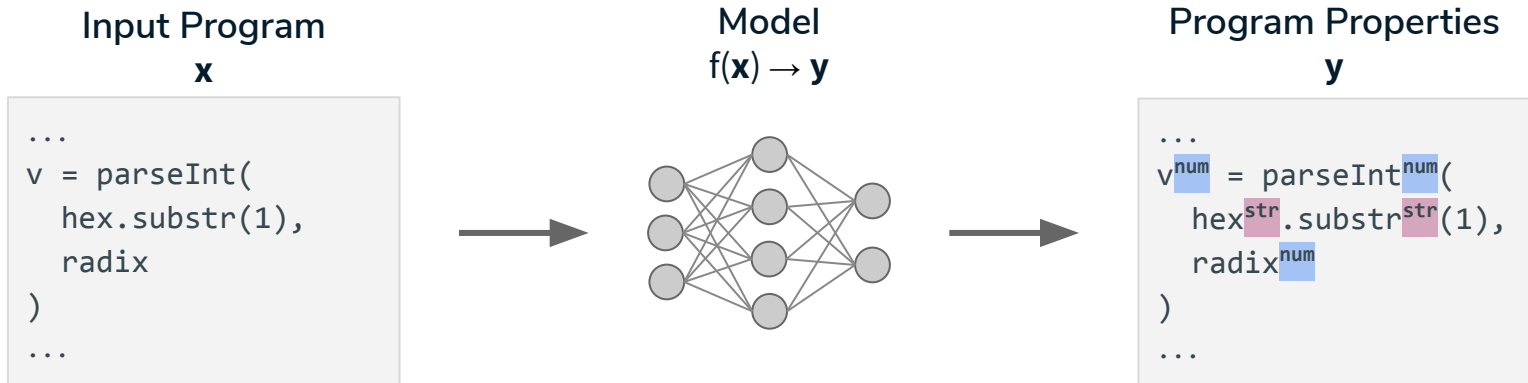constant replacement

```
...
v = parseInt(
  hex.substr(1),
  radix + 0
)
...
```
semantic equivalence

```
...
parseInt(
  hex.substr(1),
  radix
)
...
```
remove assignment

# Adversarial Robustness Example

**Input Program**
**x**

```
...
v = parseInt(
  hex.substr(1),
  radix
)
...
```

**Model**
f(**x**) → **y**

**(Type Inference)**
**Program Properties**
**y**

```
...
v^num = parseInt^num(
  hex^str.substr^str(1),
  radix^num
)
...
```

## Goal (Adversarially Robustness):
## Model is correct for *all* label preserving program transformations

```
...
v = parseInt(
  color.substr(1),
  radix
)
...
```
variable renaming

```
...
v = parseInt(
  hex.substr(42),
  radix
)
...
```
constant replacement

**S(x)**

```
...
v = parseInt(
  hex.substr(1),
  radix + 0
)
...
```
semantic equivalence

Set of valid program transformations for **x**

```
...
parseInt(
  hex.substr(1),
  radix
)
...
```
remove assignment

# How robust are existing models?

Type Inference [1]

```
...
v^num = parseInt^num(
  hex^str.substr^str(1),
  radix^num
)
...
```

[1] Adversarial Robustness for Code. ICML'20

# How robust are existing models?

### Type Inference [1]

```
...
v^num = parseInt^num(
    hex^str.substr^str(1),
    radix^num
)
...
```

### Code Captioning [2,3]

```
int indexOfTarget(Object target) {
    int i = 0;
    for (Object elem: this.elements) {
        if (elem.equals(target))
            return i
        i++;
    return -1;
}
```

[1] Adversarial Robustness for Code. ICML'20
[2] Adversarial Examples for Models of Code. ArXiv'19
[3] Semantic Robustness of Models of Source Code. ArXiv'20

# How robust are existing models?

### Type Inference [1]

```
...
v^num = parseInt^num(
  hex^str.substr^str(1),
  radix^num
)
...
```

### Code Captioning [2,3]

```
int indexOfTarget(Object target) {
  int i = 0;
  for (Object elem: this.elements) {
    if (elem.equals(target))
      return i;
    i++;
  return -1;
}
```

### Malware Detection [4]

```
push ebp
mov ebp, esp
push ebx
push edx                    yes
mov ebx, [ebp+4]            no
add ebx, 0x10
mov edx, [ebp+8]
mov [edx], ebx
```

[1] Adversarial Robustness for Code. ICML'20
[2] Adversarial Examples for Models of Code. ArXiv'19
[3] Semantic Robustness of Models of Source Code. ArXiv'20
[4] Optimization-Guided Binary Diversification to Mislead Neural Networks for Malware Detection. ArXiv'19

# How robust are existing models?

### Type Inference [1]

```
...
v^num = parseInt^num(
    hex^str.substr^str(1),
    radix^num
)
...
```

| 89% | 48% |
|:---:|:---:|
| accuracy | robustness |

### Code Captioning [2,3]

```
int indexOfTarget(Object target) {
    int i = 0;
    for (Object elem: this.elements) {
        if (elem.equals(target))
            return i;
        i++;
    }
    return -1;
}
```

| 39.2 | 19.6 |
|:---:|:---:|
| F1 | robust F1 |

### Malware Detection [4]

```
push ebp
mov ebp, esp
push ebx
push edx                        yes
mov ebx, [ebp+4]                no
add ebx, 0x10
mov edx, [ebp+8]
mov [edx], ebx
```

| 99% | ~2% |
|:---:|:---:|
| accuracy | robustness |

[1] Adversarial Robustness for Code. ICML'20
[2] Adversarial Examples for Models of Code. ArXiv'19
[3] Semantic Robustness of Models of Source Code. ArXiv'20
[4] Optimization-Guided Binary Diversification to Mislead Neural Networks for Malware Detection. ArXiv'19

# How robust are existing models?

Type Inference [1]

Code Captioning [2,3]

Malware Detection [4]

```
int indexOfTarget(Object target) {
```

**Not trained with robustness in mind**

```
        i++;
    return -1;
}
```

```
push ebp
```

```
add ebx, 0x10
mov edx, [ebp+8]
mov [edx], ebx
```

| 89% | 48% |
|-----|-----|
| accuracy | robustness |

| 39.2 | 19.6 |
|------|------|
| F1 | robust F1 |

| 99% | ~2% |
|-----|-----|
| accuracy | robustness |

[1] Adversarial Robustness for Code. ICML'20
[2] Adversarial Examples for Models of Code. ArXiv'19
[3] Semantic Robustness of Models of Source Code. ArXiv'20
[4] Optimization-Guided Binary Diversification to Mislead Neural Networks for Malware Detection. ArXiv'19

# So Far...

original programs

```
...
v^num = parseInt^num(
  hex^str.substr^str(1),
  radix^num
)
...
```

$\delta$: rename
hex $\longrightarrow$ color

transformed programs

```
...
v^num = parseInt^num(

color^num.substr^num(1),
  radix^num
)
...
```

**89%**
accuracy

**48%**
robustness

# Adversarial Training

original programs



$\delta$: rename
hex $\rightarrow$ color

transformed programs



**Standard training**

min $loss(\theta, \mathbf{x}, \mathbf{y})$

measures the model performance

**Adversarial training**

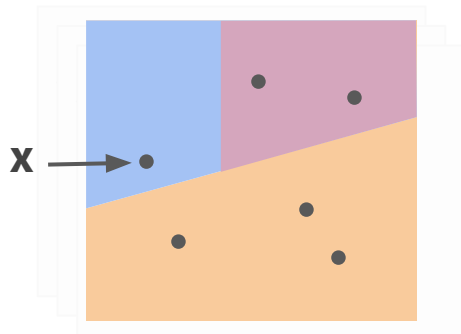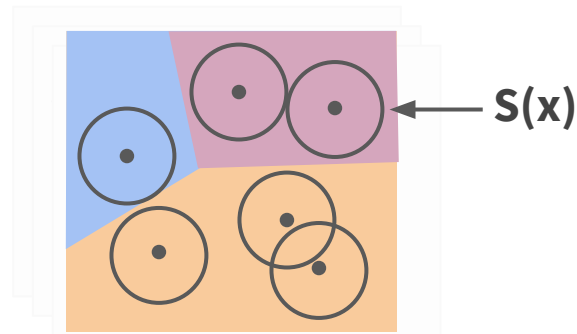min $[$max $loss(\theta, \mathbf{x} + \boldsymbol{\delta}, \mathbf{y})]$
$\boldsymbol{\delta} \in \mathbf{S(x)}$

program transformations

# Adversarial Training

original programs



$\delta$: rename
hex $\rightarrow$ color

generated programs

**Standard training**

$\min loss(\theta, \mathbf{x}, \mathbf{y})$

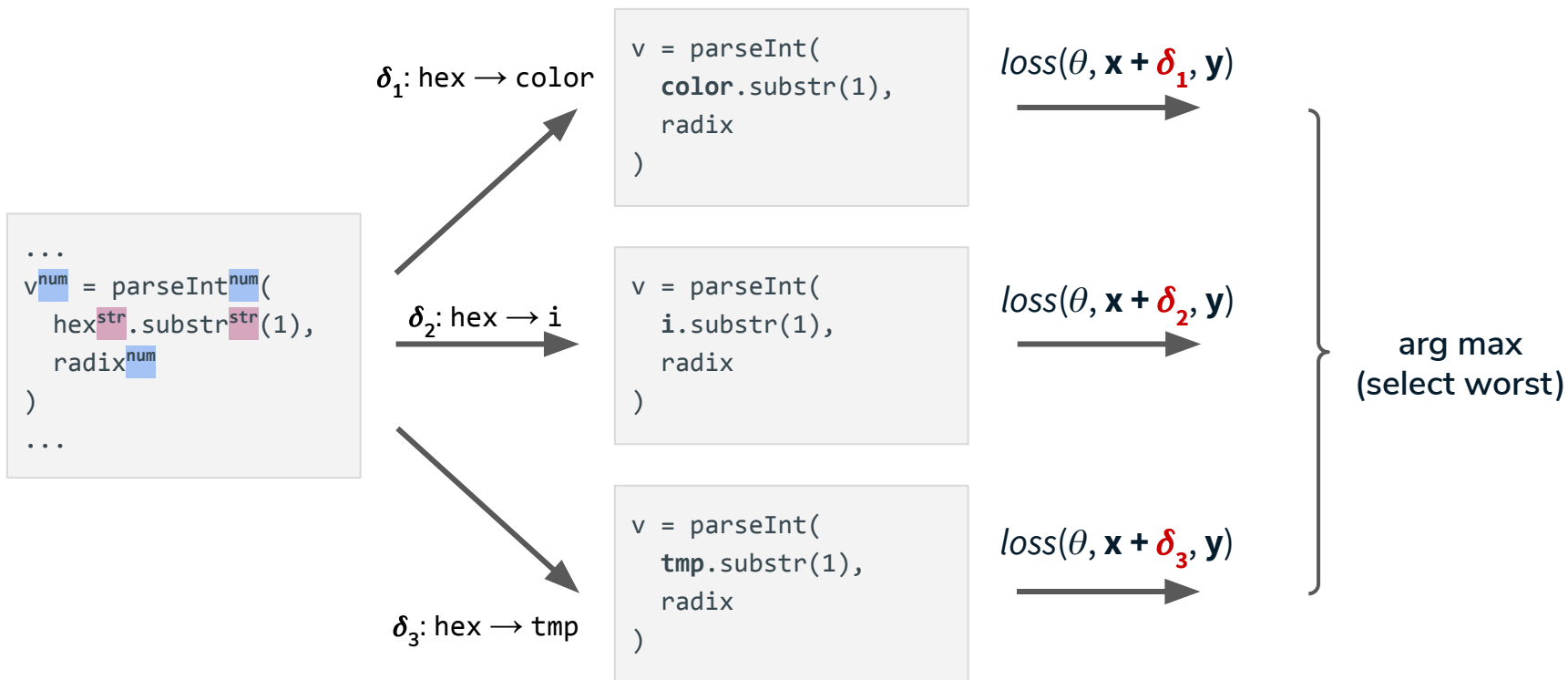**Adversarial training**

$\min [\max loss(\theta, \mathbf{x} + \delta, \mathbf{y})]$
$\delta \in \mathbf{S}(\mathbf{x})$

**trains on individual programs**

**trains on worst case generated programs**

# Finding Adversarial Examples



$\delta_1$: hex $\longrightarrow$ color

```
v = parseInt(
  color.substr(1),
  radix
)
```

$loss(\theta, \mathbf{x} + \boldsymbol{\delta_1}, \mathbf{y})$

```
...
v^num = parseInt^num(
  hex^str.substr^str(1),
  radix^num
)
...
```

$\delta_2$: hex $\longrightarrow$ i

```
v = parseInt(
  i.substr(1),
  radix
)
```

$loss(\theta, \mathbf{x} + \boldsymbol{\delta_2}, \mathbf{y})$

$\delta_3$: hex $\longrightarrow$ tmp

```
v = parseInt(
  tmp.substr(1),
  radix
)
```

$loss(\theta, \mathbf{x} + \boldsymbol{\delta_3}, \mathbf{y})$

arg max
(select worst)

**Basic approach: Try all valid modifications**

# Word Embeddings

maps each **discrete word** to a **continuous vector**



man = [1.0, 2.0, 0.0]
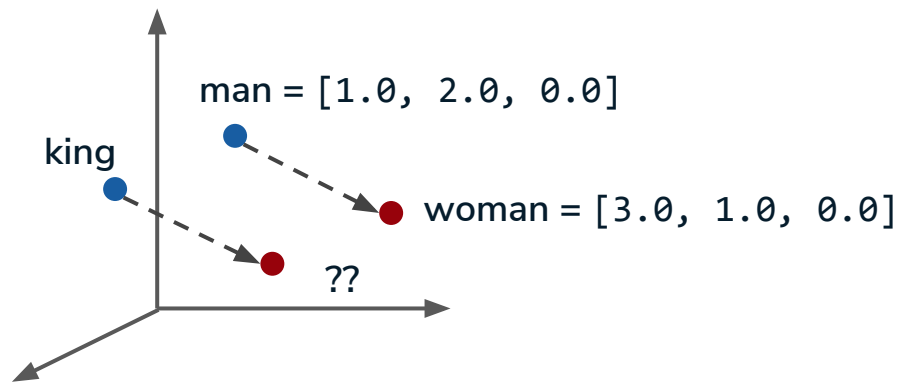
woman = [3.0, 1.0, 0.0]

**Distributed Representations of Words and Phrases and their Compositionality. Mikolov et. al. NeurIPS'13**

# Word Embeddings

maps each **discrete word** to a **continuous vector**



man = [1.0, 2.0, 0.0]

king

woman = [3.0, 1.0, 0.0]

# Word Embeddings

maps each **discrete word** to a **continuous vector**



man = [1.0, 2.0, 0.0]

king

woman = [3.0, 1.0, 0.0]

??

Distributed Representations of Words and Phrases and their Compositionality. Mikolov et. al. NeurIPS'13

# Word Embeddings

maps each **discrete word** to a **continuous vector**



man = [1.0, 2.0, 0.0]

king

woman = [3.0, 1.0, 0.0]

queen

**Distributed Representations of Words and Phrases and their Compositionality. Mikolov et. al. NeurIPS'13**

# Finding Adversarial Examples

```
...
v num = parseInt num (
    hex str .substr str (1),
    radix num
)
...
```

$\delta_1$: hex ⟶ color

```
v = parseInt(
    color.substr(1),
    radix
)
```

$\delta_2$: hex ⟶ i

```
v = parseInt(
    i.substr(1),
    radix
)
```

$\delta_3$: hex ⟶ tmp

```
v = parseInt(
    tmp.substr(1),
    radix
)
```

# Finding Adversarial Examples

renaming weight

```
...
v^num = parseInt^num(
  hex^str.substr^str(1),
  radix^num
)
...
```

$\delta_1$: hex → color

```
v = parseInt(
  ●.substr(1),
  radix
)
```

● 0.45

$\delta_2$: hex → i

```
v = parseInt(
  ●.substr(1),
  radix
)
```

◖ 0.20

$\delta_3$: hex → tmp

```
v = parseInt(
  ●.substr(1),
  radix
)
```

◖ 0.35

```
v = parseInt(
  ◖+◖+◖.substr(1),
  radix
)
```

combination of **all the renamings**

# Finding Adversarial Examples



renaming weight

$\delta_1$: hex → color

```
v = parseInt(
    ●.substr(1),
    radix
)
```

0.45

```
...
vnum = parseInt(
    hexstr.substr(
    radixnum
)
...
```

**Find weights that optimize** max $loss(\theta, \mathbf{x} + \delta, \mathbf{y})$

0.20

```
v = parseInt(
    ●+▲+●.substr(1),
    radix
)
```

**combination of all the renamings**

$\delta_3$: hex → tmp

```
v = parseInt(
    ●.substr(1),
    radix
)
```

0.35

# Solving the Inner max *loss* Efficiently

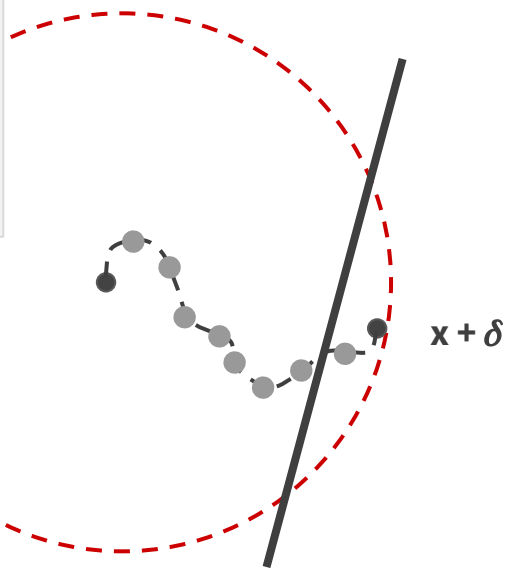**Gradient Based Optimization**

$$\theta \leftarrow \theta - \nabla_{\delta} \, loss(\theta, \mathbf{x} + \boldsymbol{\delta}, \mathbf{y})$$
$$\boldsymbol{\delta} \in \mathbf{S(x)}$$



Adversarial Examples for Models of Code.
Yefet et. al. ArXiv'20

# Robustness of Existing Models

## Type Inference [1]

```
...
v^num = parseInt^num(
  hex^str.substr^str(1),
  radix^num
)
...
```

| 89% | 45% |
|-----|-----|
| accuracy | robustness |

**Adversarial training**

| 57% |
|-----|
| robustness |

## Code Captioning [2,3]

```
int indexOfTarget(Object target) {
  int i = 0;
  for (Object elem: this.elements) {
    if (elem.equals(target))
      return i;
    i++;
  }
  return -1;
}
```

| 39.2 | 19.6 |
|------|------|
| F1 | robust F1 |

## Malware Detection [4]

```
push ebp
mov ebp, esp
push ebx
push edx                    yes
mov ebx, [ebp+4]            no
add ebx, 0x10
mov edx, [ebp+8]
mov [edx], ebx
```

| 99% | ~2% |
|-----|-----|
| accuracy | robustness |

# Robustness of Existing Models

## Type Inference [1]

```
...
v^num = parseInt^num(
    hex^str.substr^str(1),
    radix^num
)
...
```

## Code Captioning [2,3]

```
int indexOfTarget(Object target) {
    int i = 0;
    for (Object elem: this.elements) {
        if (elem.equals(target))
            return i
        i++;
    return -1;
}
```

## Malware Detection [4]

```
push ebp
mov ebp, esp
push ebx
push edx                        yes
mov ebx, [ebp+4]                no
add ebx, 0x10
mov edx, [ebp+8]
mov [edx], ebx
```

| 89% accuracy | 45% robustness |
|---|---|

| **Adversarial training** | 57% robustness |
|---|---|

| 39.2 F1 | 19.6 robust F1 |
|---|---|

| 40.6 F1 | 27.5 robust F1 |
|---|---|

| 99% accuracy | ~2% robustness |
|---|---|

# Robustness of Existing Models

## Type Inference [1]

```
...
v[num] = parseInt[num](
  hex[str].substr[str](1),
  radix[num]
)
...
```

## Code Captioning [2,3]

```
int indexOfTarget(Object target) {
  int i = 0;
  for (Object elem: this.elements) {
    if (elem.equals(target))
      return i;
    i++;
  return -1;
}
```

## Malware Detection [4]

```
push ebp
mov ebp, esp
push ebx
push edx               yes
mov ebx, [ebp+4]       no
add ebx, 0x10
mov edx, [ebp+8]
mov [edx], ebx
```

| 89% accuracy | 45% robustness |
| 39.2 F1 | 19.6 robust F1 |
| 99% accuracy | ~2% robustness |

**Adversarial training** | 57% robustness

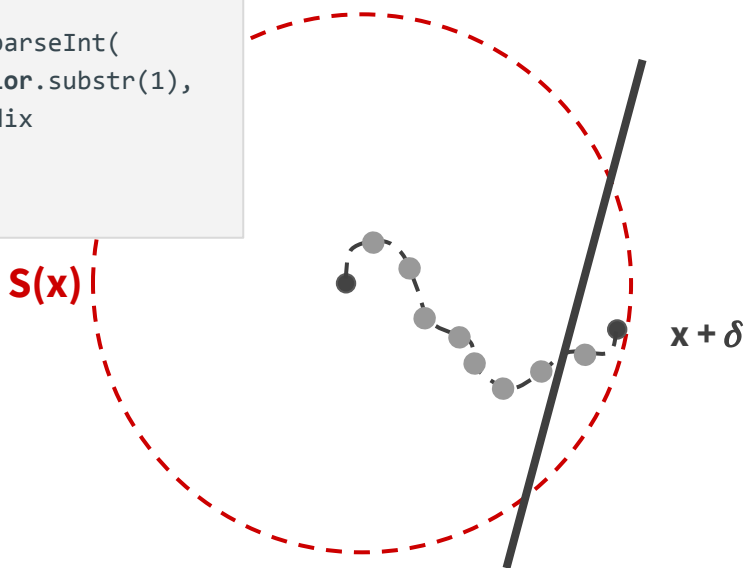| 40.6 F1 | 27.5 robust F1 |

too expensive

# Solving the Inner max loss Efficiently

## Gradient Based Optimization

$$\theta \leftarrow \theta - \nabla_\delta \, loss(\theta, \mathbf{x} + \delta, \mathbf{y})$$
$$\delta \in \mathbf{S(x)}$$

```
...
v = parseInt(
  color.substr(1),
  radix
)
...
```
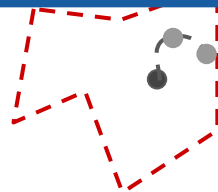
$\mathbf{S(x)}$

$\mathbf{x} + \delta$

## Refine S

$$\min \left[ \max loss(\theta, \mathbf{x} + \delta, \mathbf{y}) \right]$$
$$\delta \in \mathbf{S(\alpha(x))}$$

```
parseInt(
  _,
  _
)
```

$\mathbf{S(\alpha(x))}$

learned
representation

Adversarial Robustness for Code.
Bielik et. al. ICML'20

# Solving the Inner *max loss* Efficiently

## Gradient Based Optimization

$$\theta \leftarrow \theta - \nabla_\delta \; loss(\theta, \mathbf{x} + \delta, \mathbf{y})$$
$$\delta \in \mathbf{S(x)}$$

```
...
v = parseInt(
  color.substr(1),
  radix
)
...
```

**S(x)**

$\mathbf{x} + \delta$

## Refine **S**

$$\min \; [\max \; loss(\theta, \mathbf{x} + \delta, \mathbf{y})]$$
$$\delta \in \mathbf{S(\alpha(x))}$$

```
parseInt(
  _,
  _
)
```

**S(α(x))**

reduces the search space

leads to an easier optimization

# Solving the Inner max loss Efficiently

## Gradient Based Optimization

$$\theta \leftarrow \theta - \nabla_{\delta}\, loss(\theta, \mathbf{x} + \delta, \mathbf{y})$$
$$\delta \in \mathbf{S(x)}$$

```
...
v = parseInt(
  color.substr(1),
  radix
)
...
```

**S(x)**

$\mathbf{x} + \delta$

## Refine S

$$\min\,[\max\, loss(\theta, \mathbf{x} + \delta, \mathbf{y})]$$
$$\delta \in \mathbf{S(\alpha(x))}$$

**orthogonal to gradient optimization**

**supports all transformations**

**S(α(x))**

**reduces the search space**

**leads to an easier optimization**

# Key Techniques



Standard
Training

$\delta = \texttt{hex} \rightarrow \texttt{color}$

$\delta \in \mathbf{S(x)}$

1$^{\text{st}}$ technique

Adversarial
Training

$\alpha(\mathbf{x} + \delta)$

2$^{\text{nd}}$ technique

Representation
Learning

$y_1$    $y_2$

# Key Techniques



$\delta = \text{hex} \rightarrow \text{color}$

$\delta \in \mathbf{S(x)}$

$\alpha(\mathbf{x} + \delta)$

**1ˢᵗ technique**

**2ⁿᵈ technique**

**Standard Training**

**Adversarial Training**

**Representation Learning**

**3ʳᵈ technique**

**Abstain** = **Predict label** + **Not enough confidence for prediction**

Robust & Accurate + Robust

# Robustness of Existing Models

## Type Inference [1]

**89%** accuracy    **45%** robustness

**Adversarial training**    **57%** robustness

**+ Refinement & Abstain**    **67%** robustness

## Code Captioning [2,3]

**39.2** F1    **19.6** robust F1

**27.5** robust F1

## Malware Detection [4]

**99%** accuracy    **~2%** robustness

too expensive

# Robust Models for Source Code

# Robust Models for Source Code