

ICML 2020

Adversarial Robustness for Code

Pavol Bielik, Martin Vechev

pavol.bielik@inf.ethz.ch, martin.vechev@inf.ethz.ch

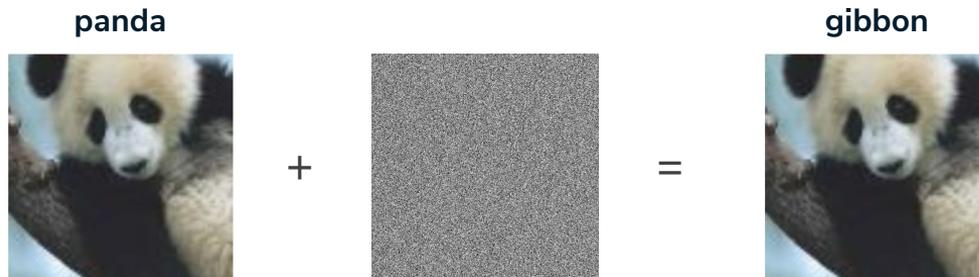
Department of Computer Science

ETH zürich

SRILAB

Adversarial Robustness

 Vision



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

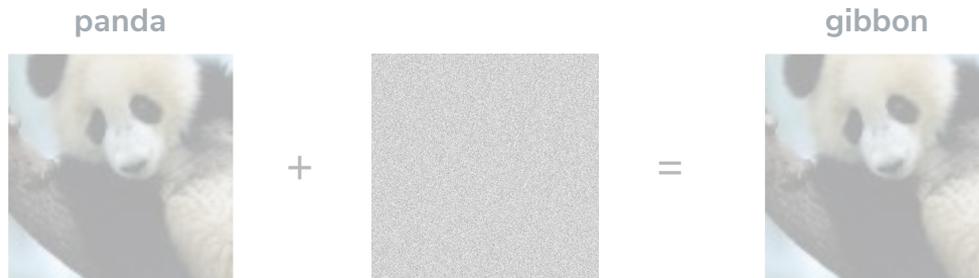
 Sound



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

Adversarial Robustness for Code

 Vision



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

 Sound

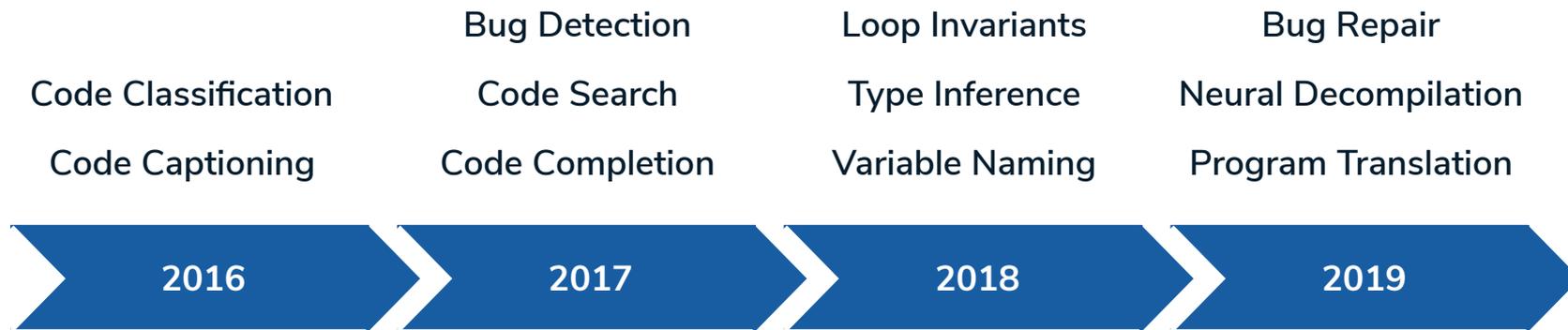


Audio Adversarial Examples: Targeted Attacks on Speech-to-Text. Carlini et. al. ICML'18 workshop

 Code



Deep Learning + Code



Prior Works

90%

Accuracy

Adversarial Robustness for Code



Prior Works

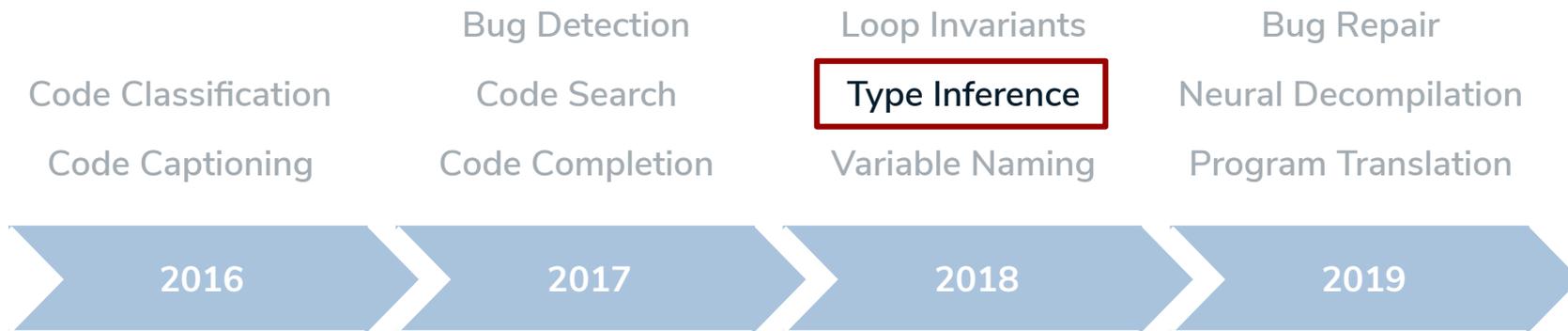
90%

Accuracy

?

Robustness

Adversarial Robustness for Code



Prior Works

90%

Accuracy

4%-50%

Robustness

This Work

88%

Accuracy

84%

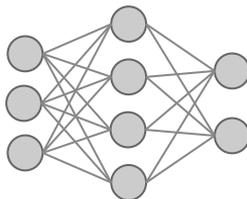
Robustness

Adversarial Robustness Example

Input Program
 x

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

Model
 $f(x) \rightarrow y$



(Type Inference)
Program Properties
 y

```
...  
vnum = parseIntnum(  
  hexstr.substrstr(1),  
  radixnum  
)  
...
```

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

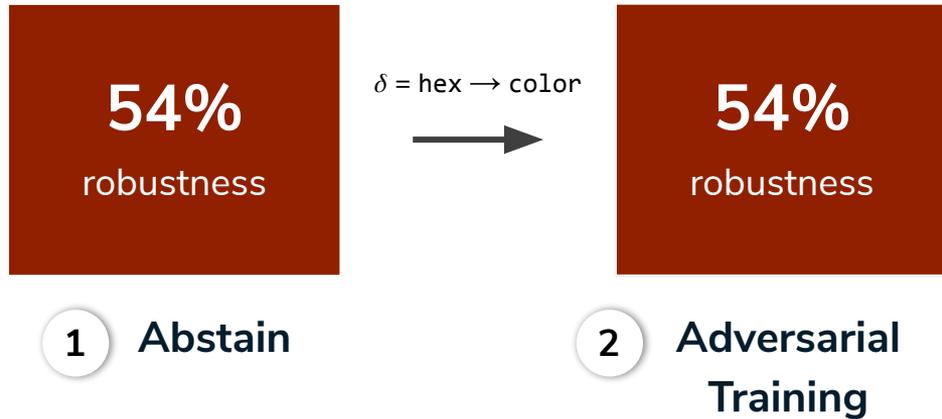
Our Work: Three Key Techniques

```
...  
v = parseInt(  
  hexabs.substrabs(1),  
  radixabs  
)  
...
```

**Allows model
not to make
a prediction
if uncertain**

1 Abstain

Our Work: Three Key Techniques



Our Work: Three Key Techniques

```
...  
v = parseInt(  
  hexabs.substrabs(1),  
  radixabs  
)  
...
```

1 Abstain

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



```
...  
vnum = parseIntnum(  
  color.substr(1),  
  radix  
)  
...
```

2 Adversarial
Training

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



```
parseIntnum(  
  -,  
  -  
)
```

3 Representation
Learning

Our Work: Three Key Techniques

```
...  
v = parseInt(  
  hexabs.substrabs(1),  
  radixabs  
)  
...
```

1 Abstain

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



```
...  
vnum = parseIntnum(  
  color.substr(1),  
  radix  
)  
...
```

2 Adversarial
Training

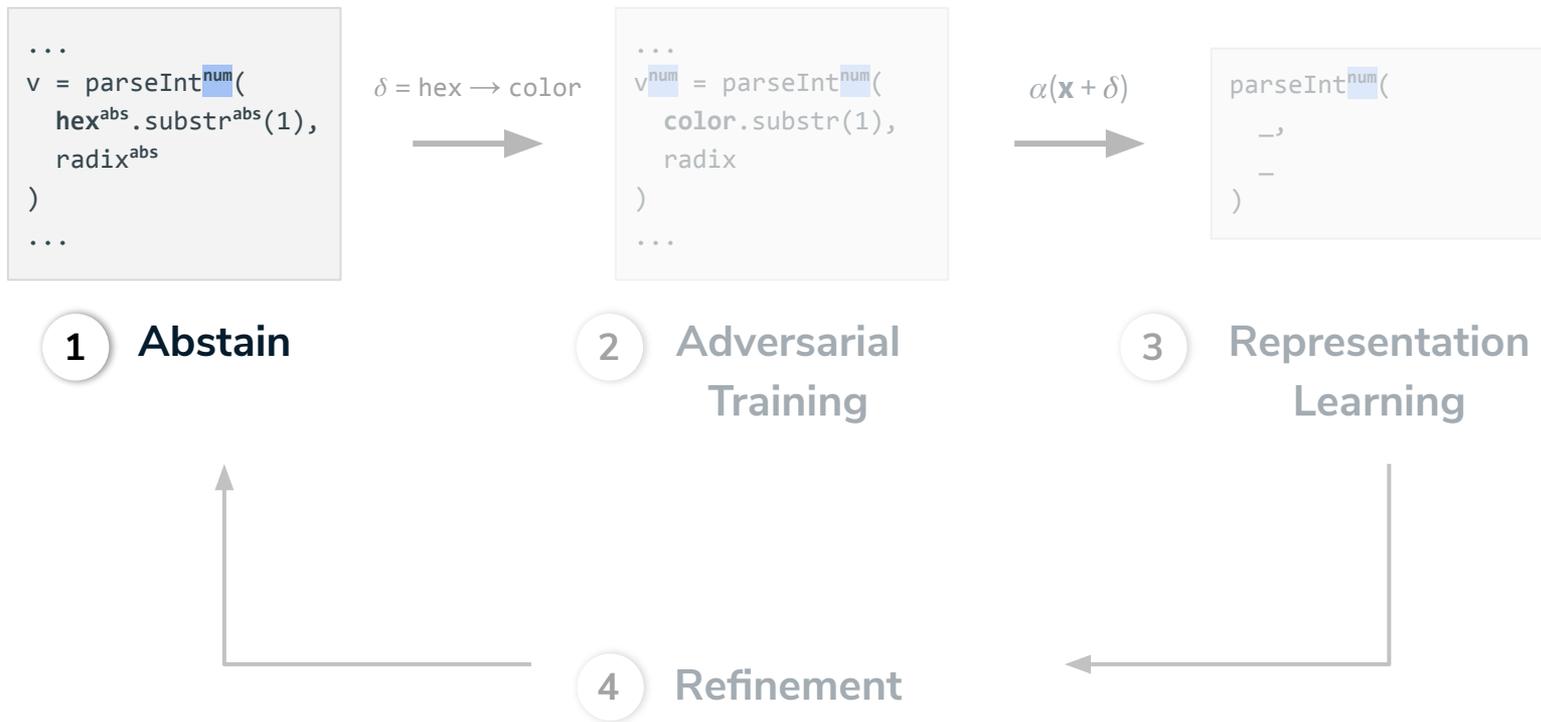
$\alpha(\mathbf{x} + \delta)$



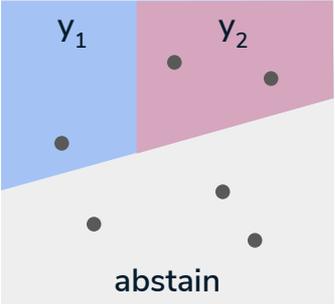
3 Representation
Learning

84%
robustness

Our Work: Three Key Techniques

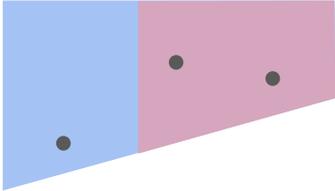


Learning to Abstain



● input x_i

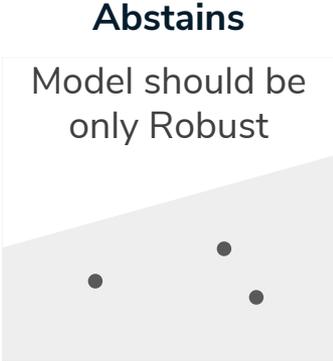
=



Model should be both Robust and Accurate

Predict Class

+



Leads to a simpler optimization problem

Property prediction problem is undecidable

Learning to Abstain

Main Insight

Combine Robustness + Learning to Abstain

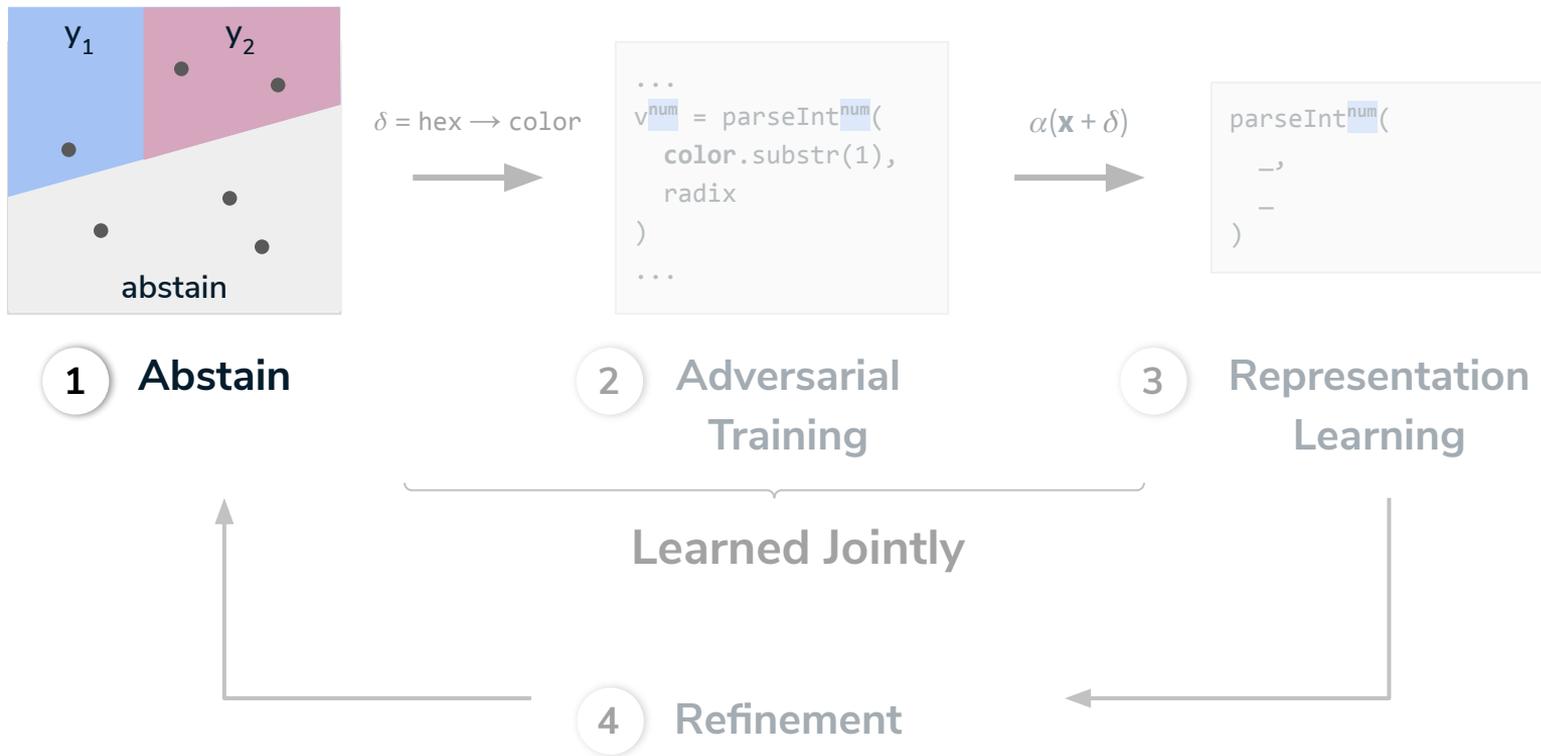
How to Abstain?

**Deep Gamblers: Learning to Abstain with Portfolio Theory.
Liu et. al. NeurIPS'19**

**Leads to a simpler
optimization problem**

**Property prediction
problem is undecidable**

Our Work: Three Key Techniques



Adversarial Training

measures the model performance ground-truth label

Standard training

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


Adversarial training

$$\min_{\theta} [\max_{\delta \in \mathbf{S}(\mathbf{x})} \text{loss}(\theta, \mathbf{x} + \delta, \mathbf{y})]$$

Label preserving program transformations



1

Define the space **S** of
program transformations

2

Solve the inner
max *loss* efficiently

Label Preserving Program Transformations

Word Substitution

tensors + δ
very fast

Constants, Binary Operators, ...

```
7  
radix + offset
```

$x + \delta$
→

```
42  
radix - offset
```

Word Renaming

tensors + δ + analysis
fast

Rename Variables, Parameters, Fields, Method Names, ...

```
def getID() {...}  
client.Name
```

$x + \delta$
→

```
def get_id() {...}  
client.name
```

Sequence Substitution

tensors → code + δ + analysis → tensors
slow

Adding Dead Code, Reordering Statements, ...

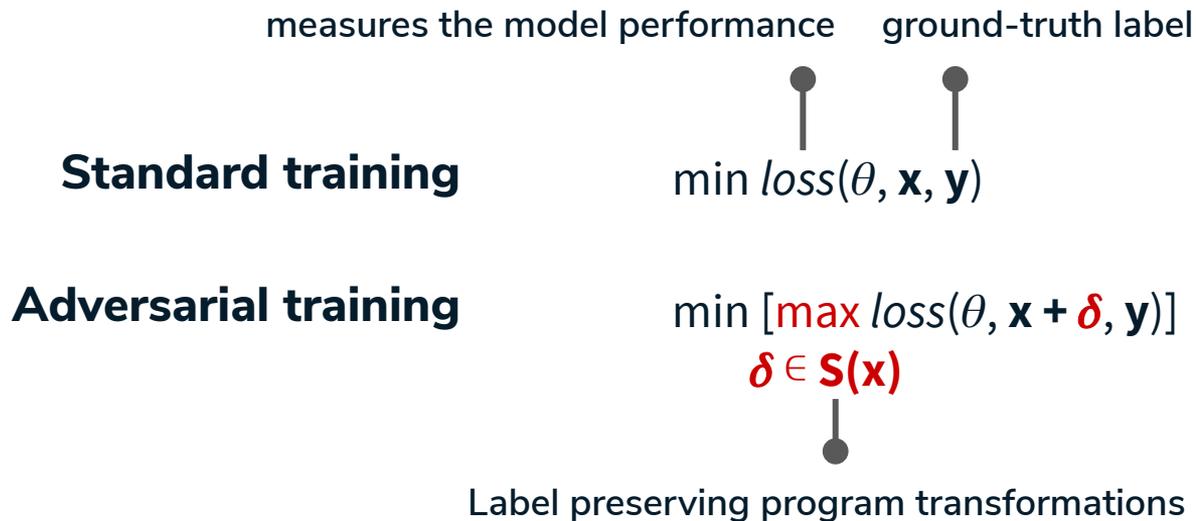
```
a = get_id()  
b = 42
```

$x + \delta$
→

```
b = 42  
a = get_id()
```



Adversarial Training



1 Define the space **S** of program transformations

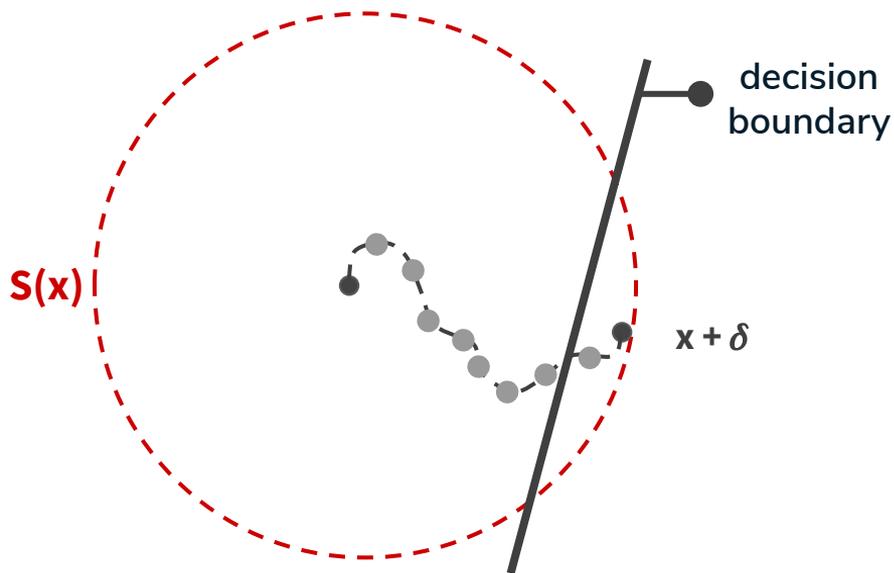
2 Solve the inner **max** loss efficiently

Solving the Inner max loss Efficiently

Gradient Based Optimization

$$\theta \leftarrow \theta - \nabla \text{loss}(\theta, \mathbf{x} + \boldsymbol{\delta}, \mathbf{y})$$

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



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

Limitations

54% → **54%**
standard adversarial

same or worse robustness

Discrete and
disruptive changes

Highly structured
and large programs

hard optimization problem

no structural transformations

Solving the Inner max loss Efficiently

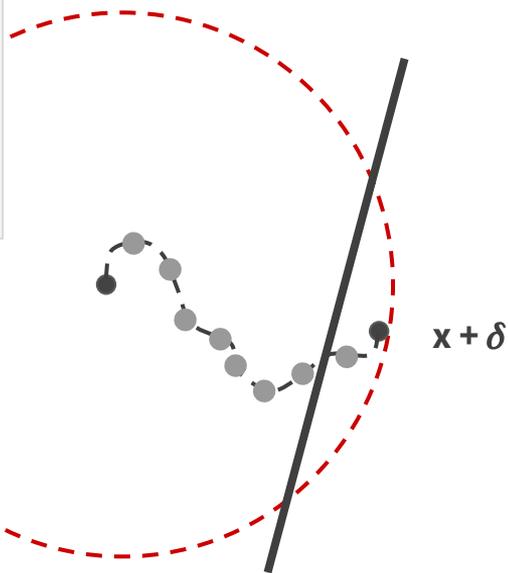
Gradient Based Optimization

$$\theta \leftarrow \theta - \nabla \text{loss}(\theta, \mathbf{x} + \boldsymbol{\delta}, \mathbf{y})$$

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

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

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



Refine \mathbf{S}

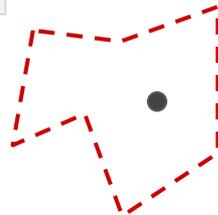
$$\min [\max \text{loss}(\theta, \mathbf{x} + \boldsymbol{\delta}, \mathbf{y})]$$

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

```
parseInt(  
  -  
  -  
)
```

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

learned
representation



Solving the Inner max loss Efficiently

Gradient Based Optimization

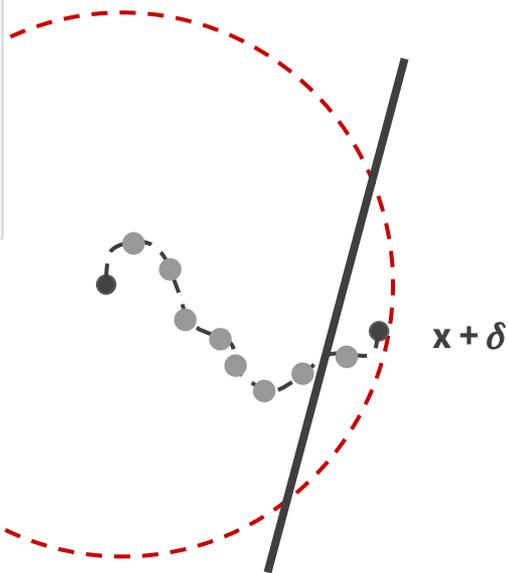
$$\theta \leftarrow \theta - \nabla \text{loss}(\theta, \mathbf{x} + \boldsymbol{\delta}, \mathbf{y})$$

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

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

```

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



Refine \mathbf{S}

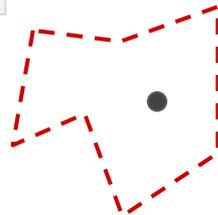
$$\min [\max \text{loss}(\theta, \mathbf{x} + \boldsymbol{\delta}, \mathbf{y})]$$

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

```
parseInt(  
  -  
  -  
)  

```

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



reduces the search space

leads to an easier optimization

Solving the Inner max loss Efficiently

Gradient Based Optimization

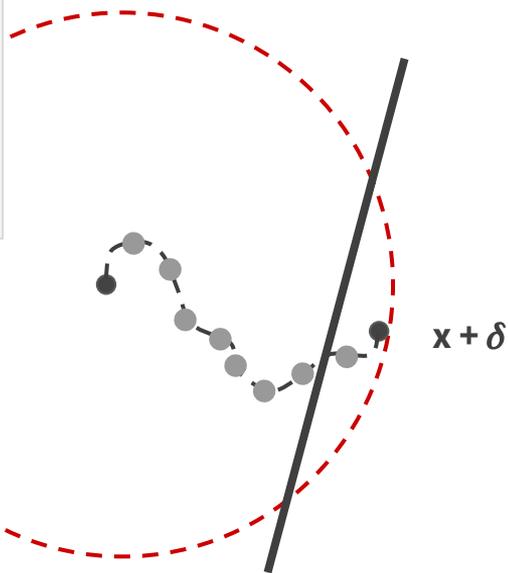
$$\theta \leftarrow \theta - \nabla \text{loss}(\theta, \mathbf{x} + \boldsymbol{\delta}, \mathbf{y})$$

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

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

```

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



Refine \mathbf{S}

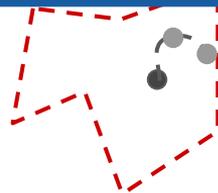
$$\min [\max \text{loss}(\theta, \mathbf{x} + \boldsymbol{\delta}, \mathbf{y})]$$

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

orthogonal to gradient optimization

supports all transformations

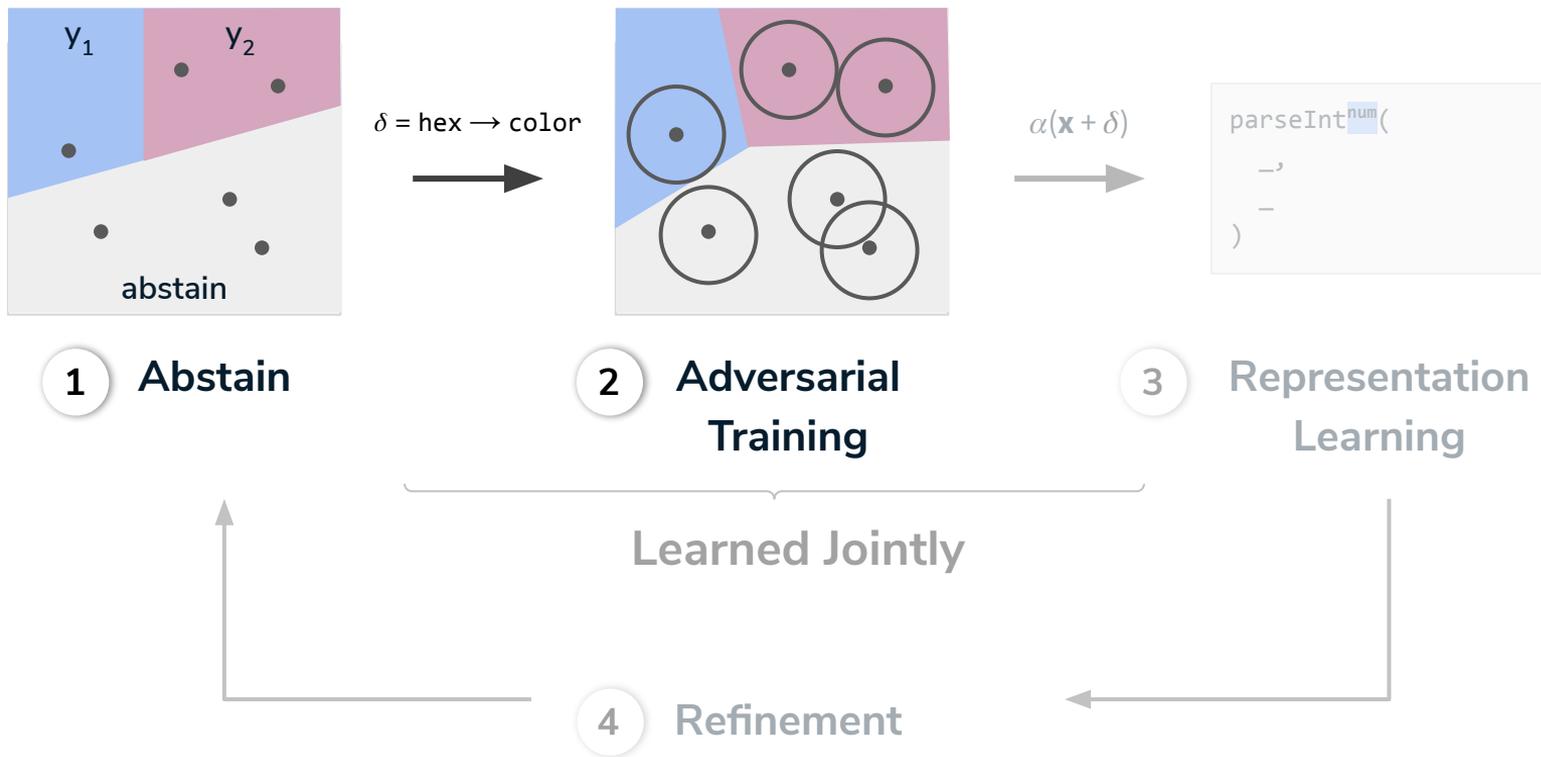
$\mathbf{S}(\alpha(\mathbf{x}))$



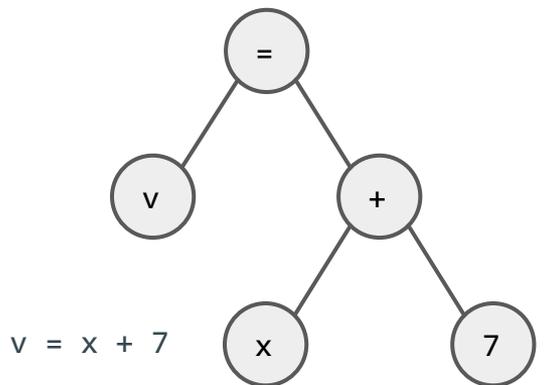
reduces the search space

leads to an easier optimization

Our Work: Three Key Techniques



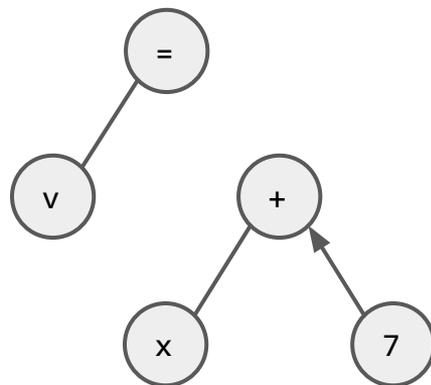
Representation Learning



nodes attributes

$$G = \langle V, E, \xi \rangle$$

edges



$$\langle V, E, \xi \rangle \rightarrow \langle V, E' \subseteq E, \xi \rangle$$

1 Programs as Graphs

Learning to Represent Programs with Graphs.

Allamanis et. al. ICLR'18

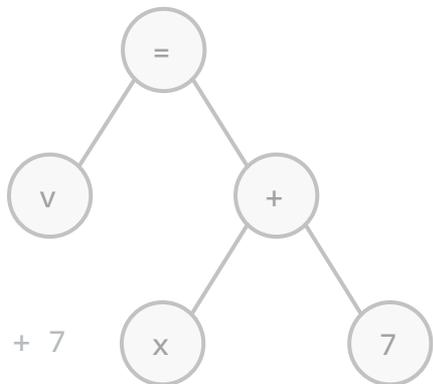
Generative Code Modeling with Graphs.

Brockschmidt et. al. ICLR'19

2 Define Refinement

Remove Graph Edges

Representation Learning



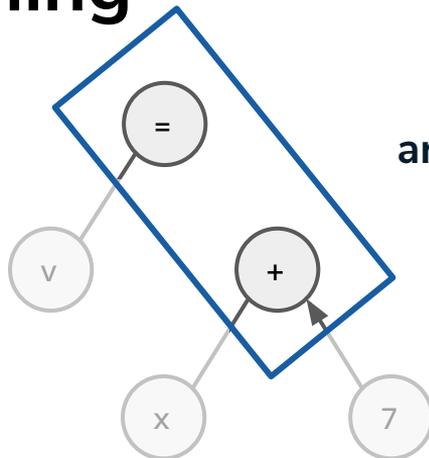
$v = x + 7$

nodes

attributes

$$G = \langle V, E, \xi \rangle$$

edges



All decisions
are made locally

$\alpha:$

$$\langle V, E, \xi \rangle \rightarrow \langle V, E' \subseteq E, \xi \rangle$$

1 Programs as Graphs

Learning to Represent Programs with Graphs.

Allamanis et. al. ICLR'18

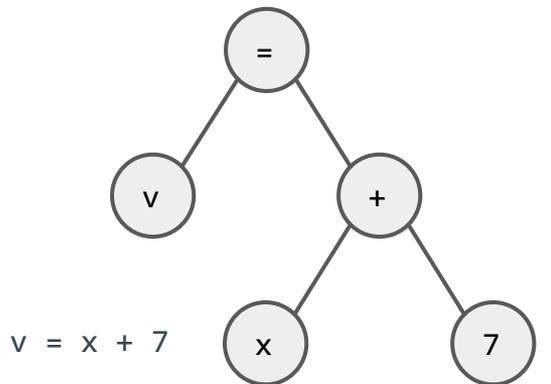
Generative Code Modeling with Graphs.

Brockschmidt et. al. ICLR'19

2 Define Refinement

Remove Graph Edges

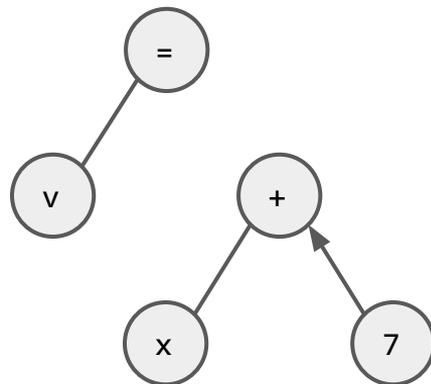
Representation Learning



nodes attributes

$$G = \langle V, E, \xi \rangle$$

edges



$$\langle V, E, \xi \rangle \rightarrow \langle V, E' \subseteq E, \xi \rangle$$

$$\arg \min_{\alpha} \sum_{(\mathbf{x}, \mathbf{y}) \in \mathcal{D}} |\alpha(\mathbf{x})|$$

subject to

$$\text{loss}(\theta, \mathbf{x}, \mathbf{y}) \approx \text{loss}(\theta, \alpha(\mathbf{x}), \mathbf{y})$$

1 Programs as Graphs

Learning to Represent Programs with Graphs.

Allamanis et. al. ICLR'18

Generative Code Modeling with Graphs.

Brockschmidt et. al. ICLR'19

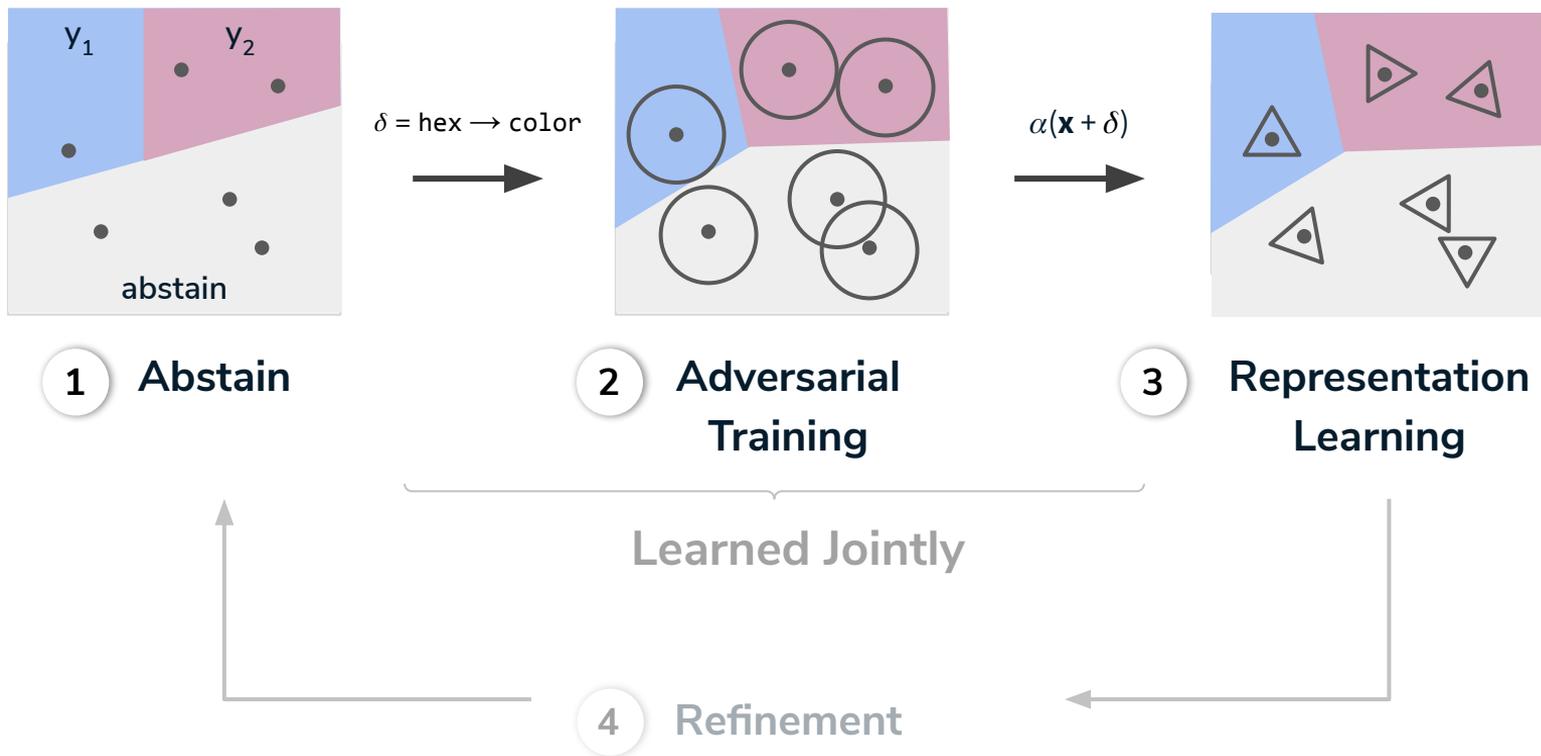
2 Define Refinement

Remove Graph Edges

3 Optimize α

Minimize Graph Size

Our Work: Three Key Techniques



Evaluation

Type Inference

Task

```
vnum = parseIntnum(  
  hexstr.substrstr(1),  
  radixnum  
)
```

target classes (y)

string, number,
boolean, void
(\Rightarrow)string, (\Rightarrow)number,
(\Rightarrow)boolean, (\Rightarrow)void
any



TypeScript



JavaScript

Models

LSTM

DeepTyper

Graph Neural Networks

LSTM + 1 layer GNN + LSTM

GNN_{Transformer}

GNN_{GCN}

GNN_{GGNN}

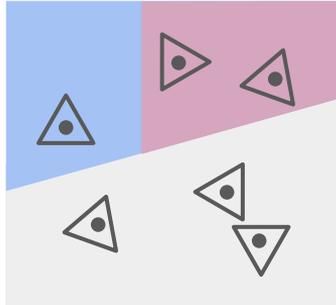
DeepTyper: Deep Learning Type Inference.
Hellendoorn et. al., FSE'18

more complex
type inference

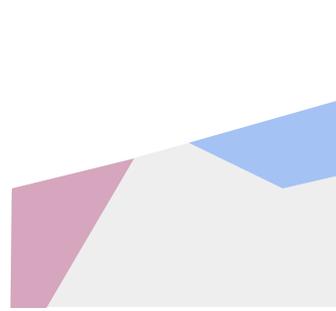
Typilus: Neural Type Hints. Allamanis et. al. PLDI'20

LambdaNet: Probabilistic Type Inference using Graph Neural Networks. Wei et. al. ICLR'20

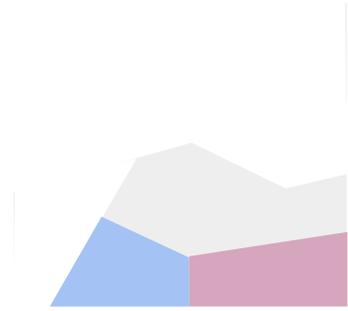
Our Work: Three Key Techniques



1st Model



2nd Model



3rd Model



4 Refinement



Evaluation

GNN_{Transformer}

	Accuracy	Robustness
Standard Training	89.3%	54.9%
Adversarial Training	90.3%	54.3%
All Components	88.4%	83.8%

-1%

Accuracy

+29%

Robustness

Evaluation

GNN_{Transformer}

Accuracy

Robustness

Abstain

Standard Training

89.3%

54.9%

-

Adversarial Training

90.3%

54.3%

-

0%

All Components

88.4%

83.8%

-

99%

All Components

99.0%

99.6%

61.3%

100%

All Components

99.9%

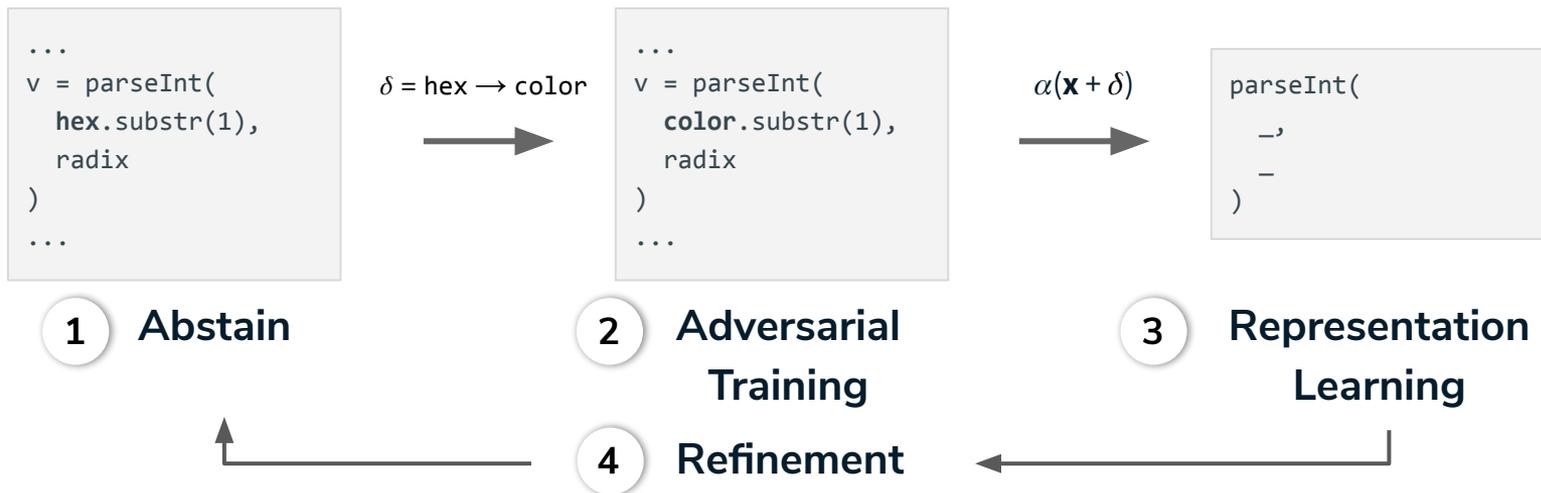
99.9%

75.9%

Target
Accuracy

Allows training highly accurate & robust models

Adversarial Robustness for Code



For more experiments and results, please refer to the extended version of our paper

We only scratched the surface, more work in domain of code is needed and is being done, e.g.:

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

Optimization-guided binary diversification to mislead neural networks for malware detection. Sharif et. al. ArXiv

Semantic Robustness of Models of Source Code. Ramakrishnan et. al., ArXiv