

# Analyzing and Securing Software via Robust and Generalizable Learning

**Kexin Pei**

Department of Computer Science



COLUMBIA UNIVERSITY  
IN THE CITY OF NEW YORK





**“Software is Eating the World”**

**- Marc Andreessen**

# Software is Plagued with Errors

“Bad software cost US businesses **\$2.41 trillion** in 2022” - **SC Media**

“**280 days** average time companies need to detect and respond to cyber attacks...” - **Skybox**

“Cybercrime is predicted to cost the world **\$7 trillion** in 2022” - **CISQ Report**



## Hackers breach energy orgs via bugs in discontinued web server

By [Sergiu Gatlan](#)

November 22, 2022 02:55 PM 0



## Florida Hack Exposes Danger to Water Systems

STATELINE ARTICLE | March 10, 2021 | By: [Jenni Bergal](#) | Read time: 7 min

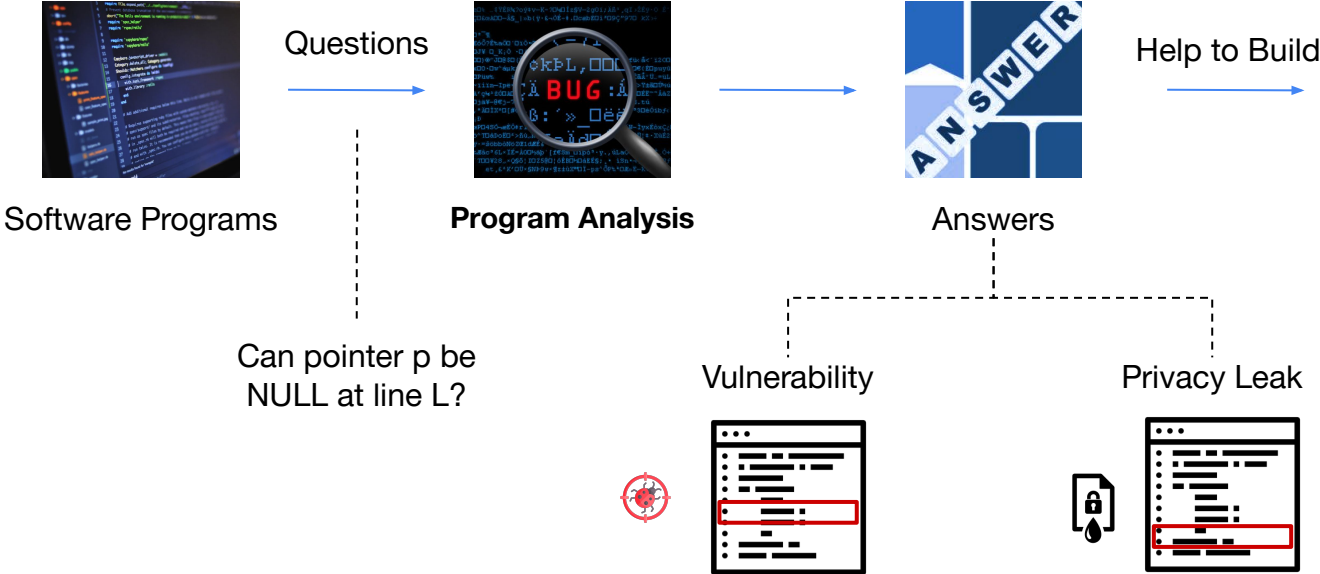


CYBERSECURITY

## Cyberattack on food supply followed years of warnings

Virtually no mandatory cybersecurity rules govern the millions of food and agriculture businesses that account for about a fifth of the U.S. economy. And now, the risk has become real.

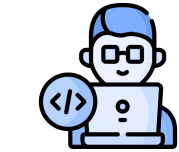
# Program Analysis is Crucial for Building Trustworthy Software



## Trustworthy Software

- Security
- Reliability
- Safety
- Privacy
- Performance

# Challenges of Traditional Program Analysis



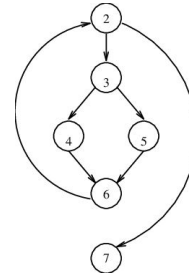
Human Expert

Hand-Curate



Rules and Heuristics

How to



Represent a program?

**Input:**  $kill(B)$  and  $gen(B)$  for every basic block  $B$ .

**Output:**  $in(B)$  and  $out(B)$  for every basic block  $B$ .

**for each**  $B$  **repeat**

$out(B) := gen(B)$

**while** changes to any  $out(B)$  occur **repeat**

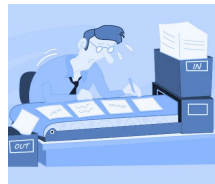
$in(B) := \bigcup_{B' \in pred(B)} out(B')$

$out(B) := gen(B) \cup (in(B) \setminus kill(B))$

Design analysis rules?

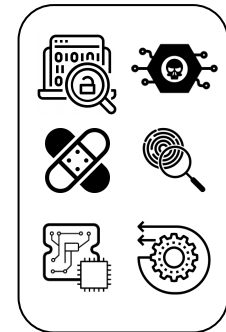
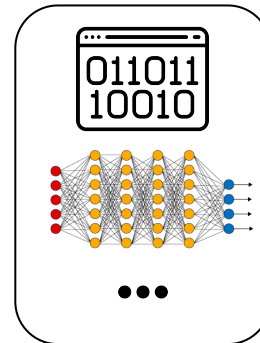
Heterogeneous Software

Various Security Applications

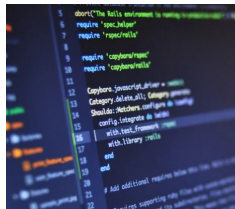
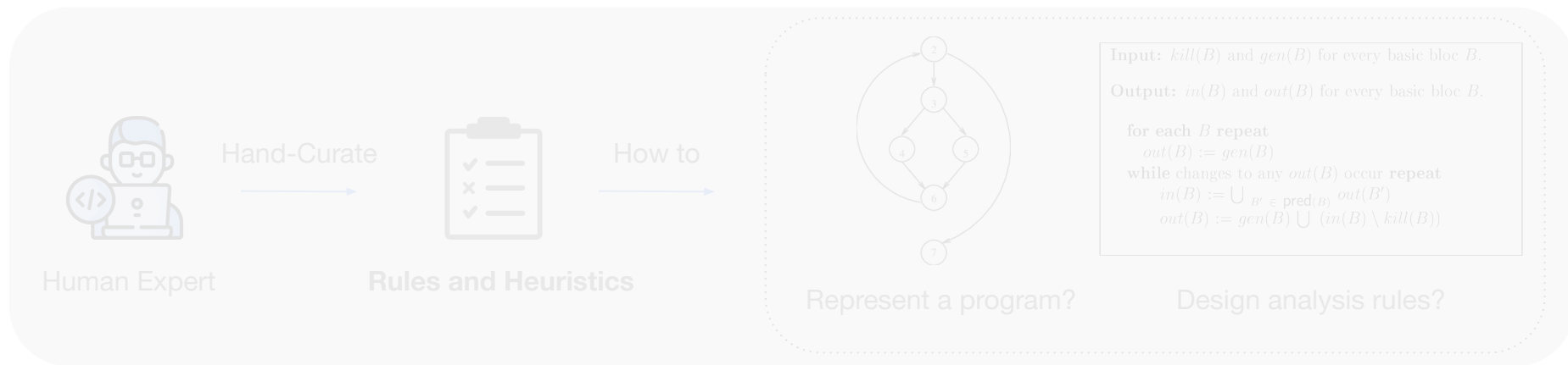


Significant Manual Effort

Tune



# Machine Learning Shows Promise for Analyzing Programs



Abundant Public Code



Abundant Compute

Learn



Program Representations

Analysis Rules

# Machine Learning Shows Promise for Analyzing Programs



## Detecting and Exploiting Vulnerabilities



devtooligan ✓  
@devtooligan



Brandon Rossi  
@0xConda

omg. seriously mind-blown 🤯

I told ChatGPT to pentest a scope of IPs and it found me \$69,420 worth of bugs.

Are there any vulnerabilities in this Solidity contract? I want to make sure no one can steal other peoples' deposits. If there are, how do I fix it?

```
contract Vault {
    mapping(address => uint256) public balances;

    function deposit() public payable {
        balances[msg.sender] += msg.value;
    }

    function withdraw(uint256 amount, address from) public {
        balances[from] -= amount;
        payable(from).transfer(amount);
    }
}
```

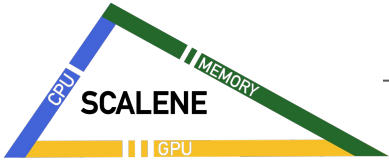


This contract appears to have a potential vulnerability in the `withdraw` function. The `withdraw` function allows anyone to specify the `from` address when calling the function, which means that a malicious user could potentially withdraw funds from an arbitrary address by calling `withdraw` and specifying a different `from` address.

To fix this issue, you could modify the `withdraw` function to only allow the caller to withdraw funds from their own account. You could do this by replacing the `from` parameter with `msg.sender`, which will always represent the caller's address. Here is an example of how you could do this:

# Machine Learning Shows Promise for Analyzing Programs

Python Profiler



Program Optimization

```
⚡ z1 = [i for i in range(0,300000)][299999]
```

---

```
15 # Proposed optimization:  
# This code can be optimized by using the built-in function max()  
z1 = max(range(0, 300000)) # ~10x faster
```



Explain Code

Translate Code

A screenshot of the GitHub Copilot 'EXPLAIN' interface. It shows a code snippet for a FizzBuzz program. Below the code is a dropdown menu with 'Explain code' selected. Underneath is an 'Advanced' section with an 'Ask Copilot' button. The 'RESULT' section contains a list of four numbered steps explaining the code's logic.

A screenshot of the GitHub Copilot 'LANGUAGE TRANSLATION' interface. It shows the same FizzBuzz code snippet. Below the code is a dropdown menu with 'python' selected. At the bottom is an 'Ask Copilot' button.



3% Code Written by ML

ML-Powered Program  
Fixing, Repair,  
Refactoring, etc.

Huge Academic Contributions:  
500+ Papers

<https://ml4code.github.io/>



# Limitations: Lack Understanding of Program Semantics

A code summarization example (Alon et al., 2019, Yefet et al., 2020, Henkel et al. 2022)

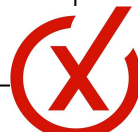
[code2vec.org](http://code2vec.org) / [code2seq.org](http://code2seq.org)

```
void f1(int[] array){
  boolean swapped = true;
  for (int i = 0;
      i < array.length && swapped; i++){
    swapped = false;
    for (int j = 0;
        j < array.length-1-i; j++) {
      if (array[j] > array[j+1]) {
        int temp = array[j];
        array[j] = array[j+1];
        array[j+1]= temp;
        swapped = true;
      }
    }
  }
}
```



Prediction: **sort** (98.54%)

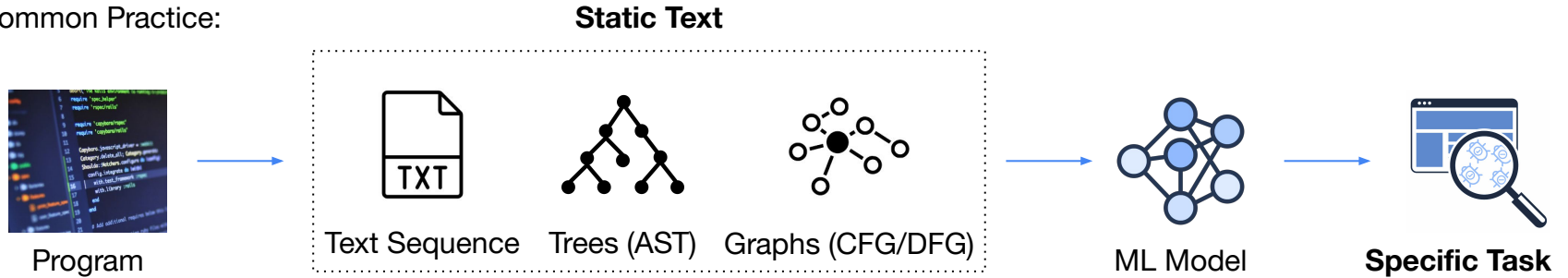
```
void f2(int[] ttypes){
  boolean swapped = true;
  for (int i = 0;
      i < ttypes.length && swapped; i++){
    swapped = false;
    for (int j = 0;
        j < ttypes.length-1-i; j++) {
      if (ttypes[j] > ttypes[j+1]) {
        int temp = ttypes[j];
        ttypes[j] = ttypes[j+1];
        ttypes[j+1]= temp;
        swapped = true;
      }
    }
  }
}
```



Prediction: **contains** (99.97%)

# Common Practice of ML on Code

Common Practice:



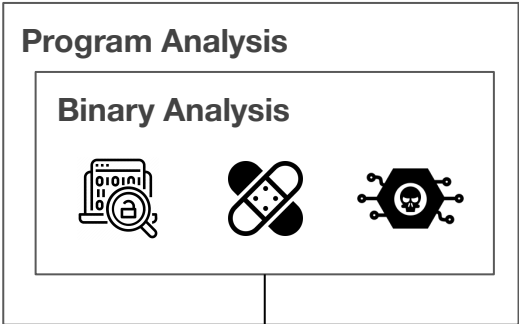
**Program semantics** does not just manifest in **static text**

**Consequences: Lacking Robustness and Generalization**

- I. **Overfit** to **spurious textual** and **task-specific** patterns
- II. **Distribution shift**: **program syntax** and **task requirement** changes

**Security Applications Require More Rigorous Understanding of Program Semantics**

# A Popular Type of Program Analysis in Security: Binary Analysis



**Stripped Binaries**  
Proprietary software,  
Third-party component  
Malware

## Additional Complications:

```
int variable_init()
{
    int var = 22;
    return var;
}
```

```
mov eax, 0x16
```






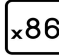

```
xor eax, eax
sub eax, -0x16
```

- Variables
- Arguments
- Types
- Data Structures
- Names
- ...

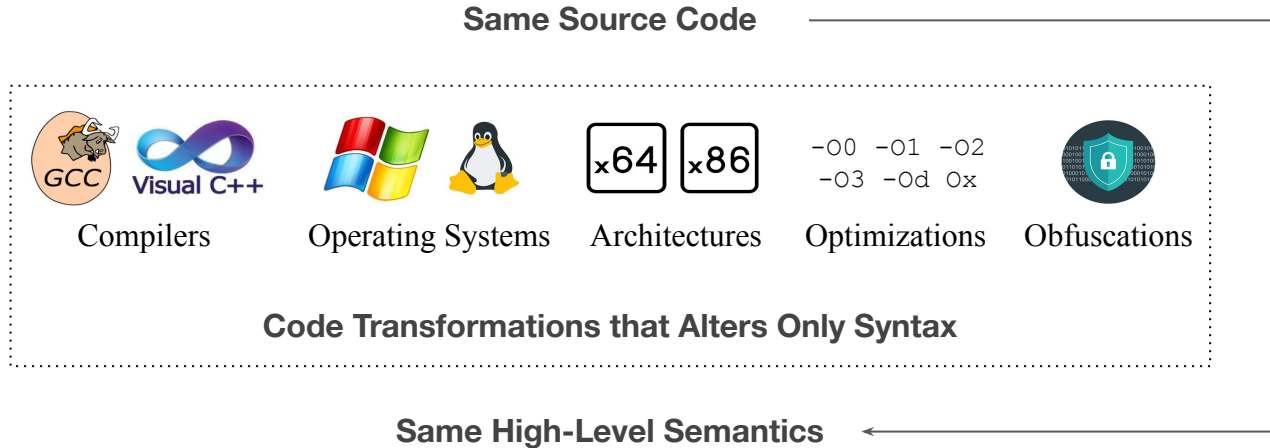
**Semantic Hints are Absent**

- Registers
- Memory Accesses

## Various Code Transformations

			
Compilers		Operating Systems	
		-O0 -O1 -O2 -O3 -Od Ox	
Architectures		Optimizations	Obfuscations

# What do Robustness and Generalization Imply in Binary Program Analysis?

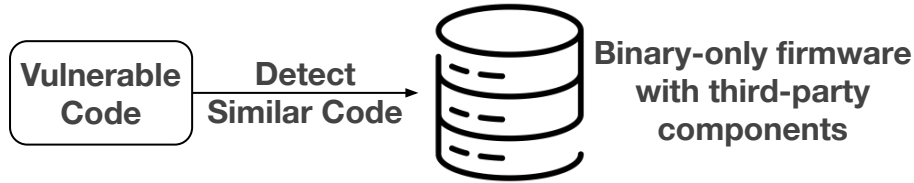


**Robustness:** Stay Invariant to **Syntactic** Changes

**Generalization:** Generalize to new **Syntactic** Changes

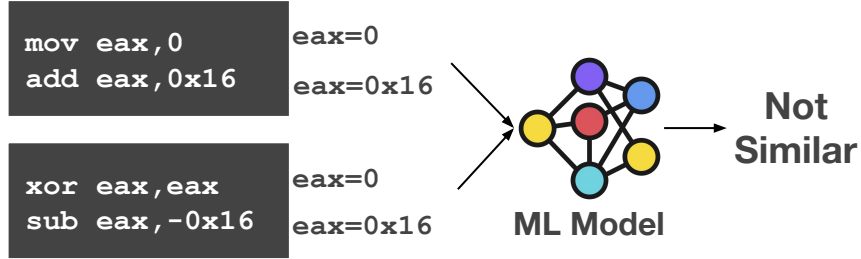
# Security Applications Require Rigorous Understanding of Program Semantics

## Detecting Binary Code Reuse Vulnerability




Without understanding **mov**, **xor**, **add**, **sub**, etc.

ML model cannot reason about program behavior to predict similarity



KP `mov eax, 0`  
`add eax, 0x16`

 This is assembly code...these instructions initialize the **EAX** register to 22.

KP `xor eax, eax`  
`sub eax, -0x16`

 No...the second set uses the "xor" and "sub" instructions to **set the value of EAX register to -22...**

# My Research

## Data-Driven Program Analysis

**Semantic Similarity**  
[TSE'22]

**Specification Inference**  
[ICML'23]

**Debug Symbol Recovery**  
[FSE'21, CCS'22]





**Memory Dependence**  
[FSE'22]

### Program Behavior


**Efficient**  
98.1X

**Precise**  
118%+

**Generalizable and Robust Across**

    -O0 -O1 -O2  
-O3 -O4 -Ox

Compilers Architectures Optimizations

 Obfuscations



**Type Inference**  
[FSE'16]

**Fuzzing via Program Smoothing**  
[Oakland S&P'19]

**Malware Analysis**  
[DSN'15]

**Attack Forensics**  
[ACSAC'16]

**Disassembly**  
[NDSS'21]

**SSL/TLS Hostname Verification**  
[Oakland S&P'17]

### Program Structure

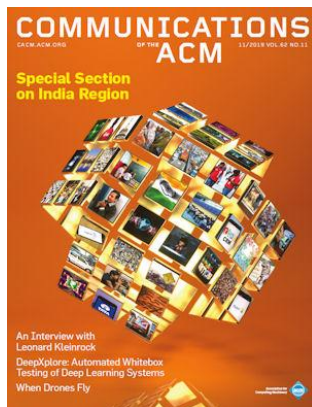
## Systematic Whitebox Testing of Neural Networks

[SOSP'17, ICSE'18]

SOSP Best Paper Award

Inspired over Thousands of Follow-Up Projects

Data-Driven  
Program Analysis



brain-research/  
tensorfuzz

A library for performing coverage guided fuzzing of neural networks

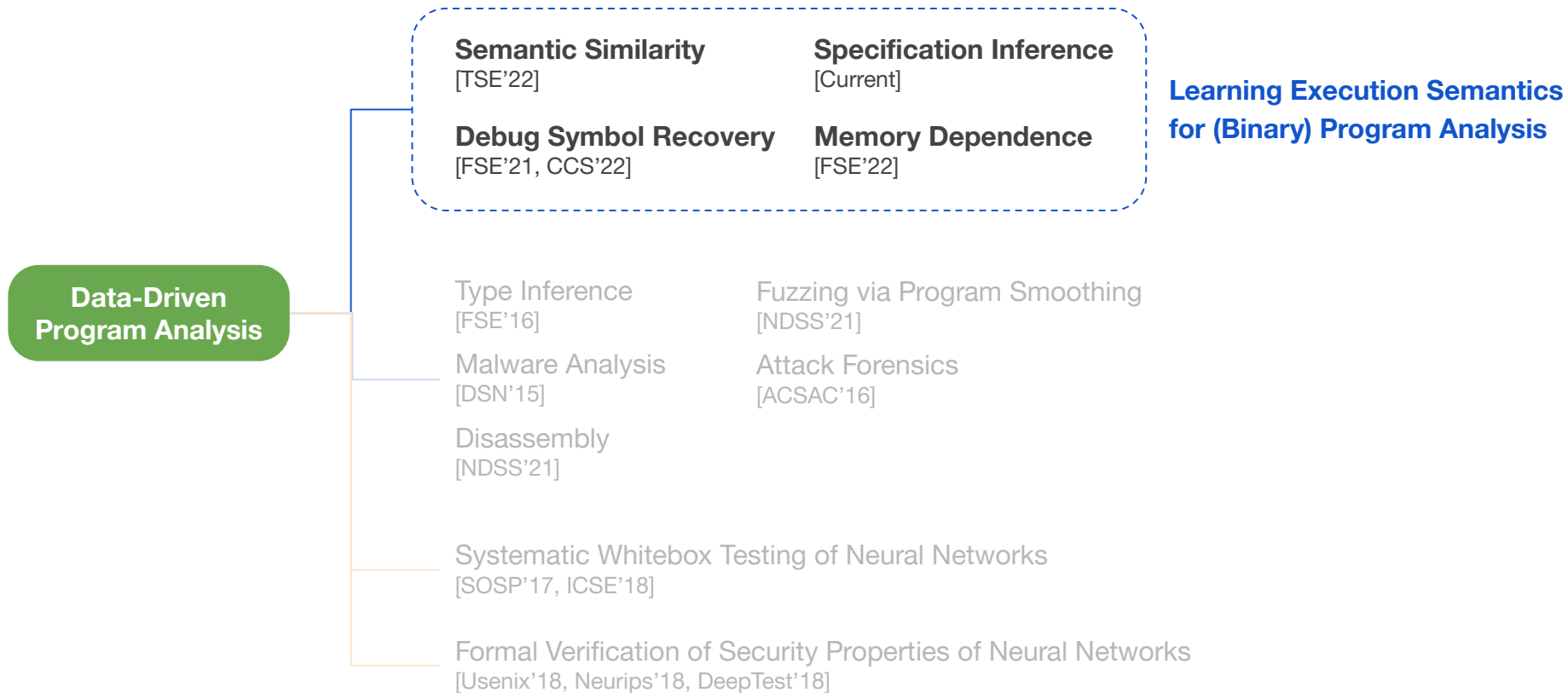
0 Contributors 7 Issues 195 Stars 56 Forks



## Formal Verification of Security Properties of Neural Networks

[Usenix'18, Neurips'18, DeepTest'18]

# My Research

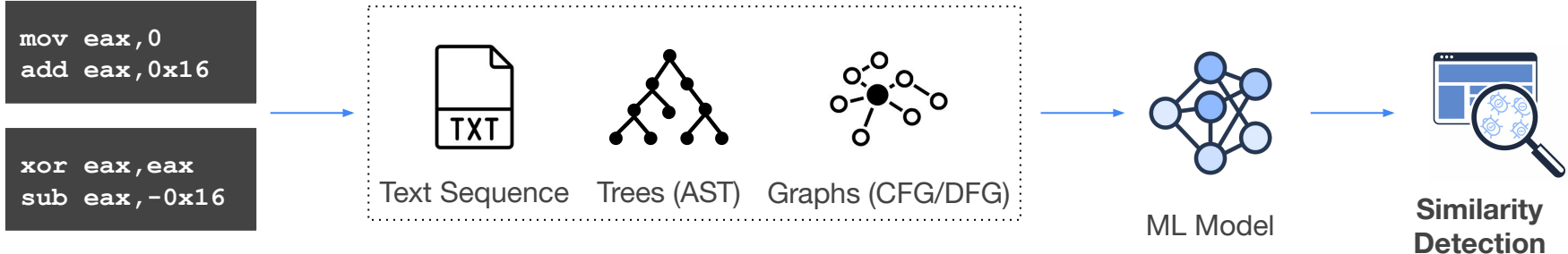




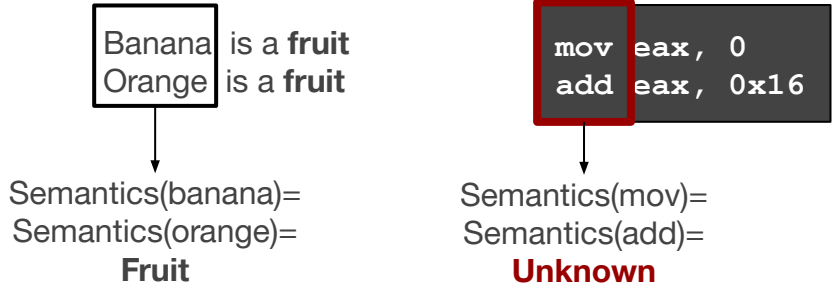
# Learning Execution Semantics for Binary Program Analysis

# Security Applications Require Rigorous Understanding of Program Semantics

Common Practice:

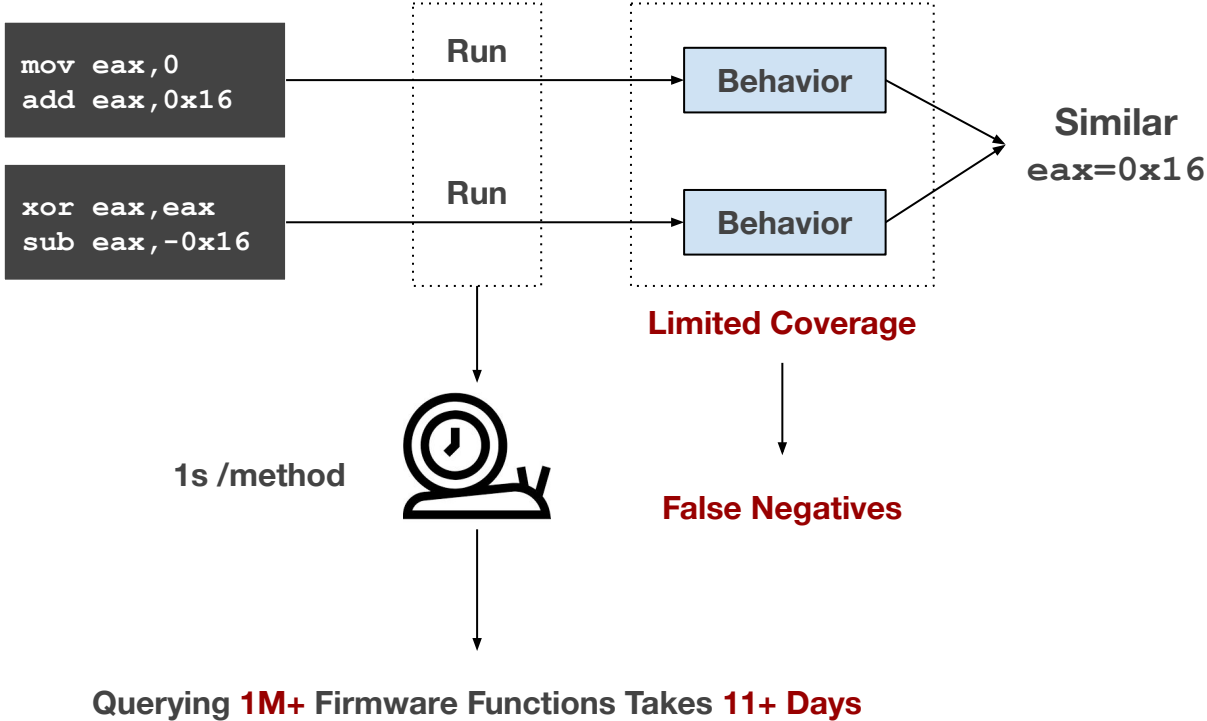


**Program semantics** does not manifest in **static text**

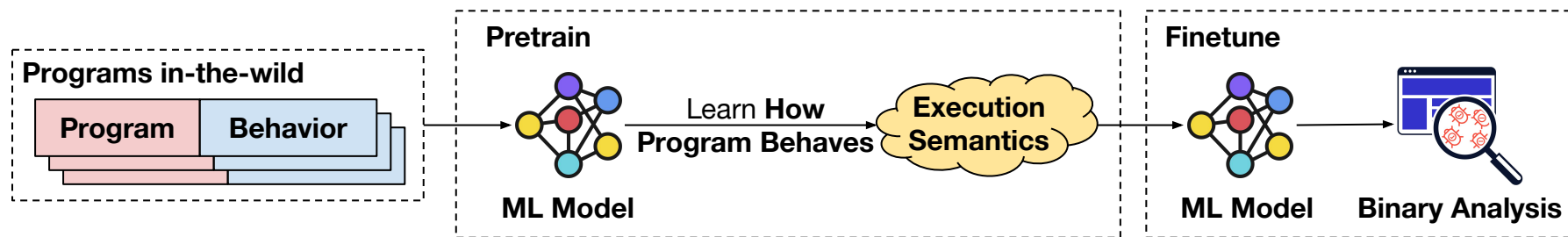
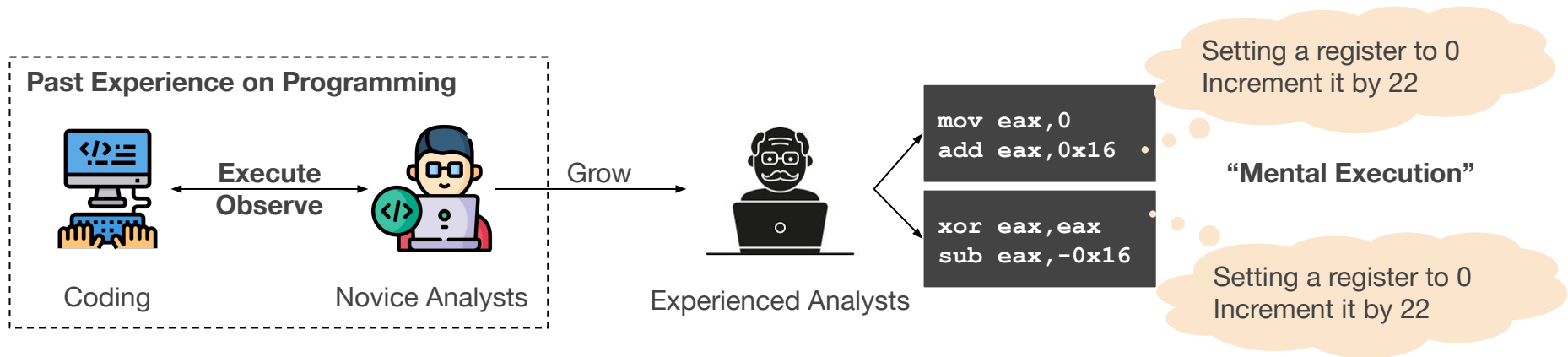


Lack understanding of instruction semantics, so cannot reason about the program behavior

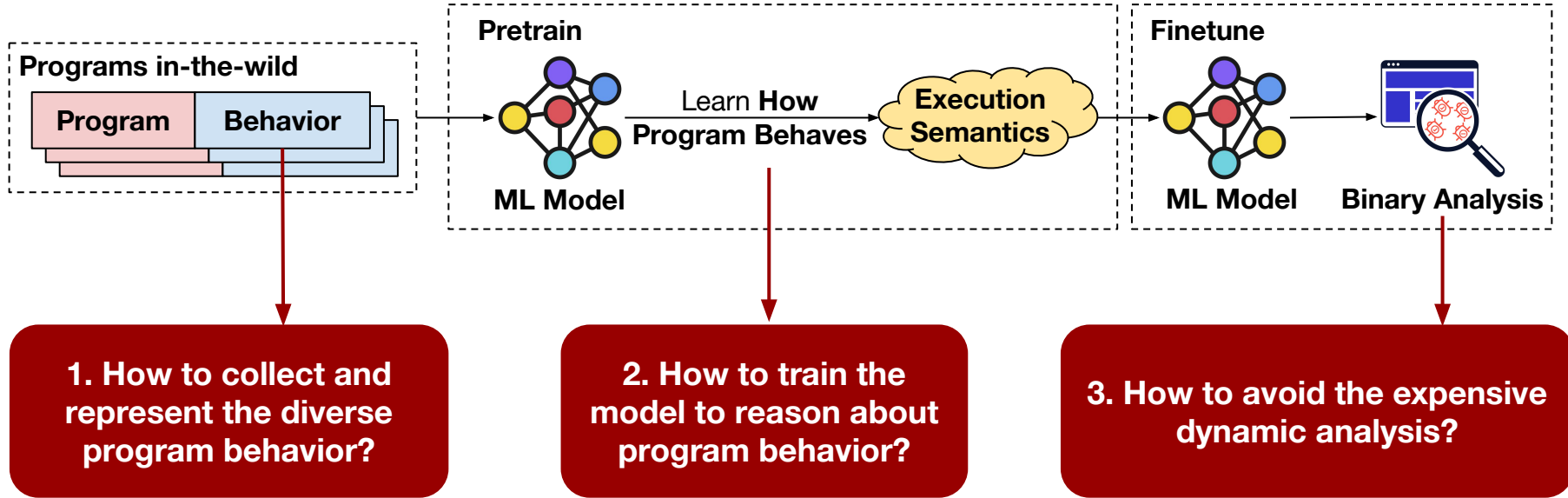
# Why not dynamic analysis?



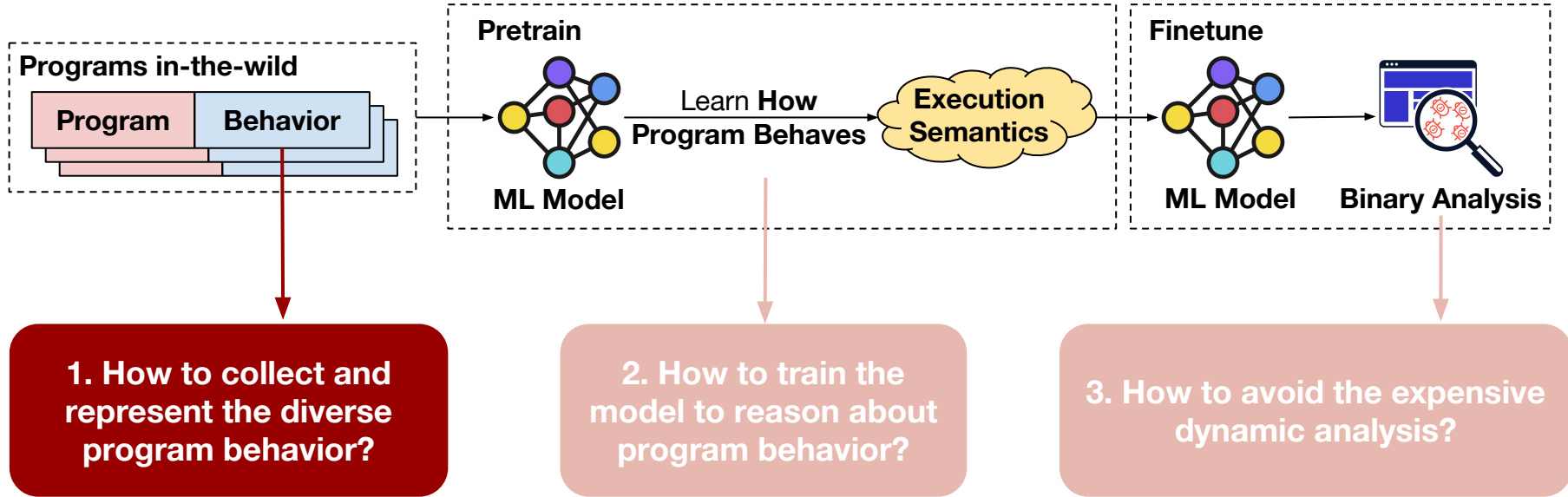
# Learning Execution Semantics and Transferring it without Dynamic Analysis



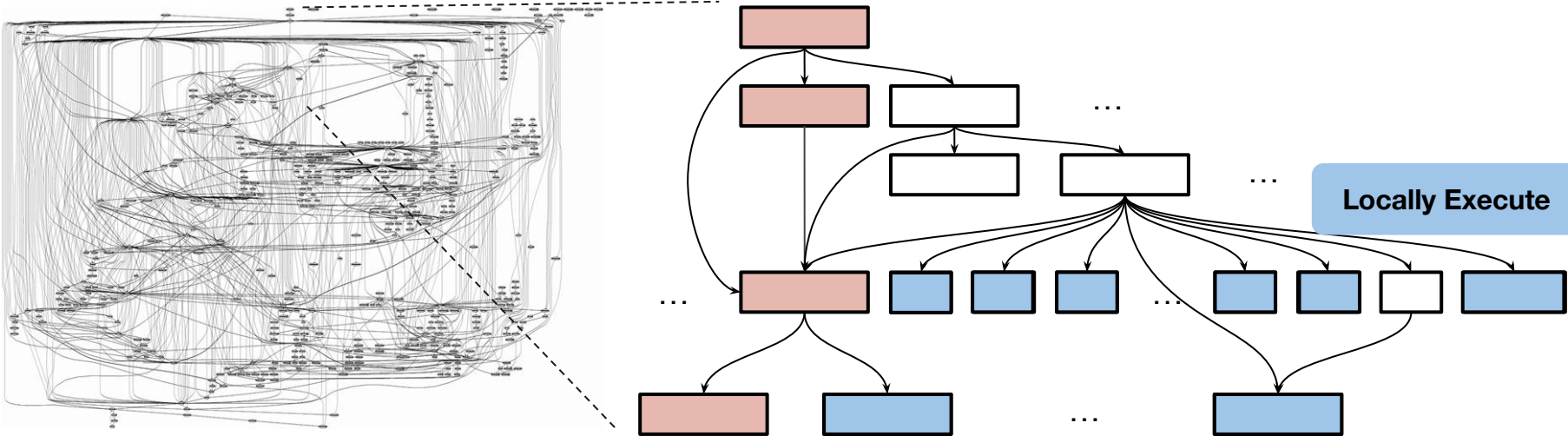
# Challenges of Learning Execution Semantics



# Challenges of Learning Execution Semantics



# How to Collect Diverse Program Behaviors?



Microsoft IIS Call Graph

**Under-Constrained Micro-Execution:  
Specify arbitrary code piece to execute**

- Expose diverse code behaviors
- Benefit large-scale pretraining on diverse execution behavior

# How to Collect Diverse Program Behaviors?

## Program instructions

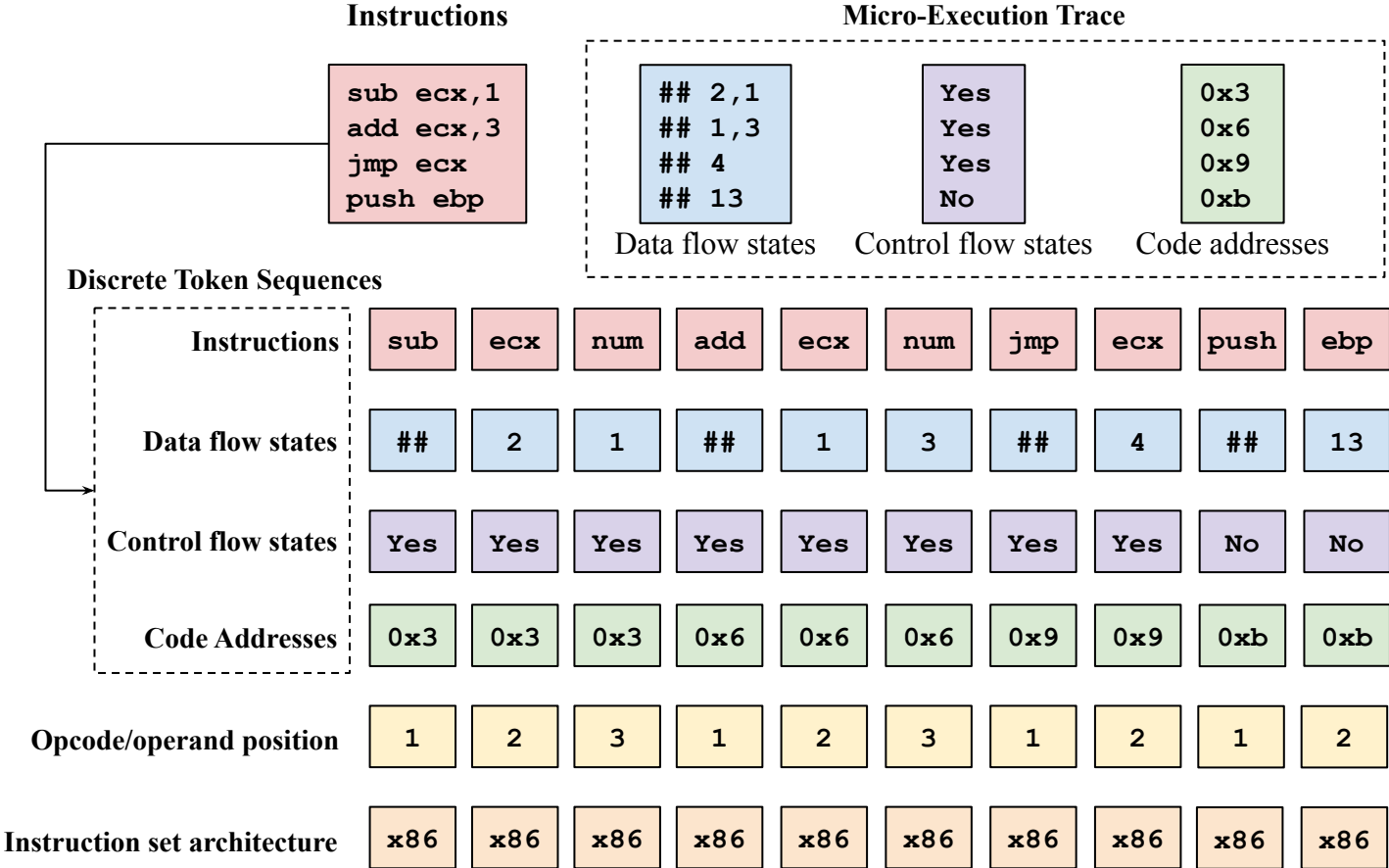
```
.....  
0x1c: mov ebp,esp  
0x1f: add [ebp+0x8],0x3  
0x26: cmp [ebp+0x8],0x2  
0x2d: jle 0x3a  
0x33: add [ebp+0x8],0x1  
0x3a: mov eax,0x1a  
.....
```

Data flow states	Control flow states	Code Addresses
..... ## 0xc,0x4	..... ✓ Yes	..... 0x1c
## [0xc+0x8],0x3	✓ Yes	0x1f
## [0xc+0x8],0x2	✓ Yes	0x26
## 0x3a	✓ Yes	0x2d
## [0xc+0x8],0x1	✗ No	0x33
## 0x1,0x1a	✗ No	0x3a
.....	.....	.....

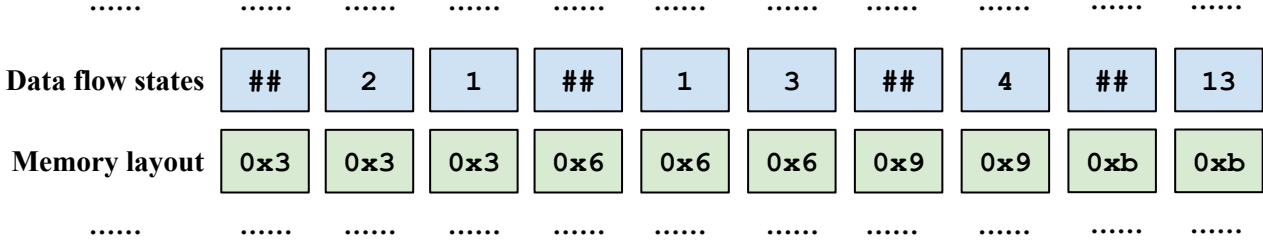
Aligned with Program Instructions



# Input Representation

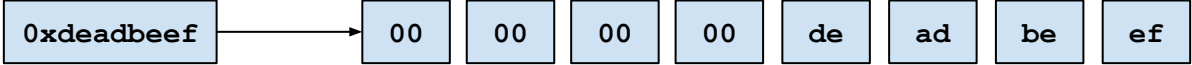


# Numerical Representation

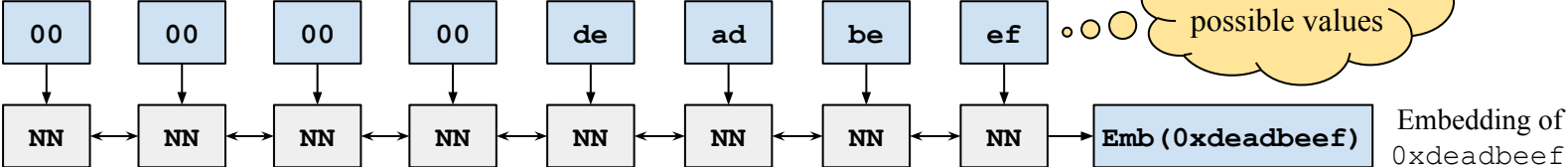


64-bit Arch:  
Vocabularies =  $2^{64}$   
**Prohibitively Large**

Pad Each Numeric Token as a Fixed-Length 8-Byte Sequence:

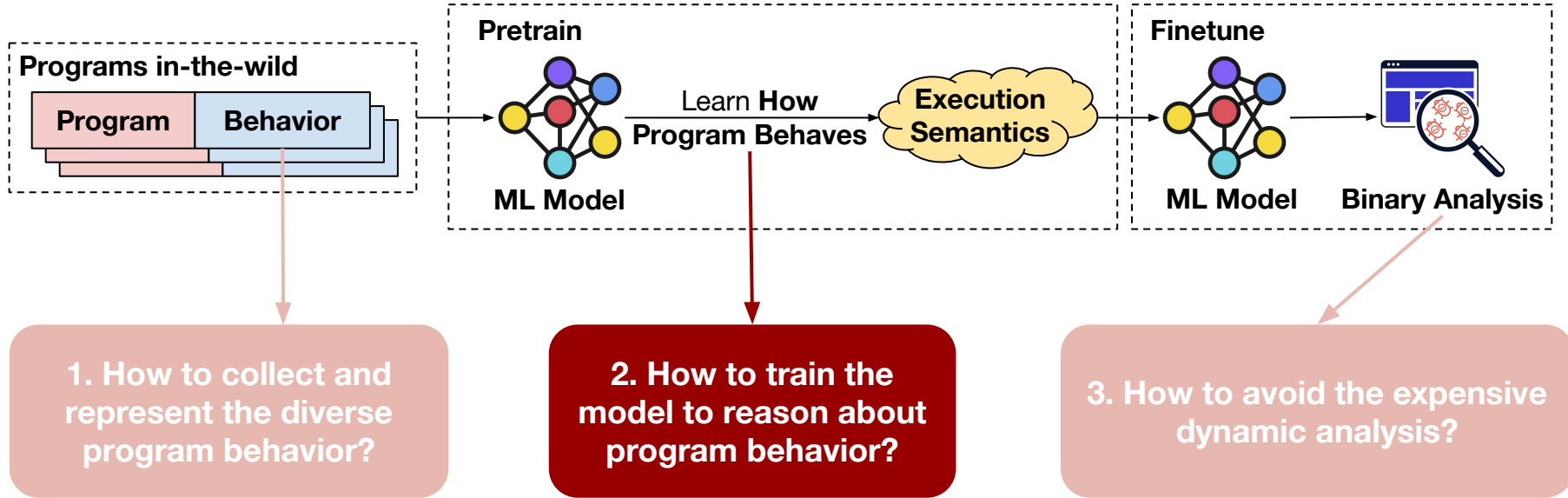


Learning to Represent the Value with a Neural Network:



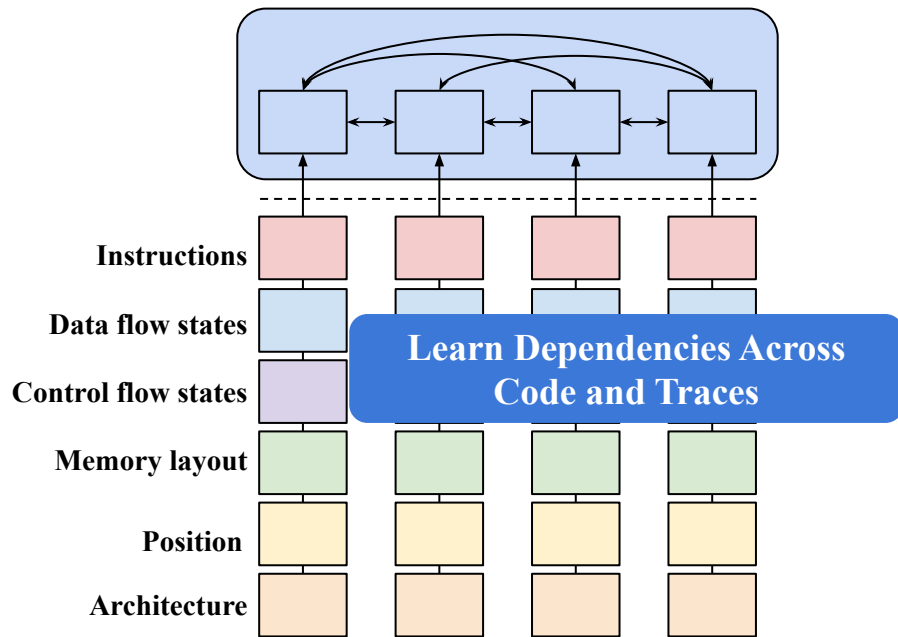
$2^{64} \rightarrow 256$   
possible values

# Challenges of Learning Execution Semantics



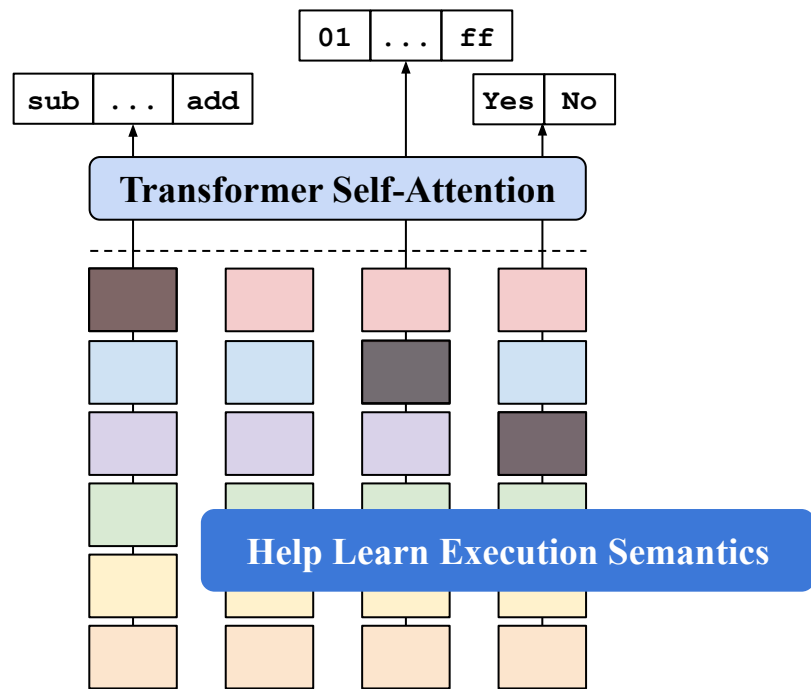
# How to train the model to reason about program behavior?

## Transformer Self-Attention



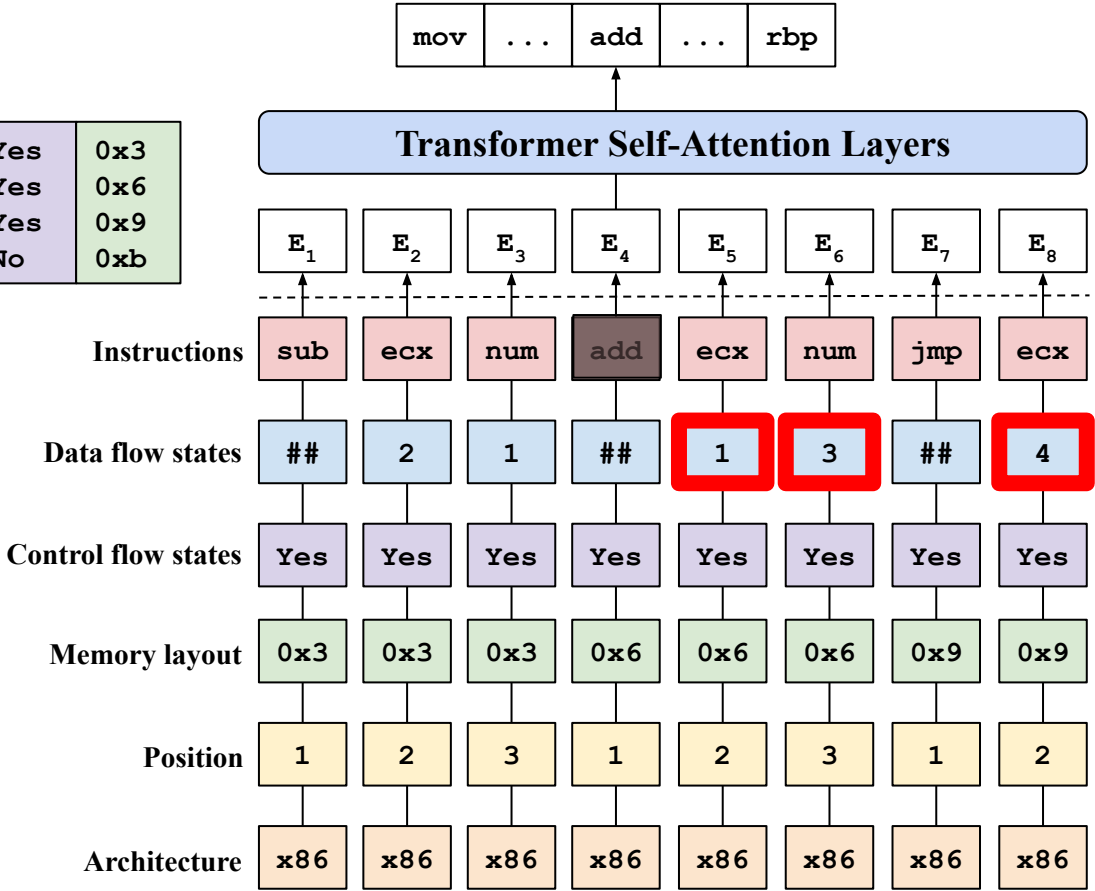
+

## Masked Language Modeling



# Motivating Example

sub ecx,1	## 2,1	Yes	0x3
add ecx,3	## 1,3	Yes	0x6
jmp ecx	## 4	Yes	0x9
	## 13	No	0xb

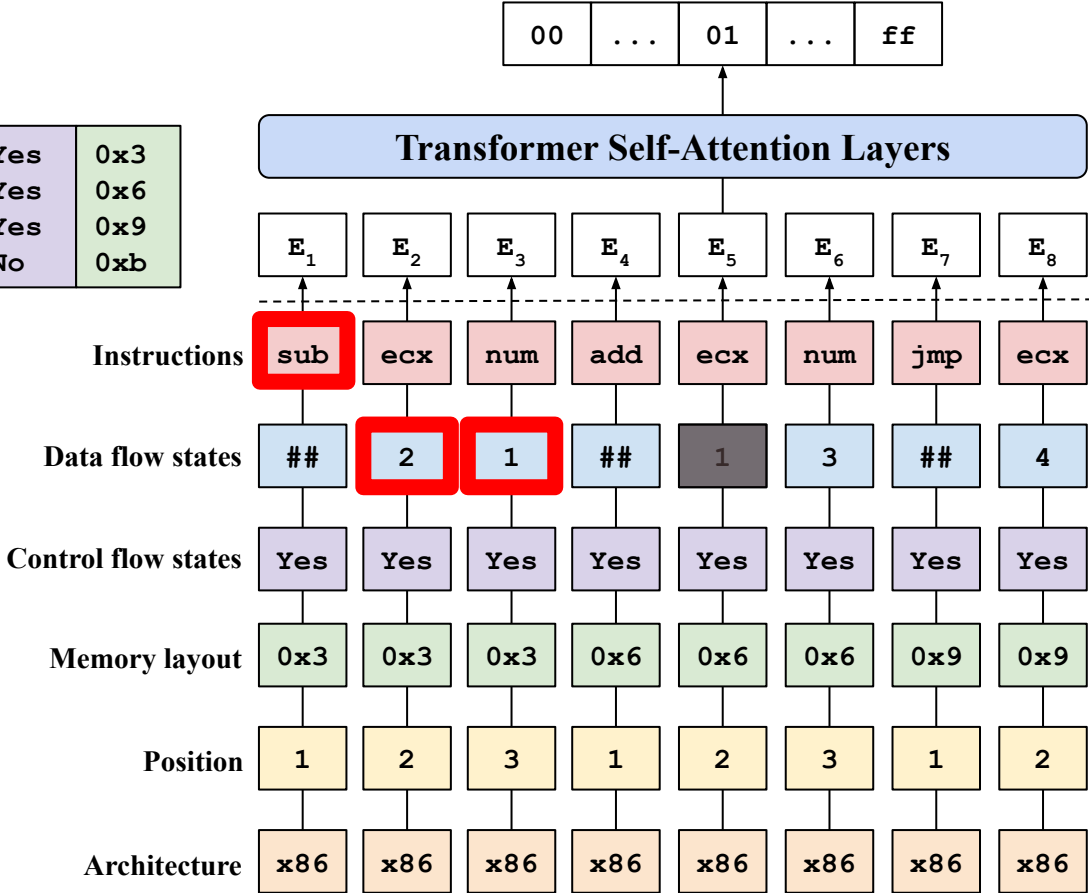


Synthesis  
1 ? 3 = 4

add

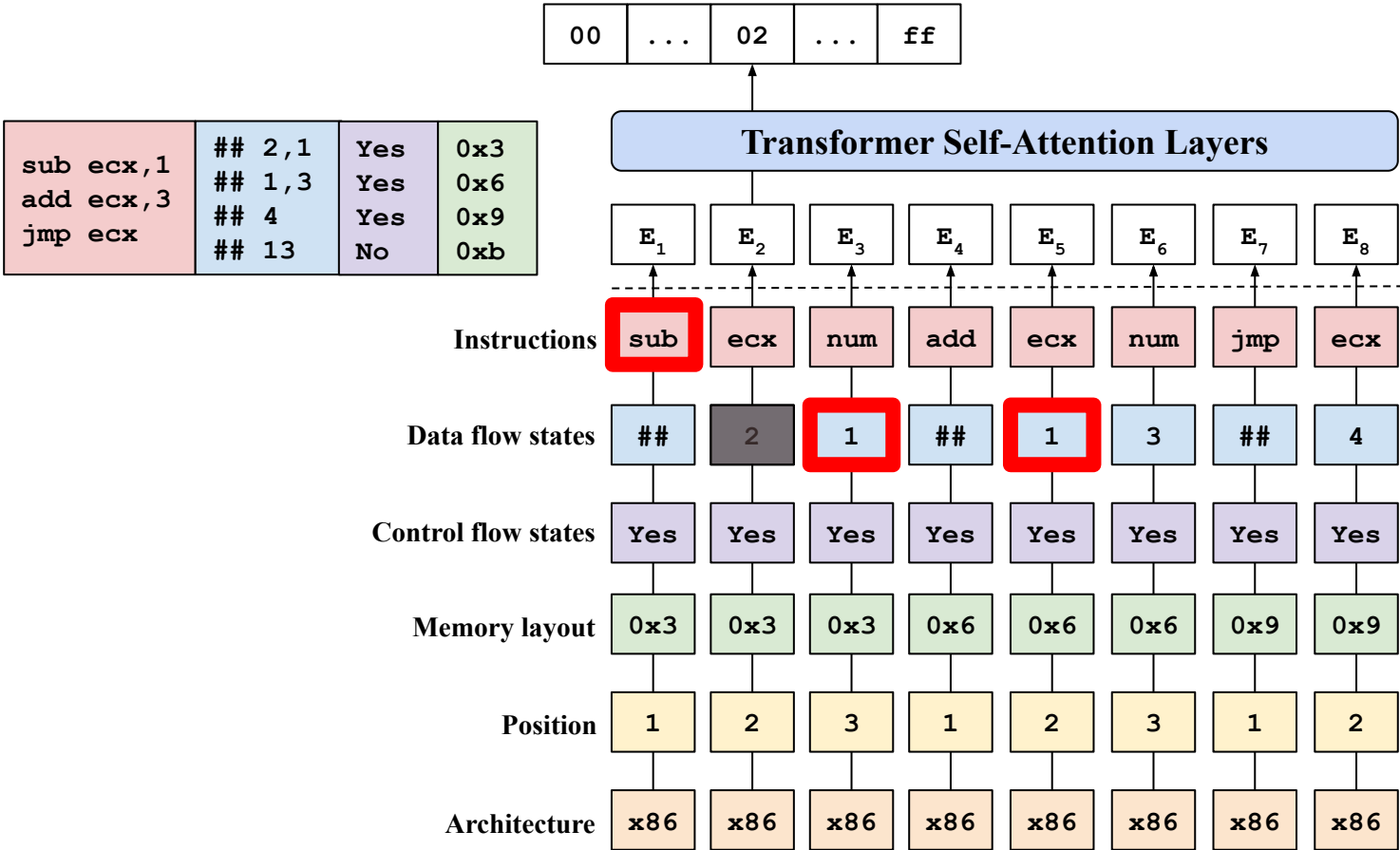
# Motivating Example

sub ecx,1	## 2,1	Yes	0x3
add ecx,3	## 1,3	Yes	0x6
jmp ecx	## 4	Yes	0x9
	## 13	No	0xb




Forward Interpretation  
 $2 - 1 = ?$

# Motivating Example

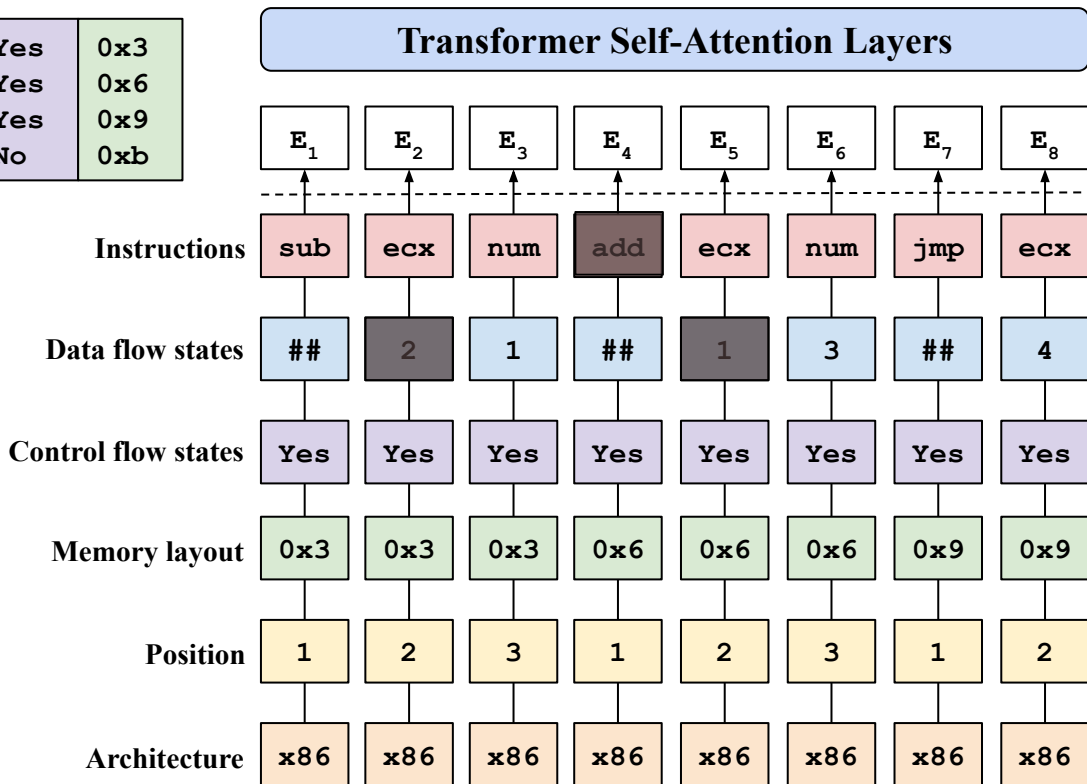


sub ecx,1	## 2,1	Yes	0x3
add ecx,3	## 1,3	Yes	0x6
jmp ecx	## 4	Yes	0x9
	## 13	No	0xb

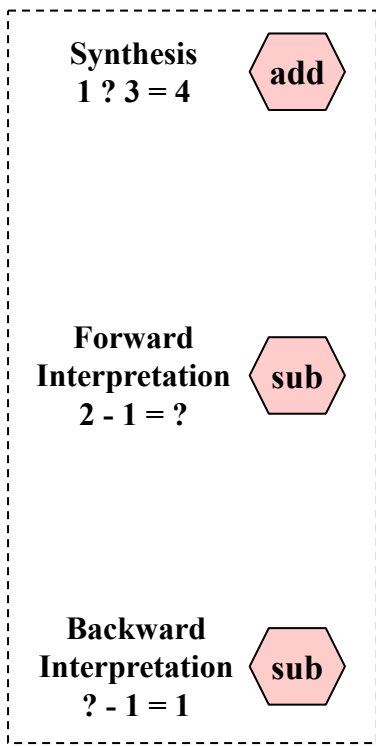
Backward Interpretation  ? - 1 = 1

# Motivating Example

sub ecx,1	## 2,1	Yes	0x3
add ecx,3	## 1,3	Yes	0x6
jmp ecx	## 4	Yes	0x9
	## 13	No	0xb



- Program Interpretation
- Program Synthesis





# How to train the model to reason about program behavior?

How to train the model to reason about **control flow**?

## Mask control flow states

add eax, 3	Yes
cmp eax, 8	Yes
jle 0x8	Yes
mov ecx, ebx	No

if (eax+3) > 8:  
  mask=Yes  
else:  
  mask=No

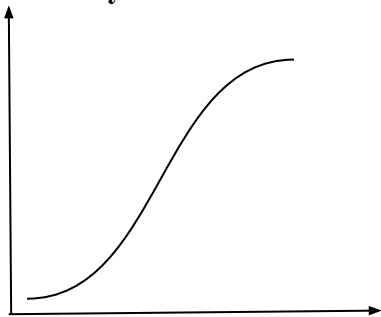
How to train the model to reason about instruction **compositions**?

## Mask more

add eax, 3	##	3, 3
sub ebx, eax	##	12, 6
sub ebx, 8	##	6, 8
mov ecx, ebx	##	0, -

ecx =  
 $(12 - (eax + 3)) - 8$

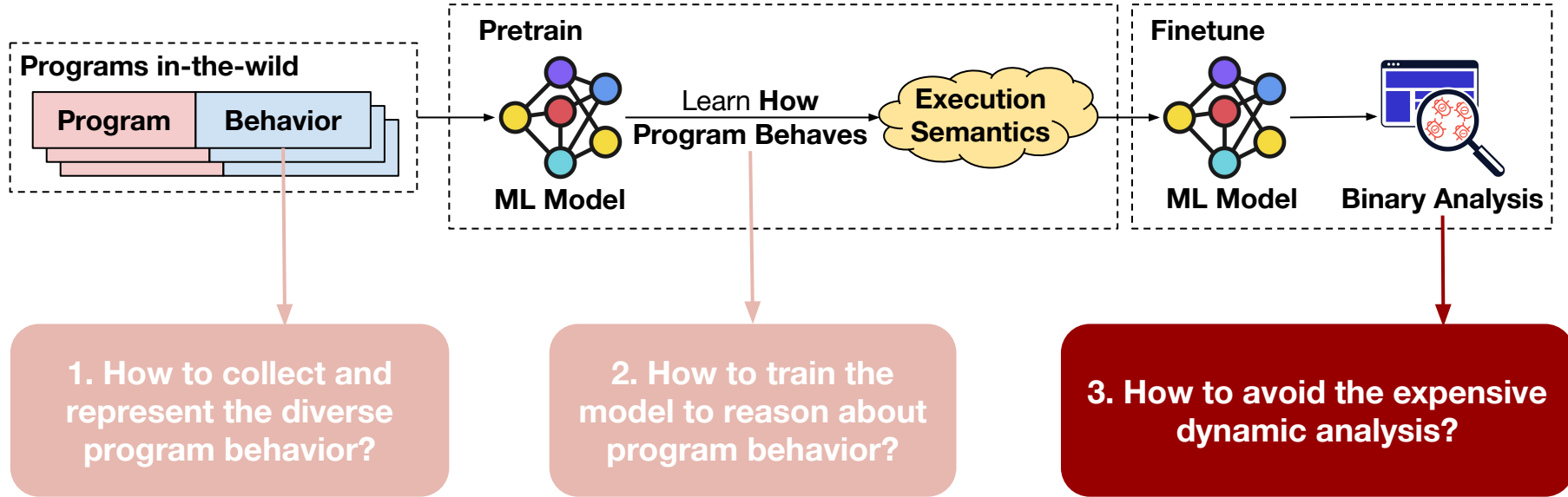
Compositionality



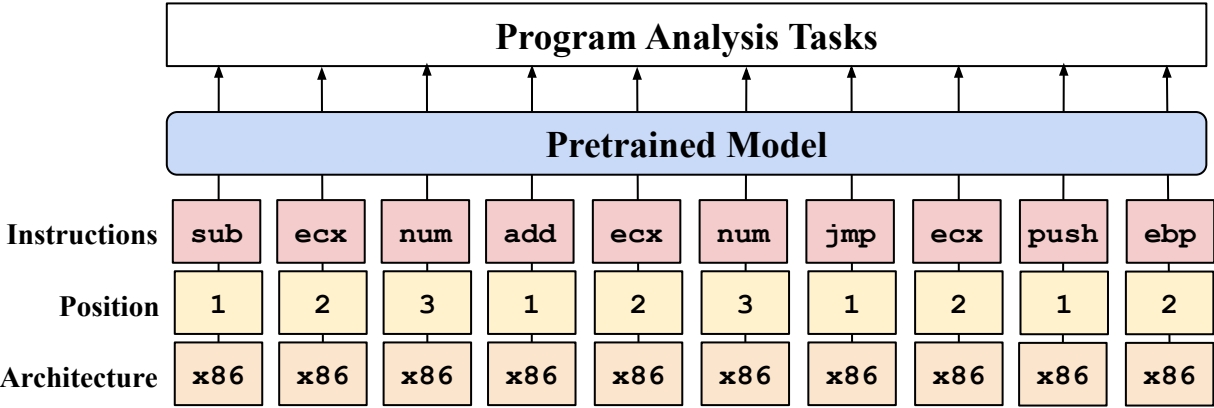
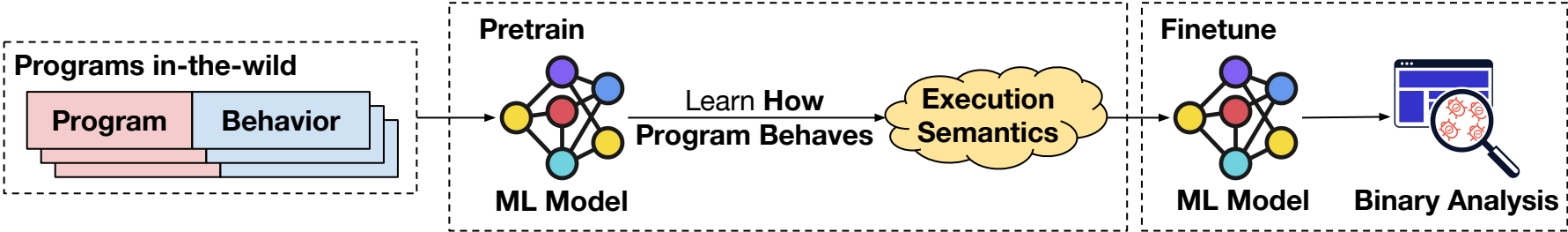
Training

- **Curriculum**: linearly increase % masks by training iterations
- **Randomized** mask selection at each iteration and samples

# Challenges of Learning Execution Semantics



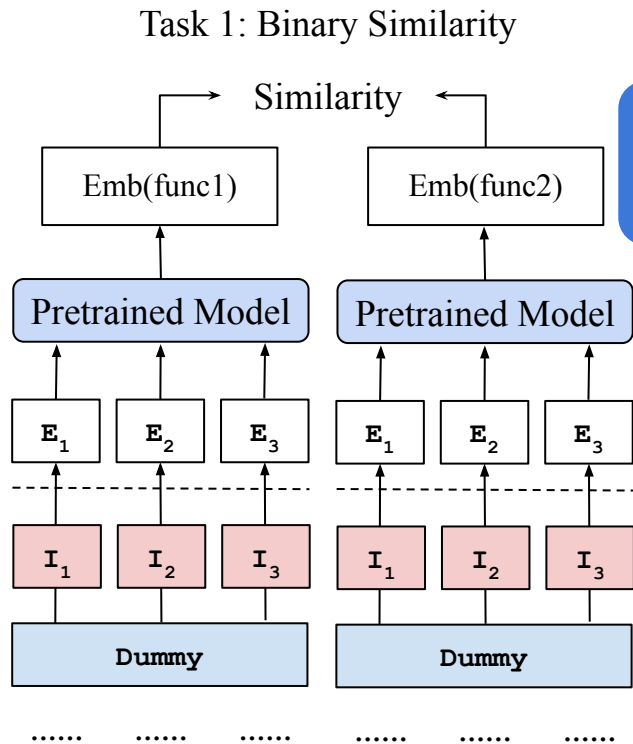
# How to Avoid the Expensive Dynamic Analysis?



Execution-Aware Code Representation

Finetuning as Static Analysis

# How Much do Learned Execution-Aware Program Representations Help?

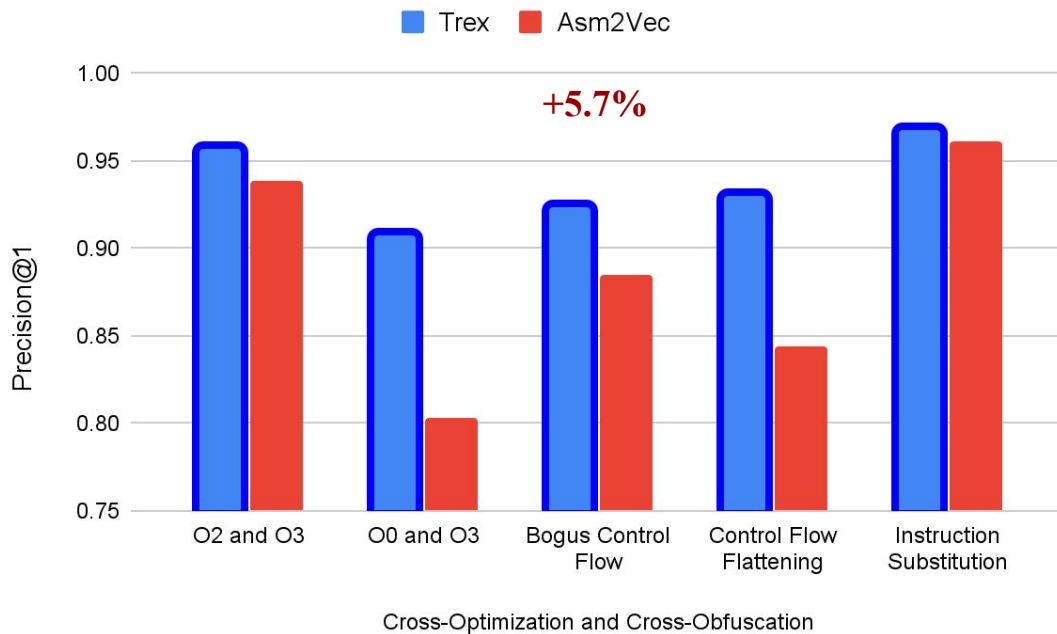
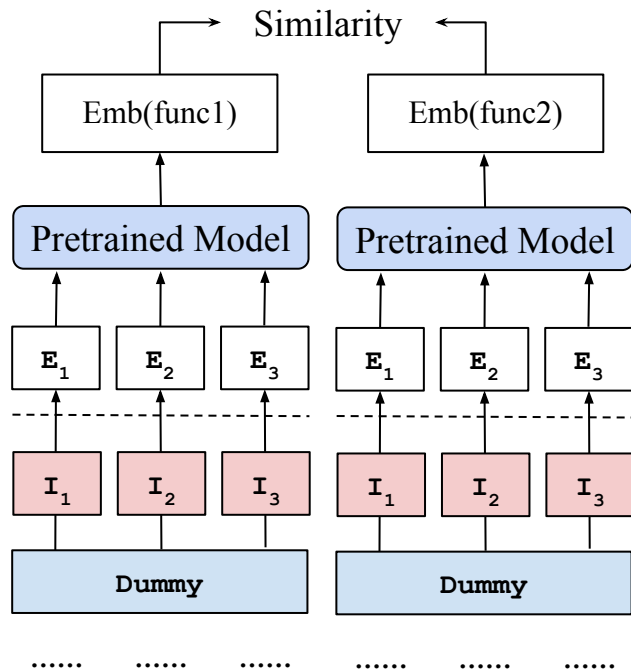


- Vulnerability search
- Malware detection
- Software patching



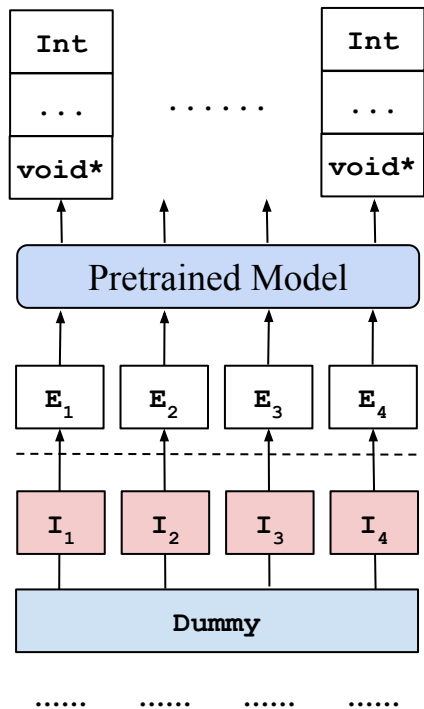
# Finetuning for Matching Semantically Similar Binary Functions

Task 1: Binary Similarity

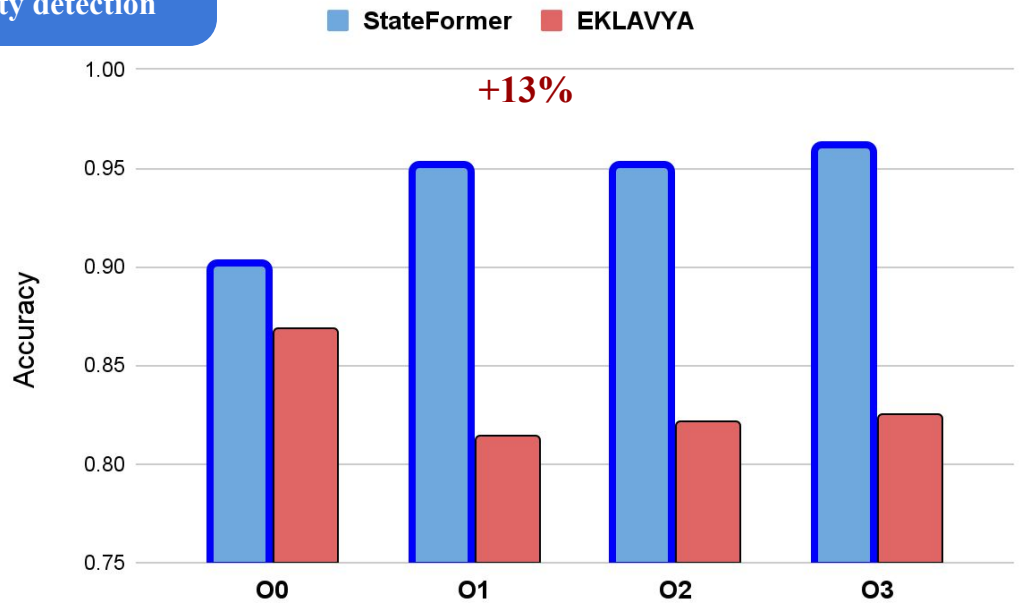


# Finetuning for Predicting Function Signatures and Type Inference

Task 2: Type Inference



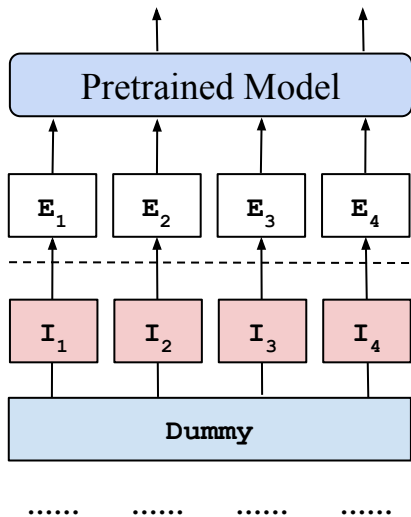
- Security retrofitting
- Decompilation
- Vulnerability detection



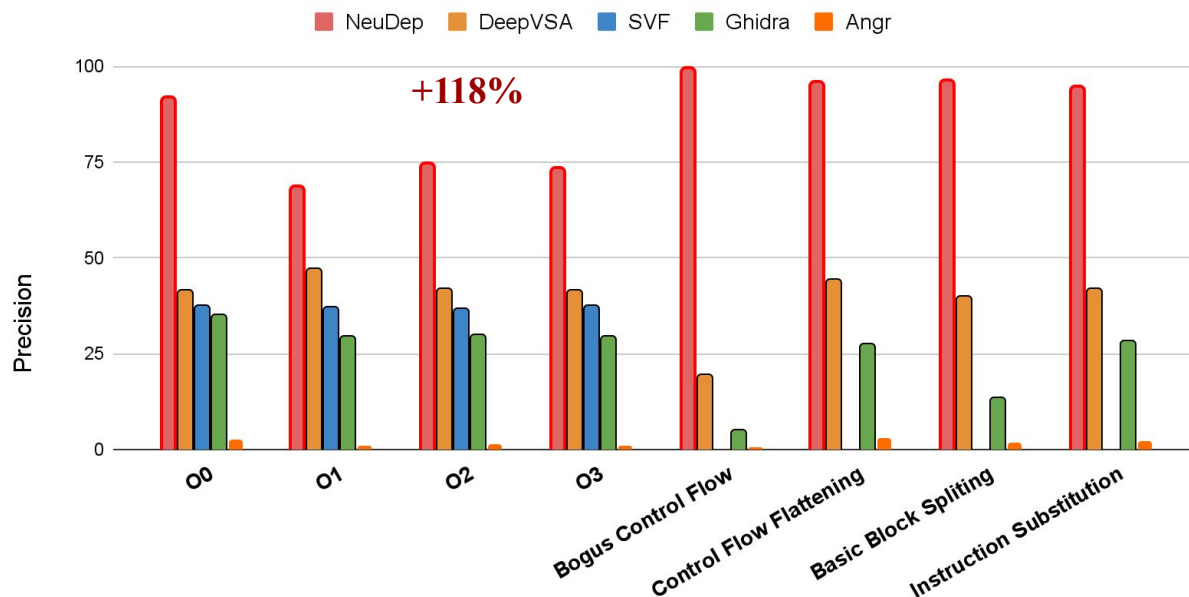
# Finetuning for Analyzing Memory Dependence

## Task 3: Memory Dependence Analysis

$P(I_2 \text{ and } I_4 \text{ is dependent}) = 0.8$



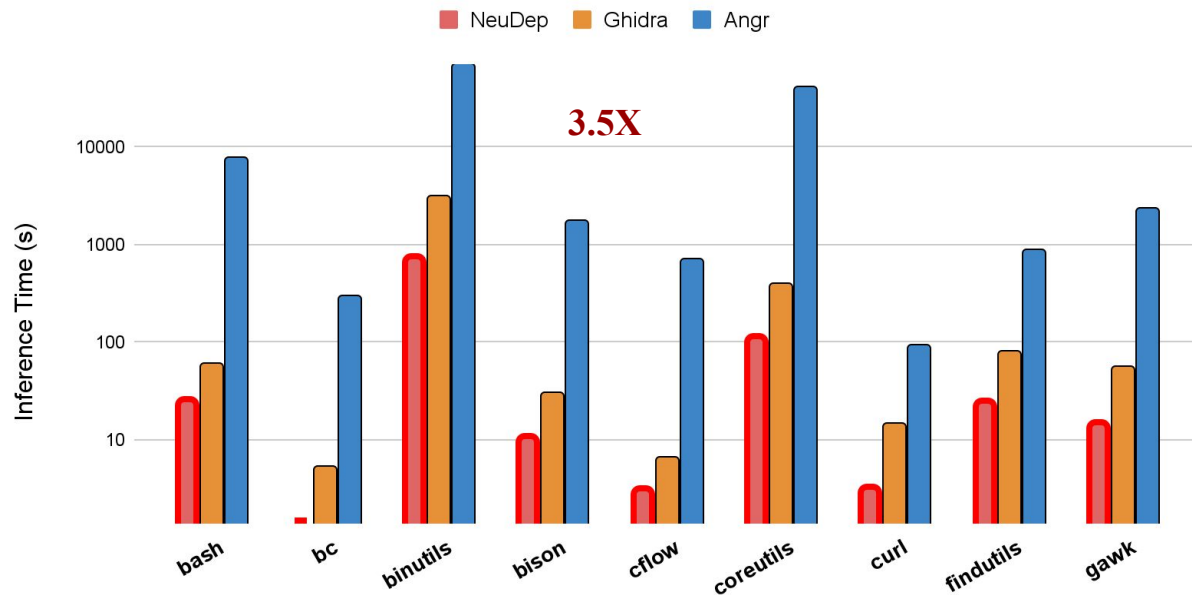
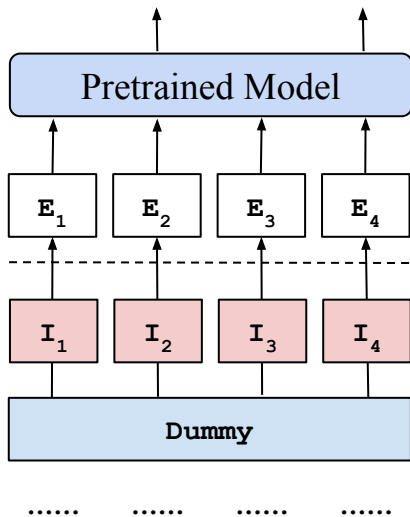
- Taint analysis
- Malware Analysis



# Finetuning for Analyzing Memory Dependence: Inference Time

## Task 3: Memory Dependence Analysis

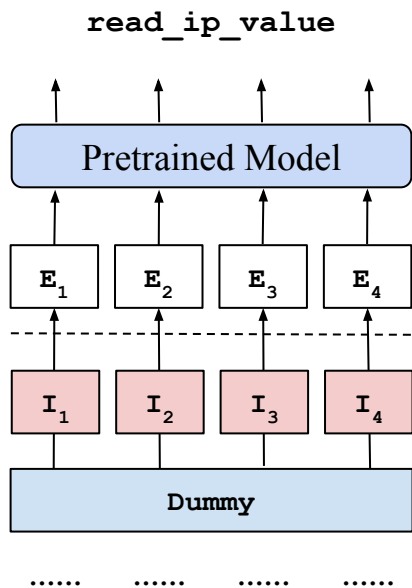
$P(I_2 \text{ and } I_4 \text{ is dependent}) = 0.8$





# Finetuning for Function Name Prediction and Memory Region Prediction

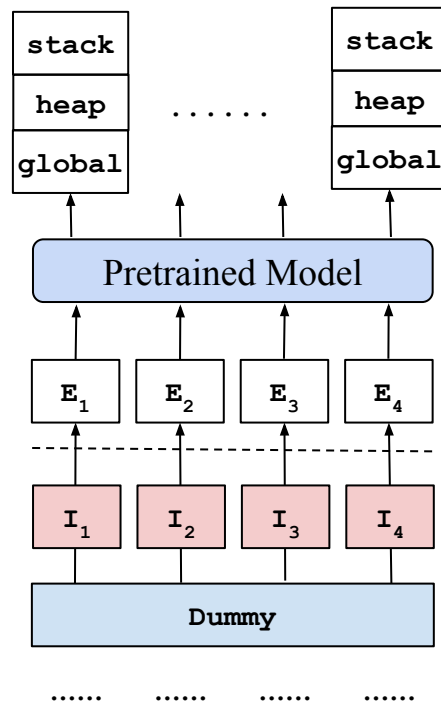
## Task 4: Function Name Prediction



- Decompilation
- Reverse engineering

**+35% Across All Architectures and Optimizations**

## Task 5: Memory Region Prediction



- Value set analysis
- Taint analysis

**+7.9% Across All Types of Memory Regions**

# Case Studies: Vulnerability Search in Firmware

CVE	Library	Description
CVE-2019-1563	OpenSSL	Decrypt encrypted message
CVE-2017-16544	BusyBox	Allow executing arbitrary code
CVE-2016-6303	OpenSSL	Integer overflow
CVE-2016-6302	OpenSSL	Allows denial-of-service
CVE-2016-2842	OpenSSL	Allows denial-of-service
CVE-2016-2182	OpenSSL	Allows denial-of-service
CVE-2016-2180	OpenSSL	Out-of-bounds read
CVE-2016-2178	OpenSSL	Leak DSA private key
CVE-2016-2176	OpenSSL	Buffer over-read
CVE-2016-2109	OpenSSL	Allows denial-of-service
CVE-2016-2106	OpenSSL	Integer overflow
CVE-2016-2105	OpenSSL	Integer overflow
CVE-2016-0799	OpenSSL	Out-of-bounds read
CVE-2016-0798	OpenSSL	Allows denial-of-service
CVE-2016-0797	OpenSSL	NULL pointer dereference
CVE-2016-0705	OpenSSL	Memory corruption

16 Vulnerabilities (Compiled in x86)

Search

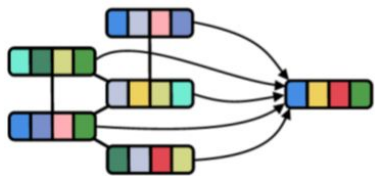


15 CVEs

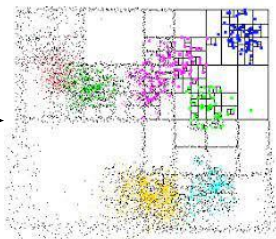
16 CVEs

7 CVEs

14 CVEs



Learned Function Embeddings

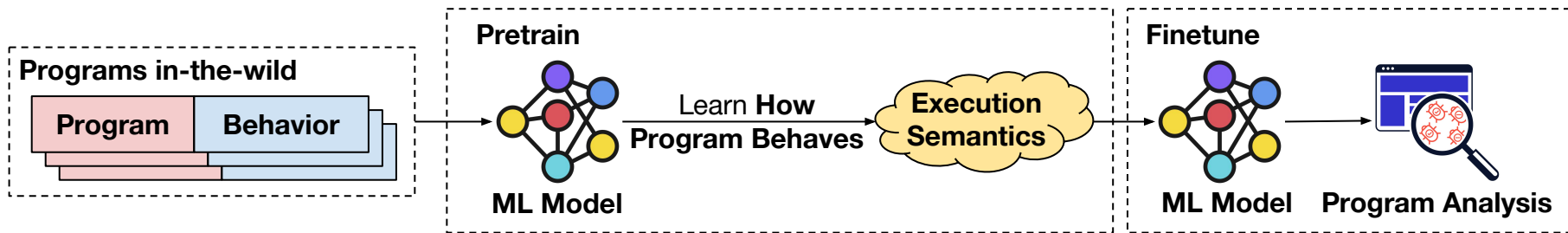


Approximate Nearest Neighbors



Search over **1.4M** functions within **6.3 seconds**

# Summary: Learning Program Semantics via Execution-Aware Pre-training Improves Program Analysis



**Precise:** Outperforms the state-of-the-art by up to **118%**

**Efficient:** Speedup over the off-the-shelf tool by up to **98.1x**



**Generalizable and Robust:** Remains accurate across



Compilers



Architectures

-O0 -O1 -O2  
-O3 -Od Ox

Optimizations



Obfuscations

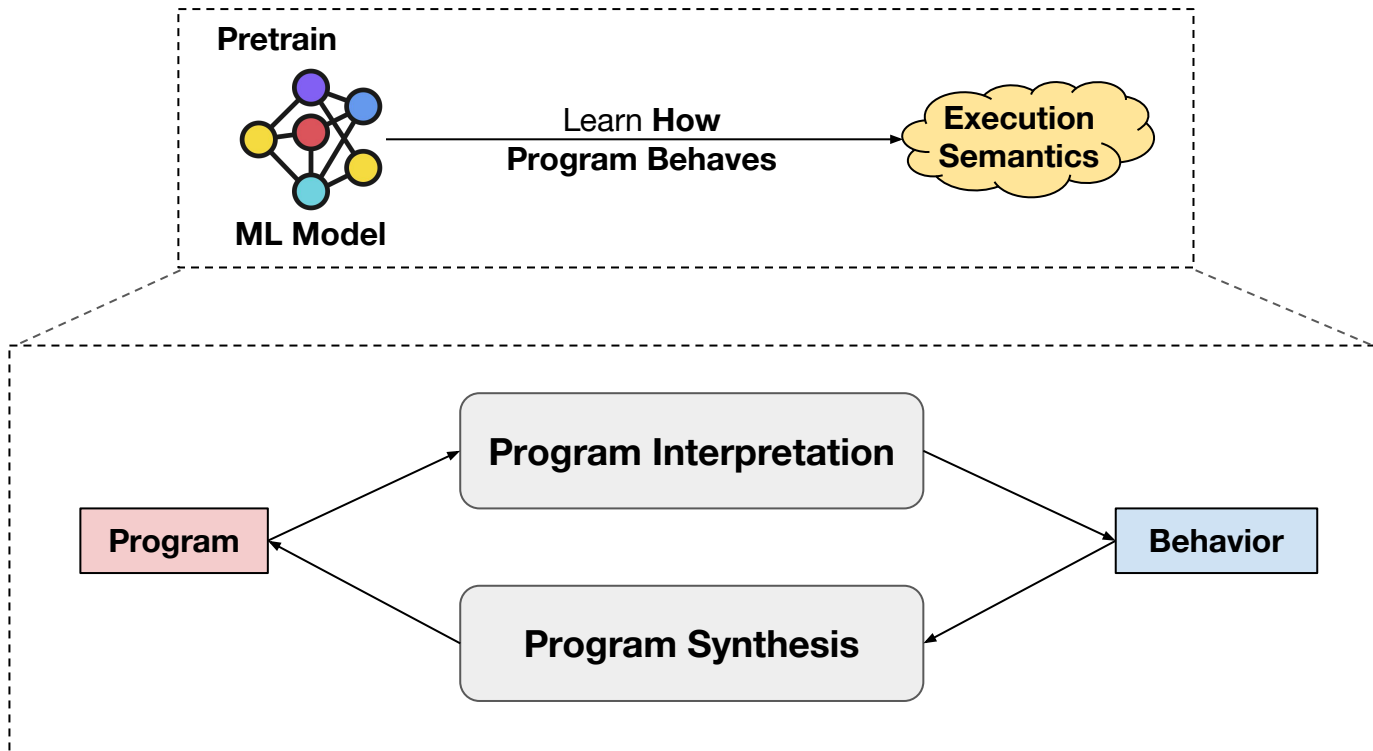


Google



GHIDRA

# Limitation: Learning Execution-Aware Program Representations is Challenging



**Extremely challenging to learn precise semantics**

## Limitation: What the Model has Learned during Pretraining?

Instructions	Dataflow states	Instructions	Dataflow states
..... sub ecx,0x3 add ecx,0x4 .....	..... ## 0x5,0x3 ## 0x2,0x4 .....	..... sub ecx,0x3 add ecx,0x4 .....	..... ## 0x43,0x3 ## 0x3d,0x4 .....

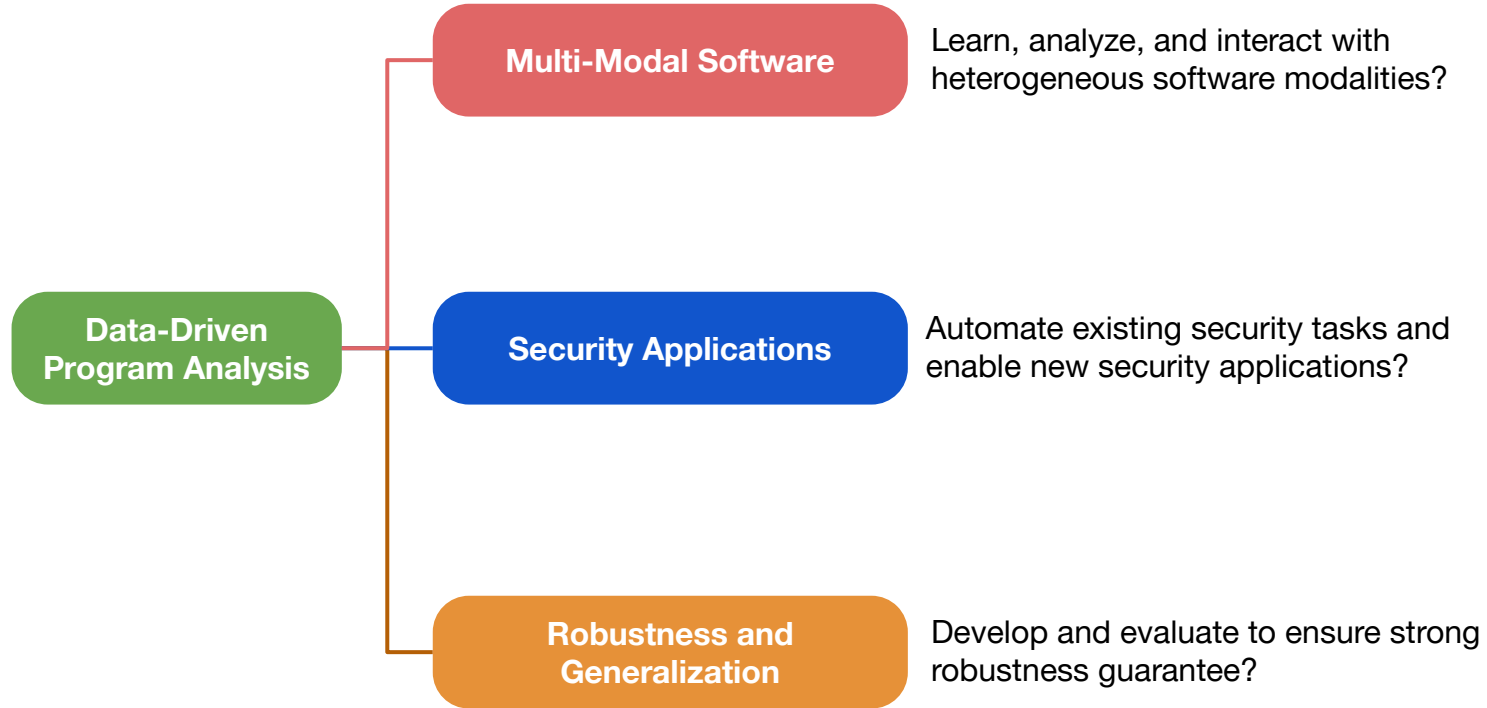
Perturb dataflow states from 0x5 to 0x43

Ground-truth	Top-1	Top-2
0x2	0x2 (98%)	0x3 (2%)
0x3d	0x3a (28%)	0x33 (13%)

**Does not extrapolate well**

# Exciting Future Work

Democratize the development of secure software system via

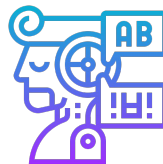
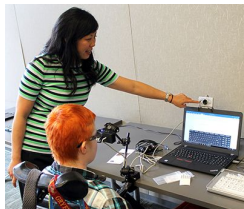


# Software is Inherently Heterogeneous and Multi-Modal

Noisy



Physical Interaction



Natural Language

```

    """
    Returns the year and
    quarter of an input date.
    Args:
        date: XXXX-XX-XX
    Returns:
        year: XXXX
        quarter: XXX
    """
    
```

```

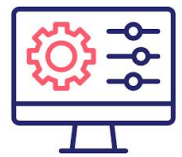
while i < n
  invariant i <= n
{
  i = i + 1;
}
    
```

```

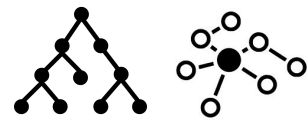
port=81
rootpath=...
execution=True
verification=True
    
```



Software Program



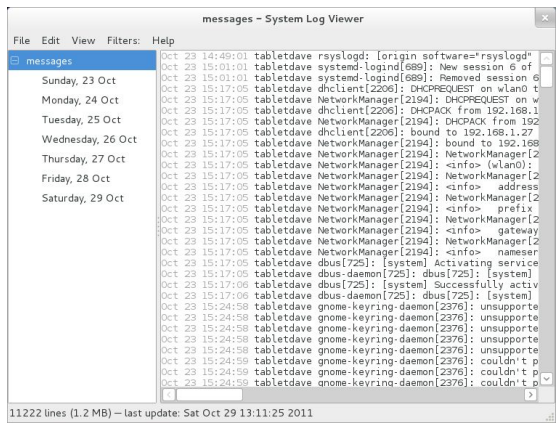
Configurations



Program Activity



Specifications



# How to learn from heterogeneous software modalities?



Physical Interaction



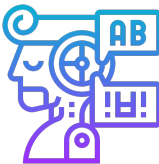
Specifications



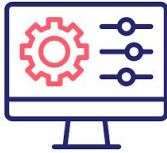
Software Program



Program Activity



Natural Language



Configurations



# How to learn from heterogeneous software modalities?



Physical Interaction



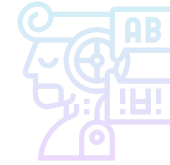
Specifications



Software Program



Program Activity



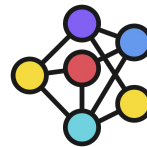
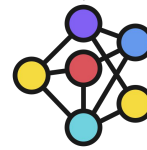
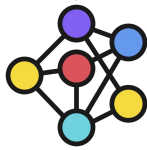
Natural Language



Configurations

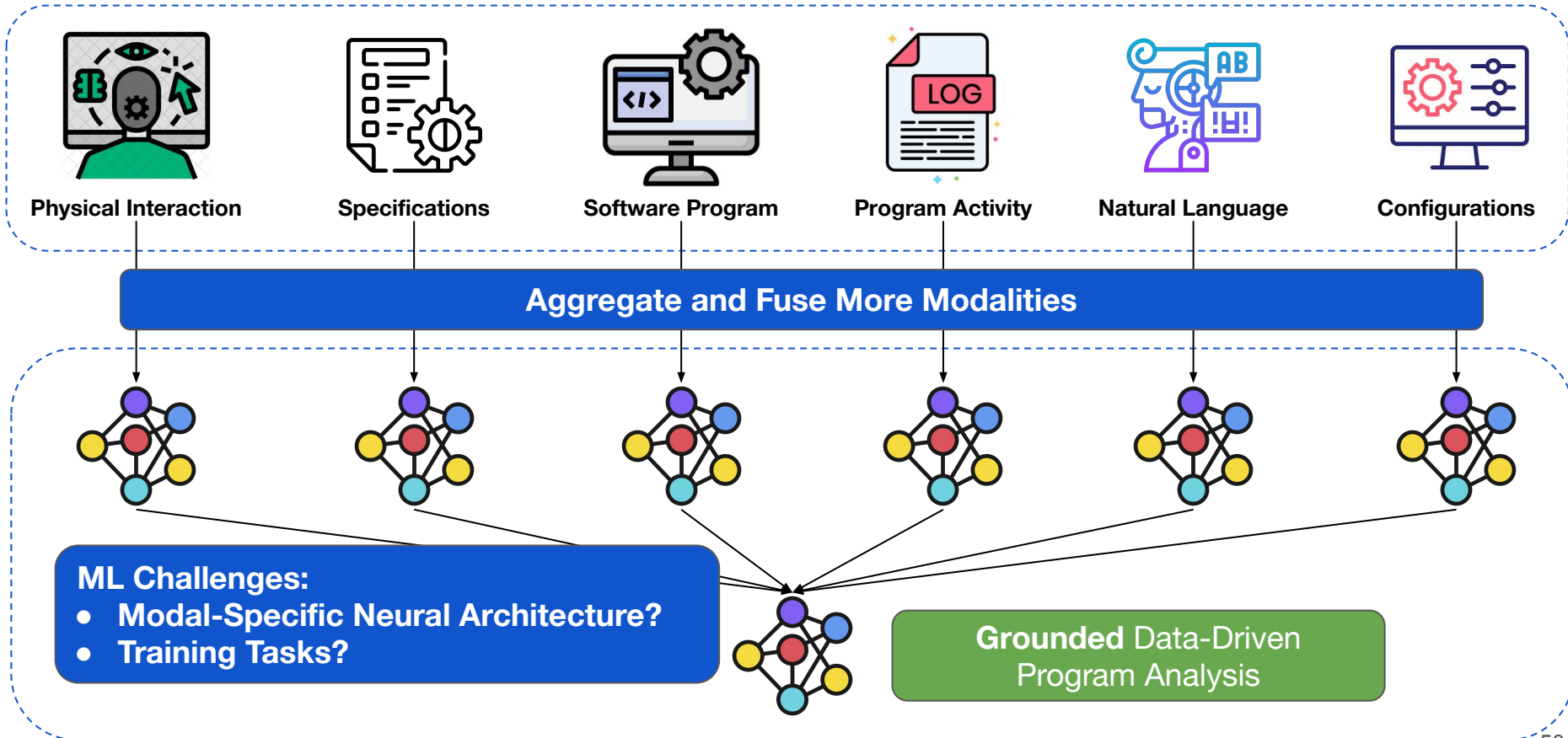
**Execution-Aware Program Representation**

Pei et al., FSE'21, FSE'22, TSE'22, CCS'22

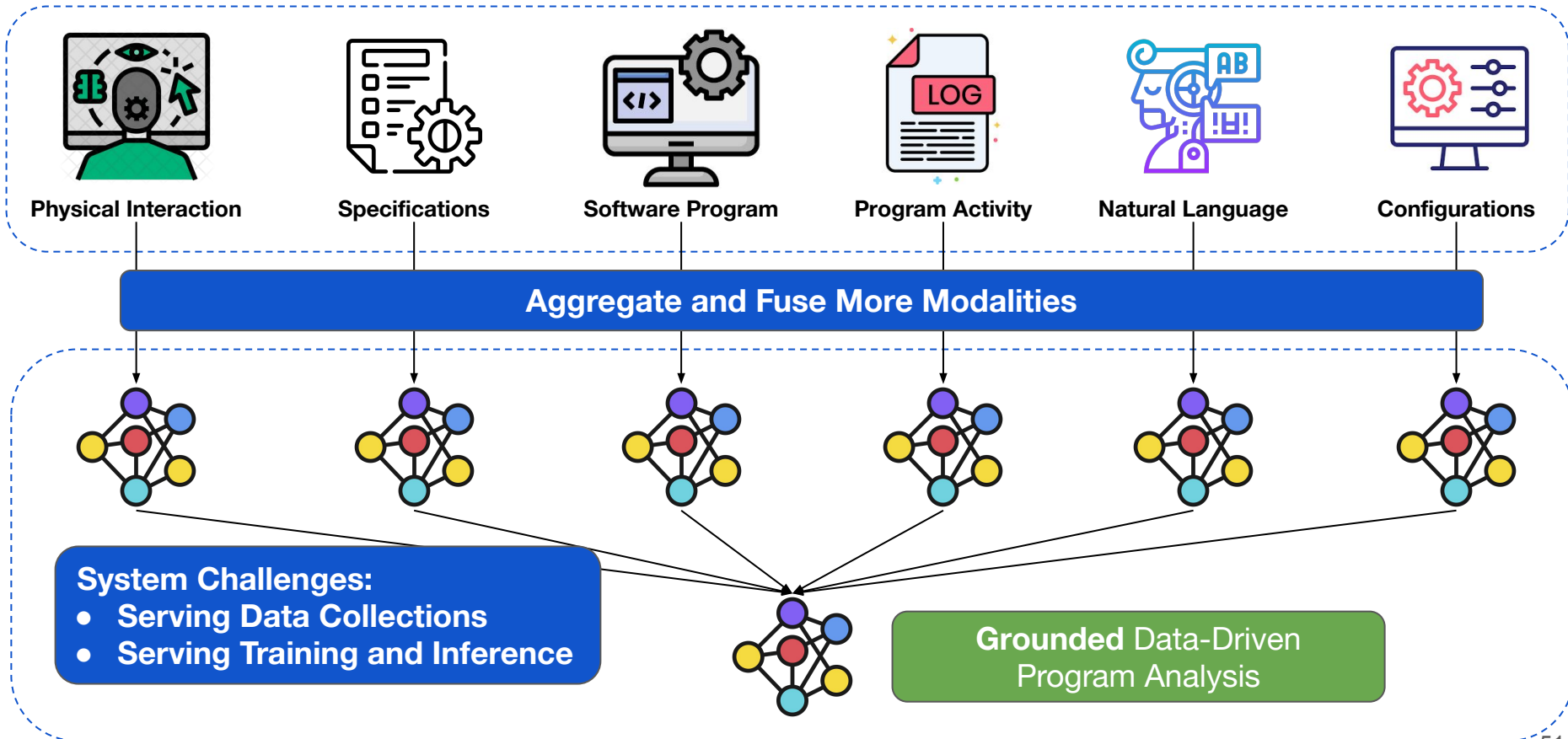


**Execution-Aware Data-Driven  
Program Analysis**

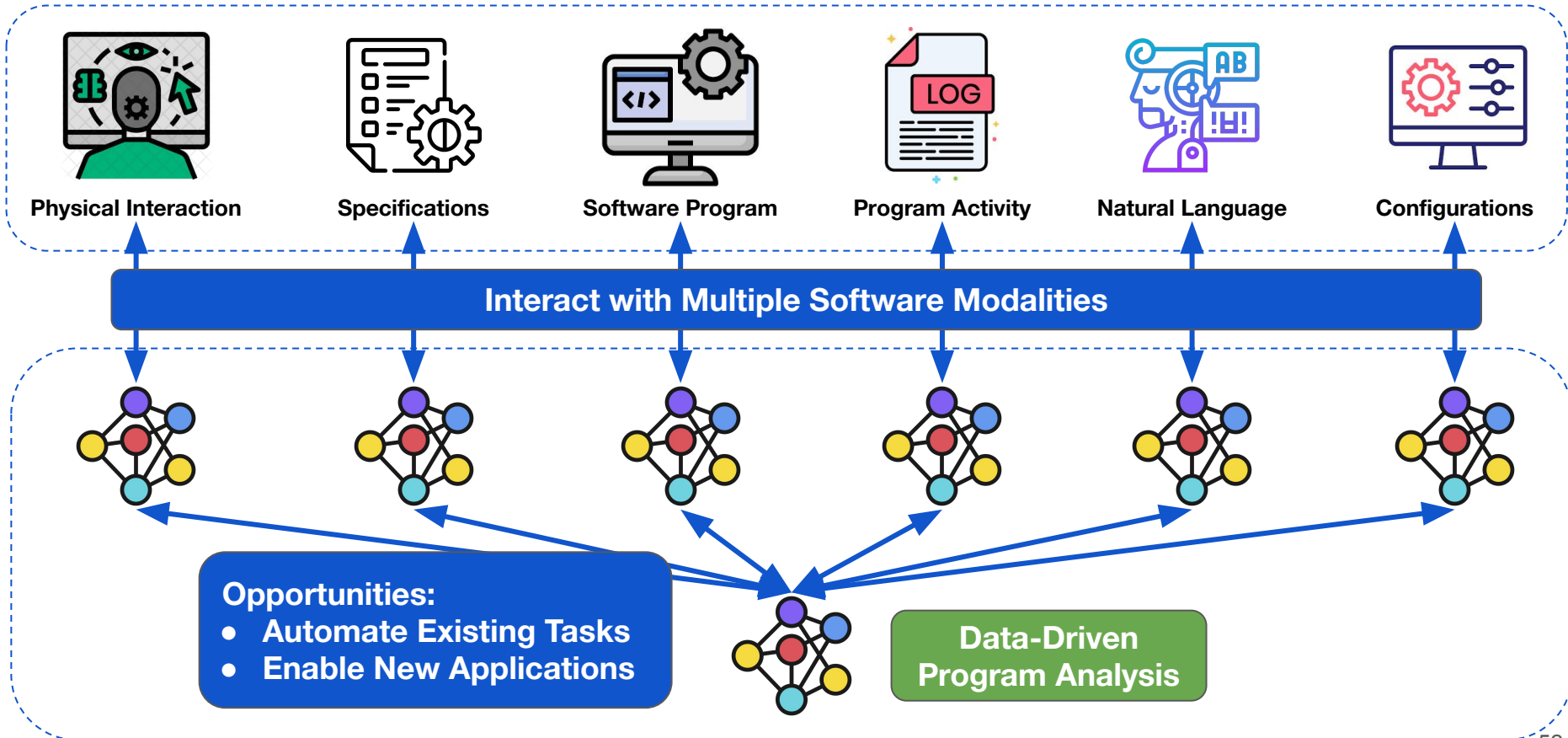
# How to learn from heterogeneous software modalities?



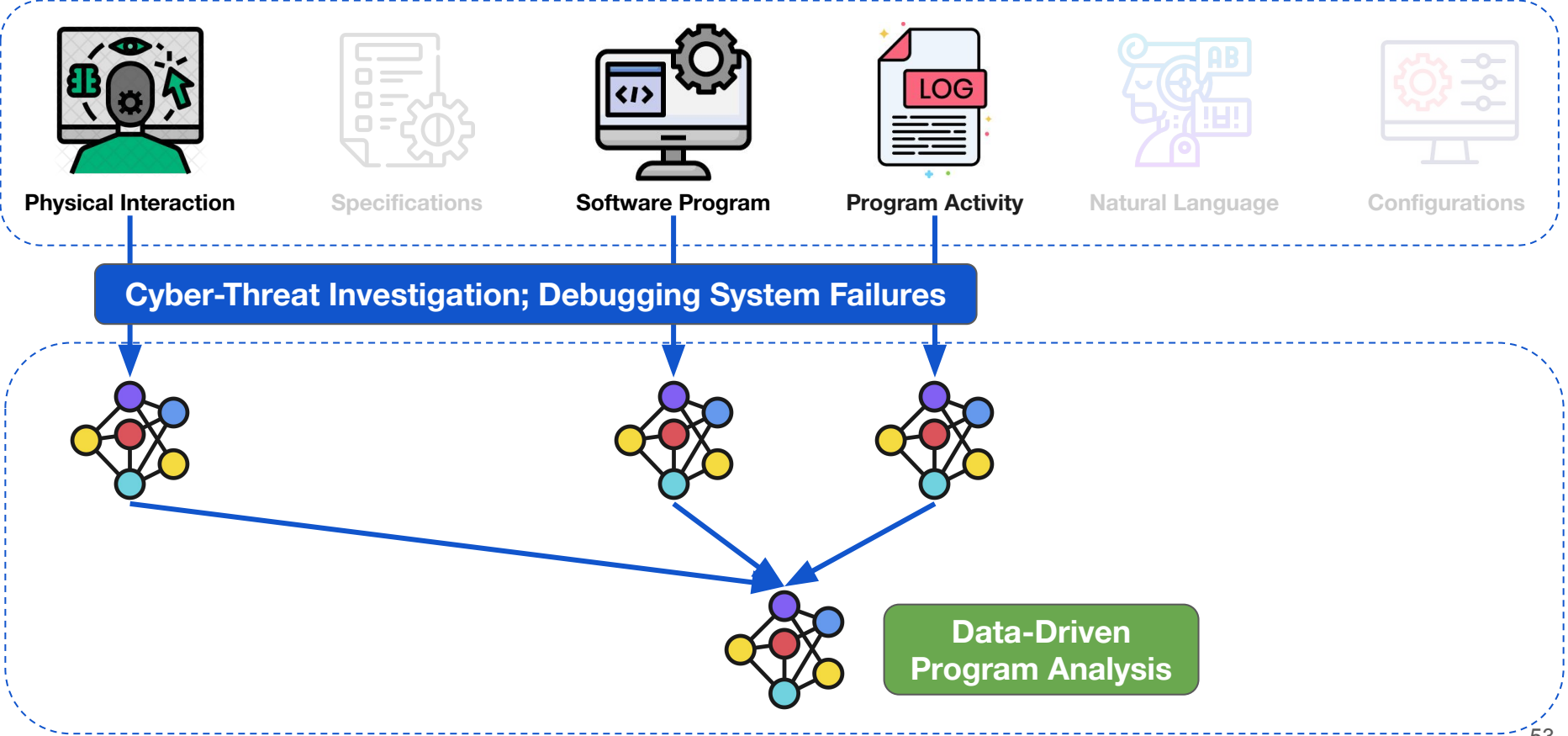
# How to learn from heterogeneous software modalities?



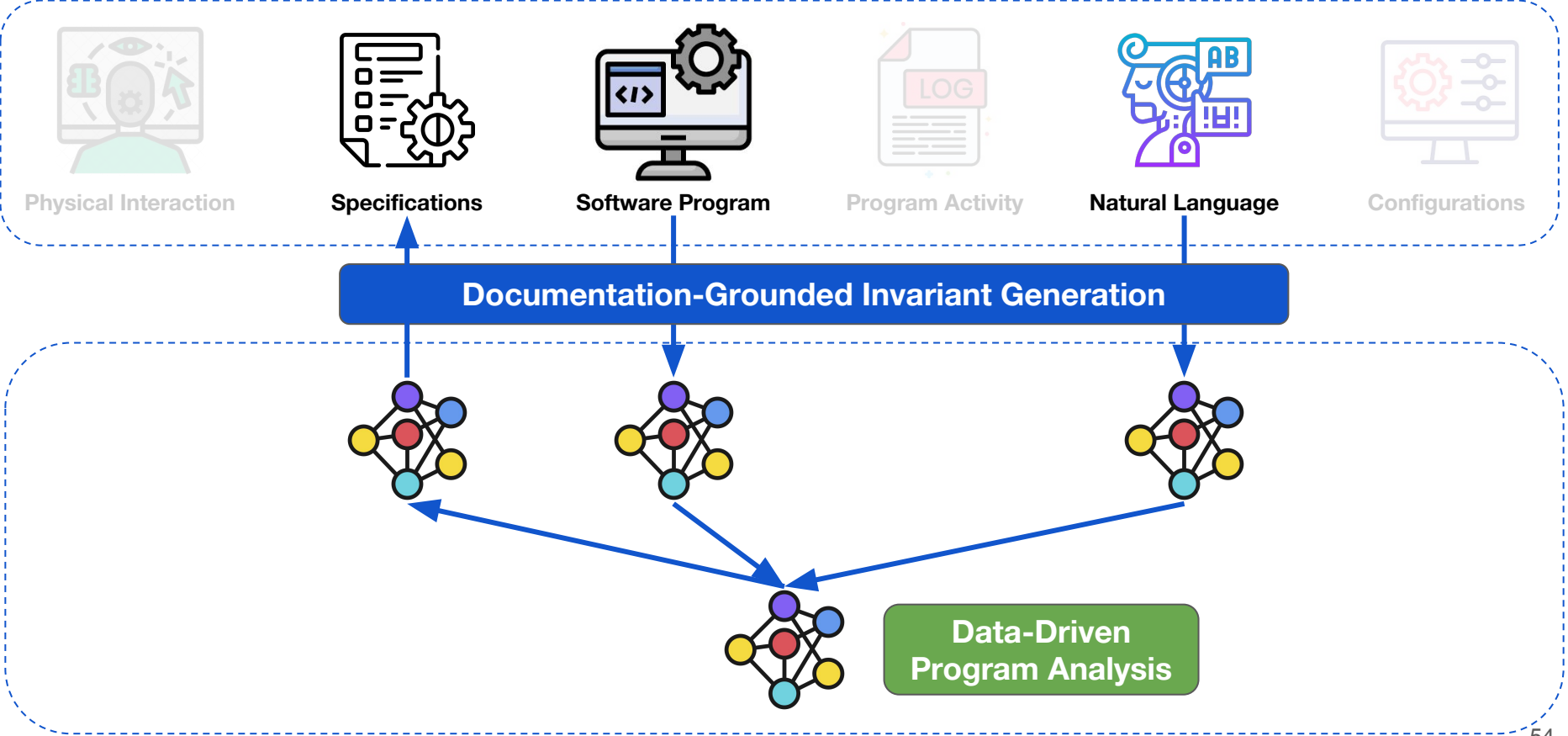
# How to interact with heterogeneous software modalities?



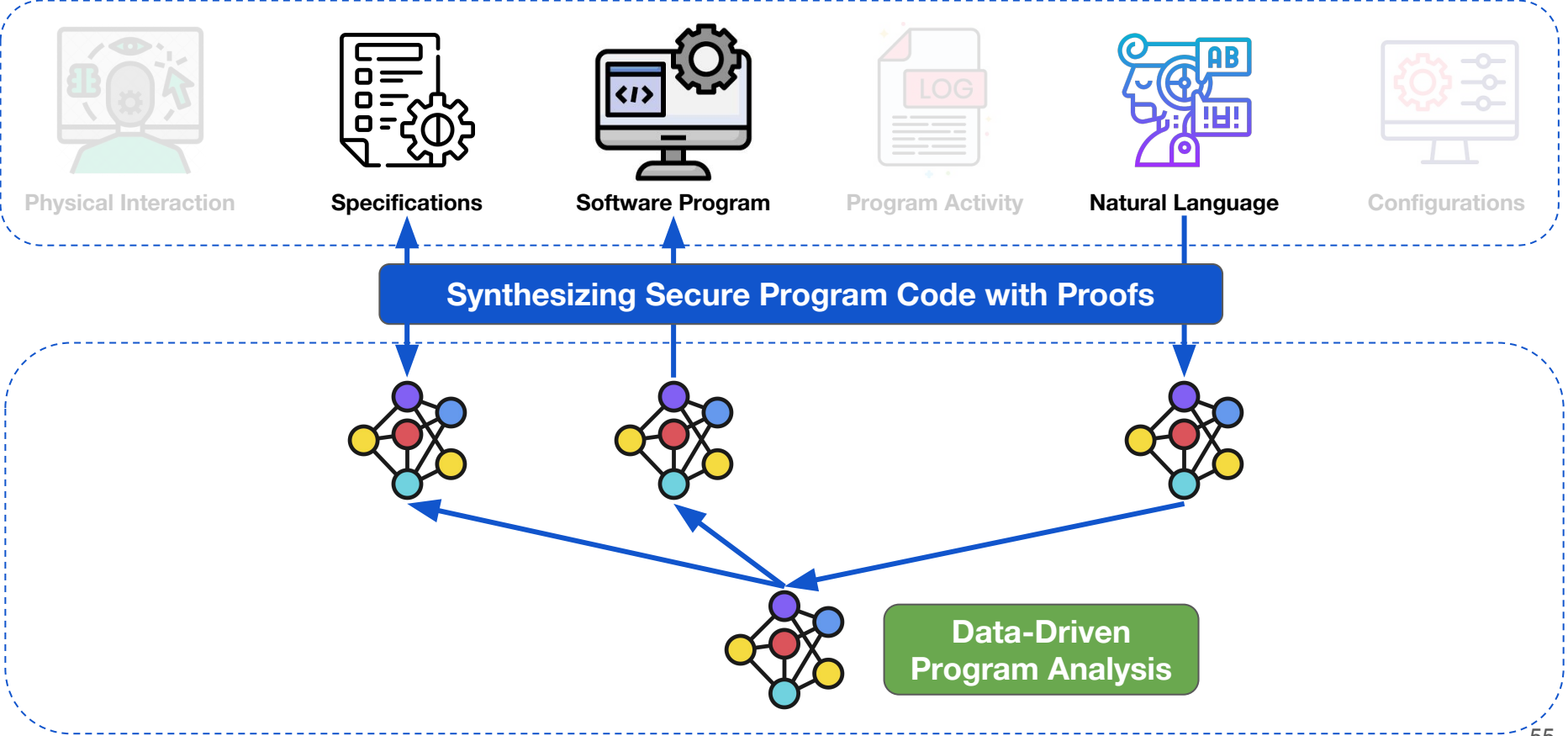
# Automating Existing Security Applications



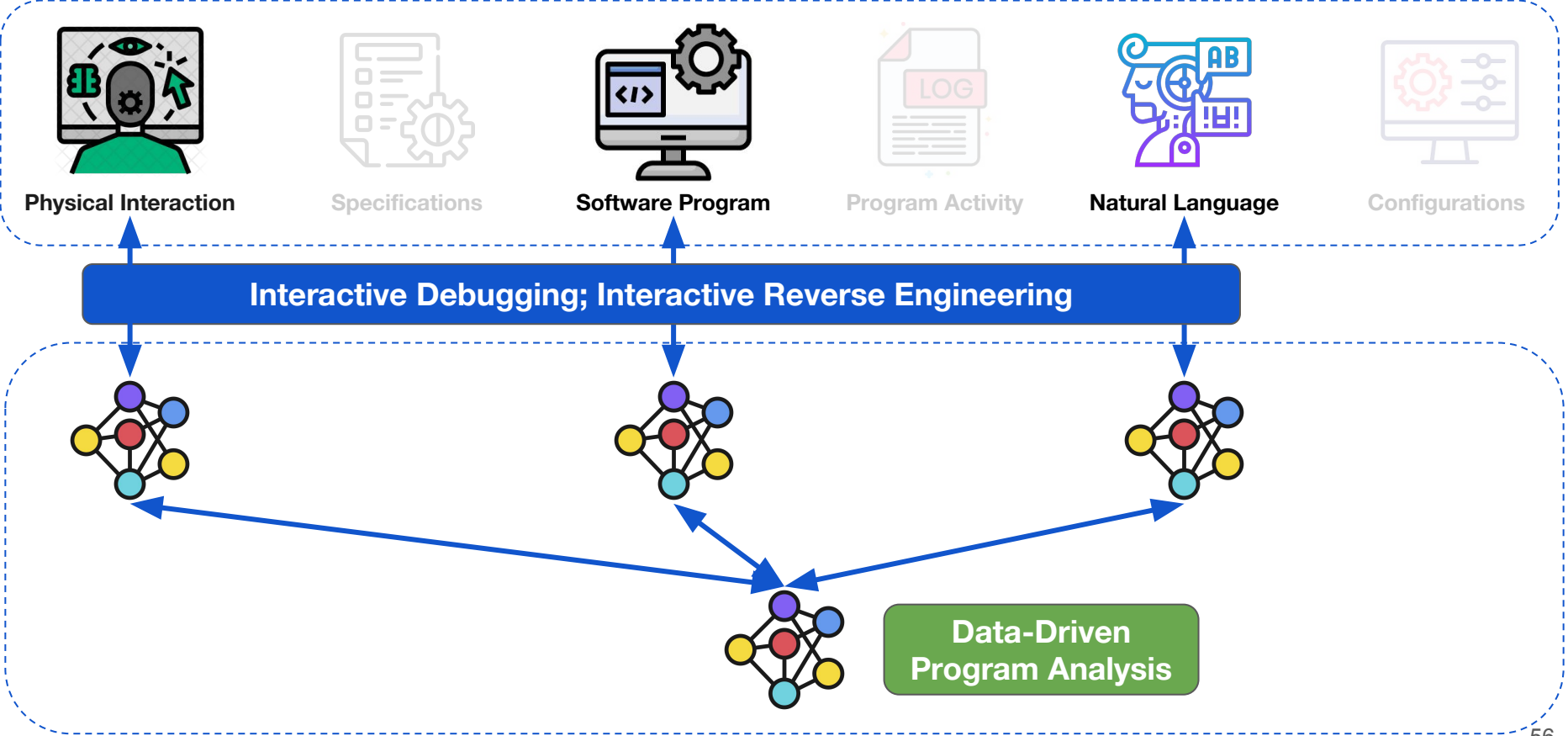
# Enabling New Security Applications



# Enabling New Security Applications



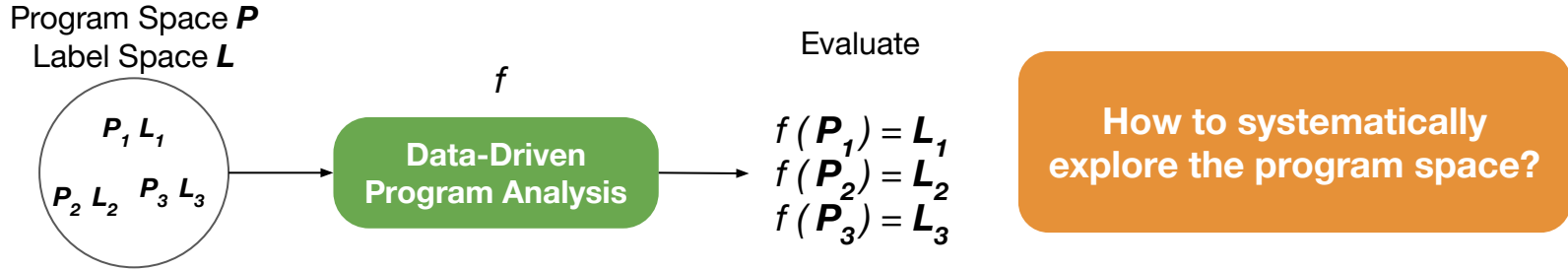
# Enabling New Security Applications



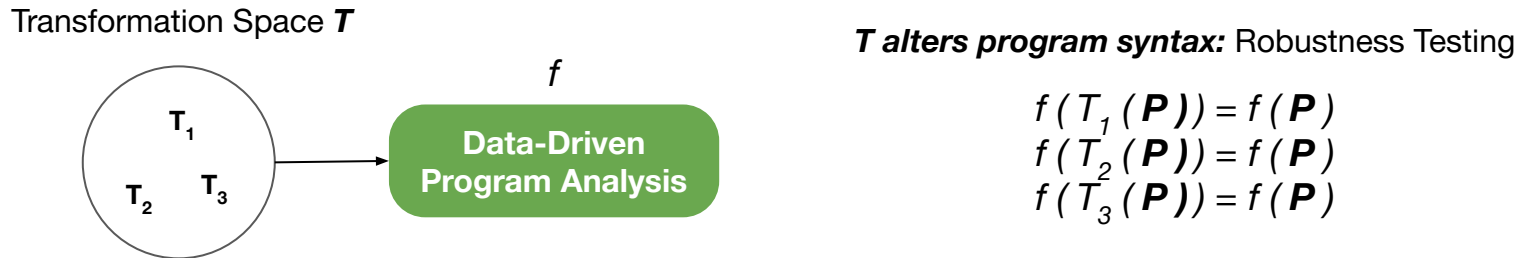


# Principled Robustness Measurement

Current testing of data-driven program analysis: **Random Testing**



Future testing of data-driven program analysis: **Transformation-Oriented Testing**



# Systematic Testing and Verification of Neural Networks

## Systematic Whitebox Testing

[SOSP'17 Best Paper Award, ICSE'18]

## Formal Verification of Security Properties

[Usenix'18, Neurips'18, DeepTest'18]

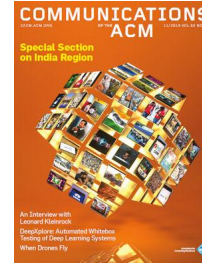


**Newsweek**



## Visual Transformations

### Data-Driven Program Analysis



### brain-research/ tensorfuzz

A library for performing coverage guided fuzzing of neural networks

🔍 0 Contributors   📄 7 Issues   ⭐ 195 Stars   🍴 56 Forks



### Translation

$$\begin{bmatrix} 1, & 0, & dx \\ 0, & 1, & dy \end{bmatrix}$$

### Scaling

$$\begin{bmatrix} s, & 0, & 0 \\ 0, & s, & 0 \end{bmatrix}$$

### Aspect ratio change

$$\begin{bmatrix} r, & 0, & 0 \\ 0, & 1, & 0 \end{bmatrix}$$

### Shear

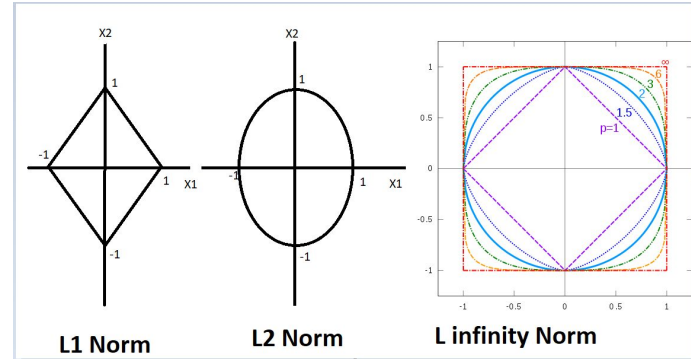
$$\begin{bmatrix} 1, & c, & 0 \\ d, & 1, & 0 \end{bmatrix}$$

### Reflect

$$\begin{bmatrix} -1(1), & 0, & 0 \\ 0, & 1(-1), & 0 \end{bmatrix}$$

### Rotation

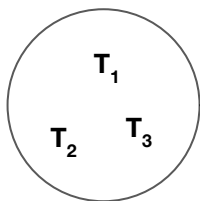
$$\begin{bmatrix} \cos(\theta), & -\sin(\theta), & 0 \\ \sin(\theta), & \cos(\theta), & 0 \end{bmatrix}$$



## Formal verification of all possible transformations

# Data-Driven Program Analysis with Provable Robustness by Construction

Transformation Space  $\mathcal{T}$



Robust  $f$  by Construction

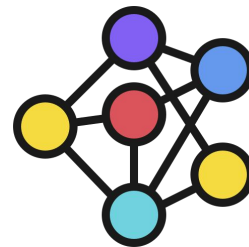
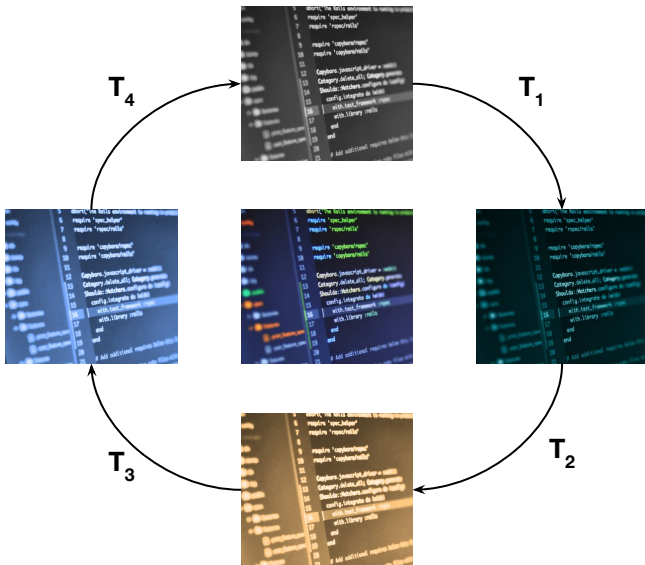


$$f(T_i(P)) = f(P) \quad \forall T_i \in \mathcal{T}$$

Permutation Group  $\mathcal{T}$

1: $x=5$	Instruction	1: $y=6$
2: $y=6$	Reordering	2: $x=5$
3: $z=x+y$		3: $z=x+y$

⋮



**Symmetry-Preserving**  
Model Architectures  
e.g., self-attention,  
graph NN

Thanks!