

Machine Learning Security

Kixin Pei

Department of Computer Science
The University of Chicago



“Software is Eating the World”

- Marc Andreessen

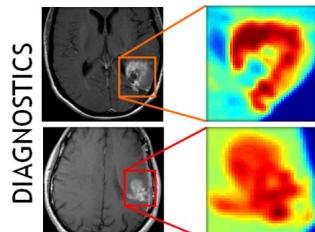
“ML is Software 2.0”

- Andrej Karpathy

ML is Increasingly Used in Security and Safety-Critical Systems



Self-Driving Cars



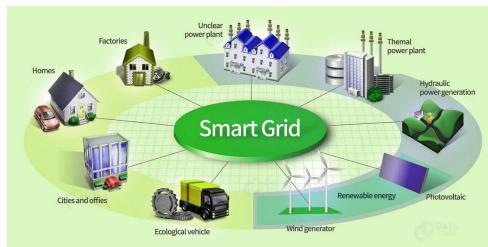
Medical Diagnosis



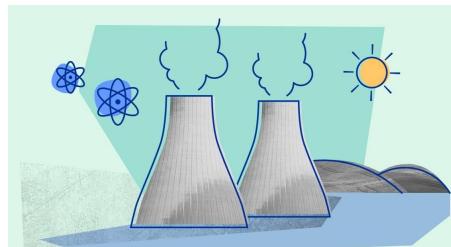
Fraud/Intrusion/Malware Detection



Flight Control Systems



Smart Grids Energy Distribution



Nuclear Energy Monitor System



Forecasting Natural Disasters

Machine Learning is Vulnerable



"panda"



"vulture"



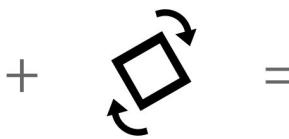
"not hotdog"

Adversarial Noise



"gibbon"

Adversarial Rotation



"orangutan"

Adversarial Photographer



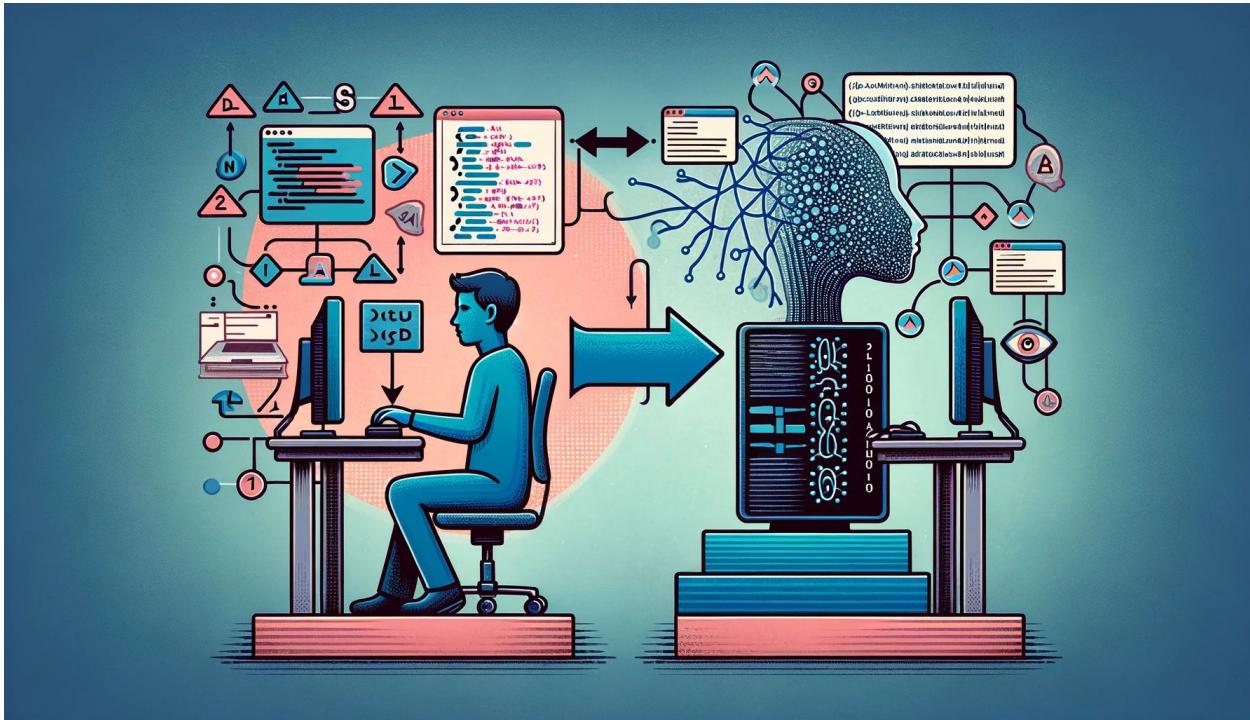
"hotdog"

Adversarial Examples

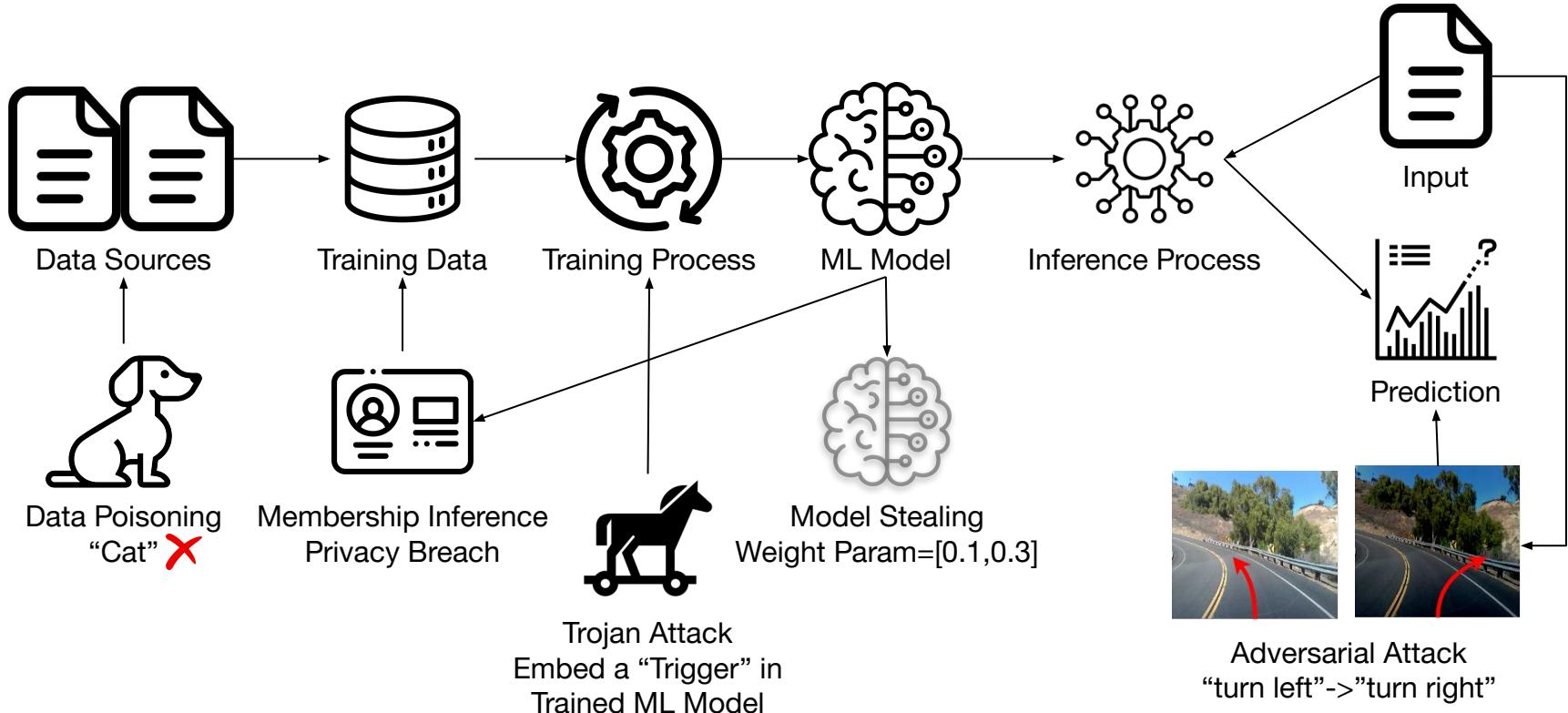


**Tesla Autopilot Failed to Recognize a White Truck
against Bright Sky Leading to Fatal Crash**

Paradigm Shift of Software Development



Threat Landscape Exposed by ML



ML Security in a Nutshell

Security of Machine Learning



Machine Learning for Security?

Machine Learning has Automated Many Security Applications

[Research](#) [Threat intelligence](#) [Vulnerabilities and exploits](#) · 4 min read

Microsoft researchers work with Intel Labs to explore new deep learning approaches for malware classification

By [Microsoft Threat Intelligence](#)

Learn / Microsoft Defender for Cloud Apps /

Anomaly detection policies in Defender for Cloud Apps

environment and triggers alerts with respect to a baseline that was learned on your organization's activity. These detections also use [machine-learning algorithms](#) designed to profile the users and sign in pattern to reduce false positives.

09 MARCH 2021 [TECH TOPICS](#)

Detecting threats in AWS Cloudtrail logs using machine learning

By [Craig Chamberlain](#)

Share



A Smarter Way to Detect Suspicious Cloud Logins

10/25/2021 | [ActZero](#) | 4 minutes

How Proofpoint Aegis Uses Machine Learning

SHARE WITH YOUR NETWORK!

FEBRUARY 21, 2023 | ADAM STARR AND LUIS BLANDO



The advent of transformer-based natural language processing (NLP) models, large language models (LLMs) and generative models has made advanced machine learning (ML) systems accessible to everyone with internet access through tools like ChatGPT. Unfortunately, malicious actors can also use these models to create [phishing emails](#).

In this post, we'll explain how Proofpoint defends against [artificial intelligence \(AI\)](#)-generated phishing threats and deploys advanced ML models throughout our detection systems.

Defending against AI-generated phishing threats

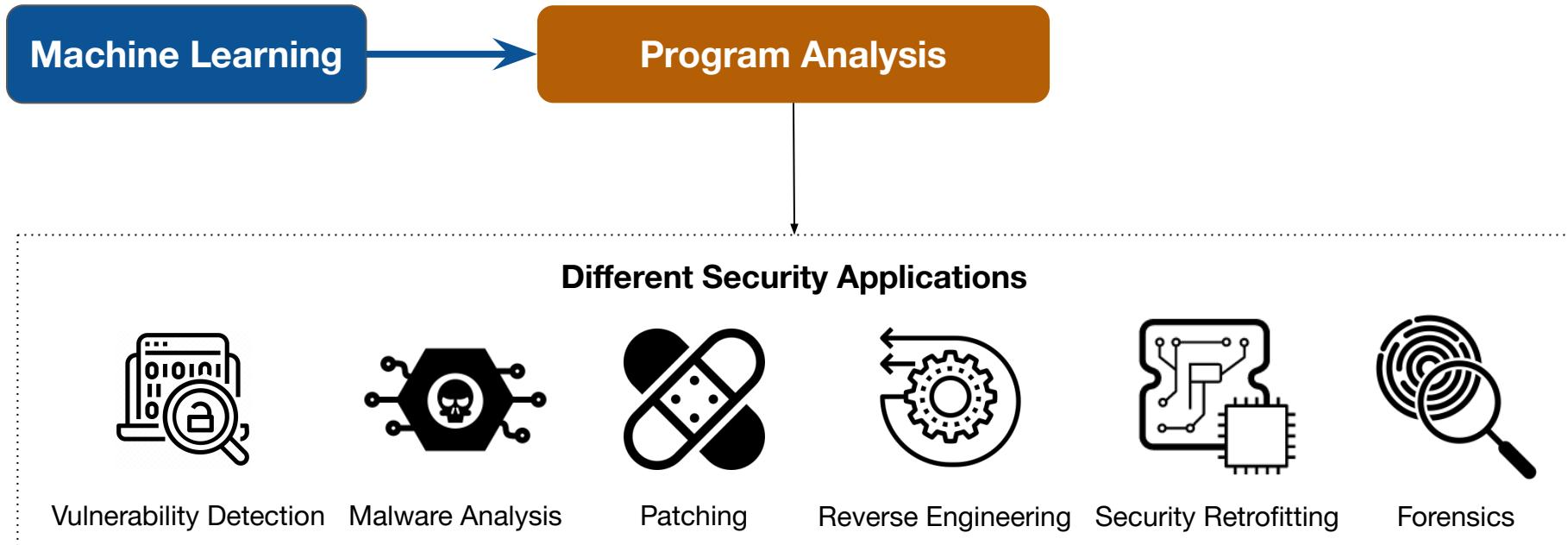


Engineering Security

Leveraging machine learning to find security vulnerabilities

A behind-the-scenes peek into the machine learning framework powering new code scanning security alerts.

My Research: Machine Learning for Software Security



Program Analysis is Crucial for Building Trustworthy Software



Software Programs

Questions



Program Analysis

Can pointer p be
NULL at line L?



Answers

Help to Build

Vulnerability



Privacy Leak



Trustworthy Software

Security

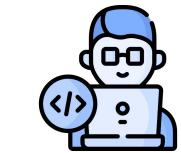
Reliability

Safety

Privacy

Performance

Challenges of Traditional Program Analysis



Hand-Curate



How to

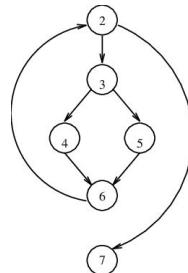
Human Expert

Rules and Heuristics



Tune

Significant Manual Effort



Represent a program?

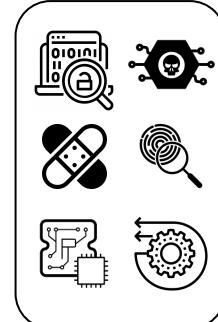
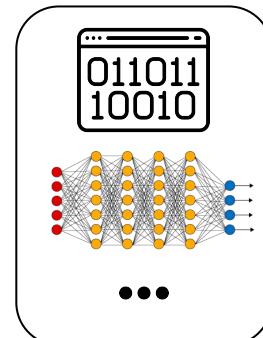
Input: $kill(B)$ and $gen(B)$ for every basic bloc B .

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

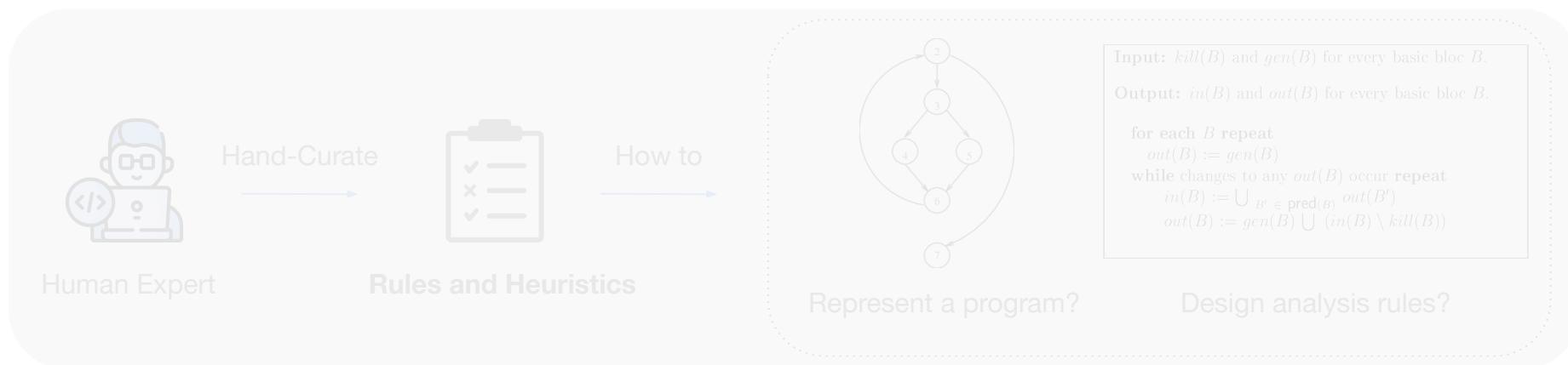
```
for each  $B$  repeat  
     $out(B) := gen(B)$   
    while changes to any  $out(B)$  occur repeat  
         $in(B) := \bigcup_{B' \in pred(B)} out(B')$   
         $out(B) := gen(B) \cup (in(B) \setminus kill(B))$ 
```

Design analysis rules?

Heterogeneous Software Various Security Applications



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

omg. seriously mind-blown 😱



Brandon Rossi
@0xConda

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);  
    }  
}
```



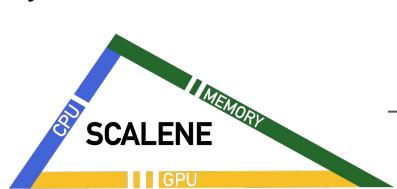
This contract appears to have a potential vulnerability in the `withdraw` function. The `withdraw` function allows anyone to specify the `from` address when calling the function, which means that a malicious user could potentially withdraw funds from an arbitrary address by calling `withdraw` and specifying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



GPT-4

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

EXPLAIN

```
const donut = Number(process.argv[2])  
  
if ((donut % 3 === 0) && (donut % 5 === 0))  
{  
    console.log('FizzBuzz')  
}
```

Explain code

Advanced

Ask Copilot

RESULT

```
/* Here is the explanation for the code above:  
1. We are taking the argument from the command line and converting it to a number.  
2. We then use an if statement to check if the number is divisible by 3 and 5.  
3. If the number is divisible by 3 and 5, we output FizzBuzz.  
4. If the number is not divisible by 3 and 5, we output the number. */
```

LANGUAGE TRANSLATION

```
const donut = Number(process.argv[2])  
  
if ((donut % 3 === 0) && (donut % 5 === 0))  
{  
    console.log('FizzBuzz')  
}
```

Translate code into:

python

Ask Copilot

Translate Code



3% Code Written by ML

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

**Huge Academic Contributions:
1,000+ 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 / 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%)

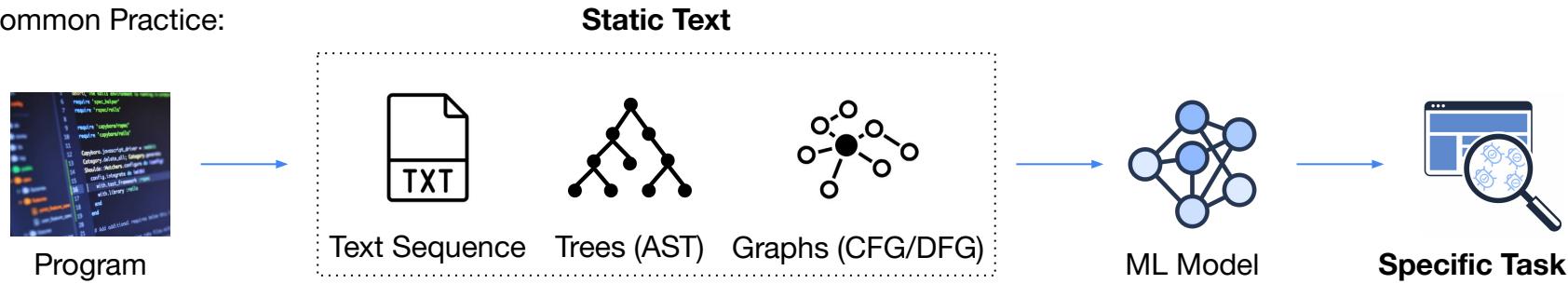
```
void f2(int[] ttypes){  
    boolean swapped = true;  
    for (int i = 0;  
        i < ttypes.length && swapped; i++) {  
        swapped = false;  
        for (int j = 0;  
            j < ttypes.length-1-i; j++) {  
            if (ttypes[j] > ttypes[j+1]) {  
                int temp = ttypes[j];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



Prediction: **contains** (99.97%)

Common Practice of ML on Code

Common Practice:



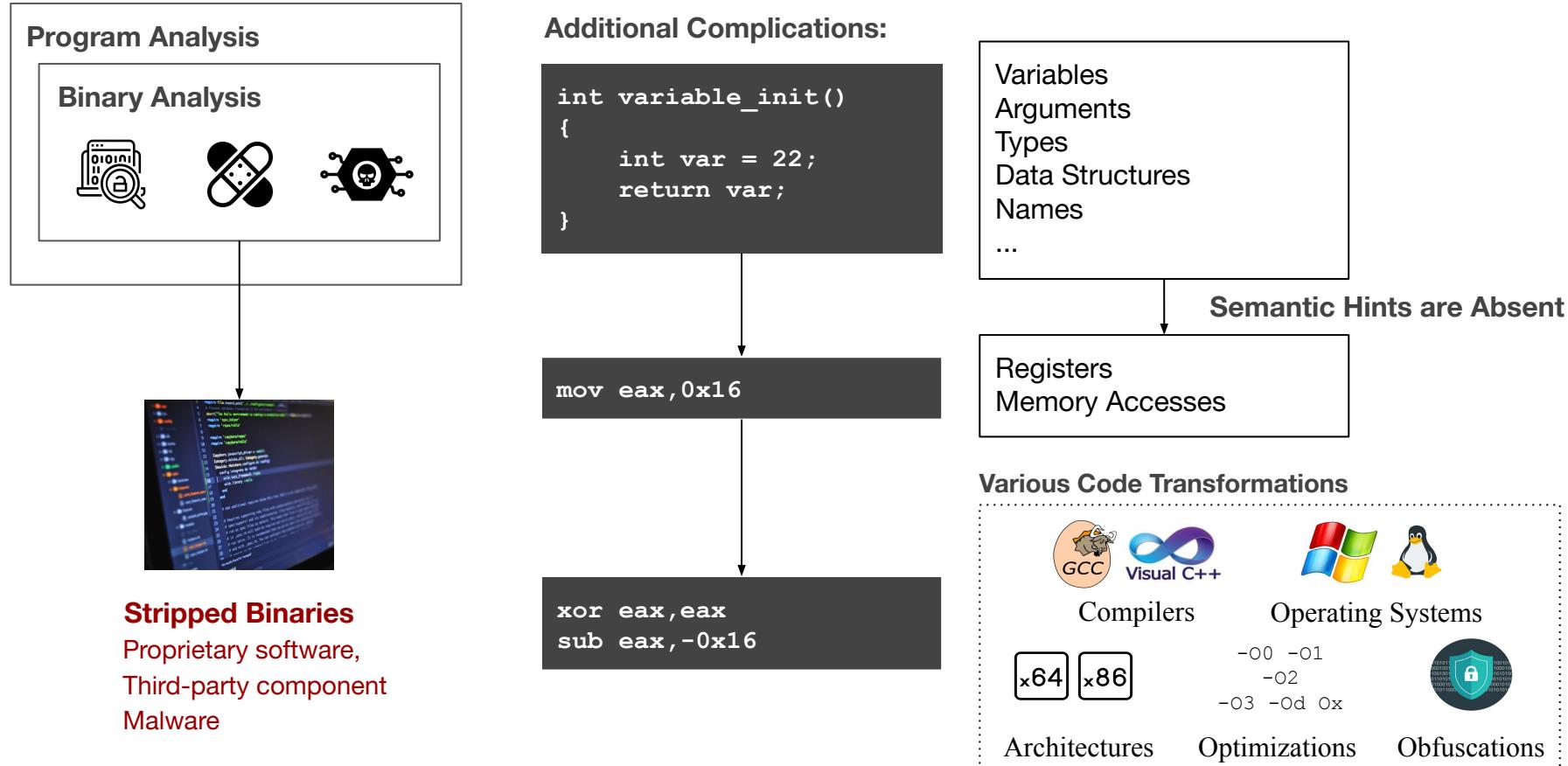
Program semantics does not just manifest in **static text**

Consequences: Lacking Robustness and Generalization

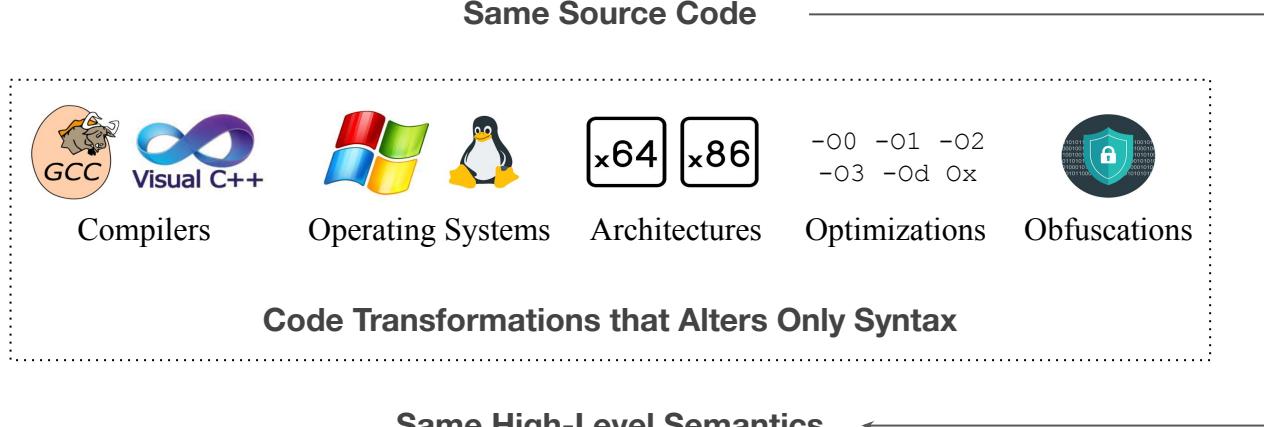
- I. Overfit to spurious textual and task-specific patterns
- II. Distribution shift: program syntax and task requirement changes

Security Applications Require More Rigorous Understanding of Program Semantics

A Popular Type of Program Analysis in Security: Binary Analysis



What do Robustness and Generalization Imply in Binary Program Analysis?

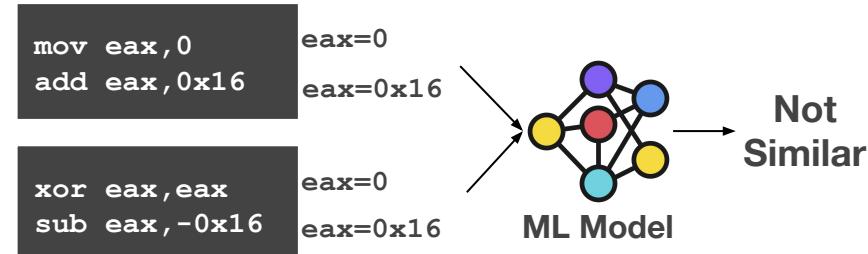
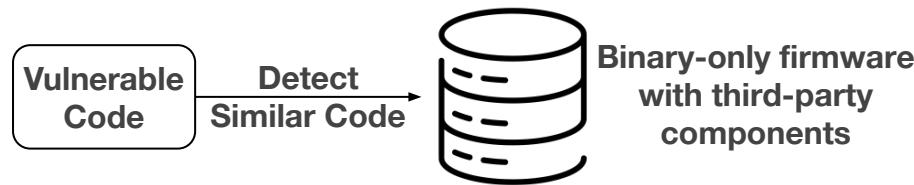


Robustness: Stay Invariant to **Syntactic** Changes

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



mov eax, 0
add eax, 0x16



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

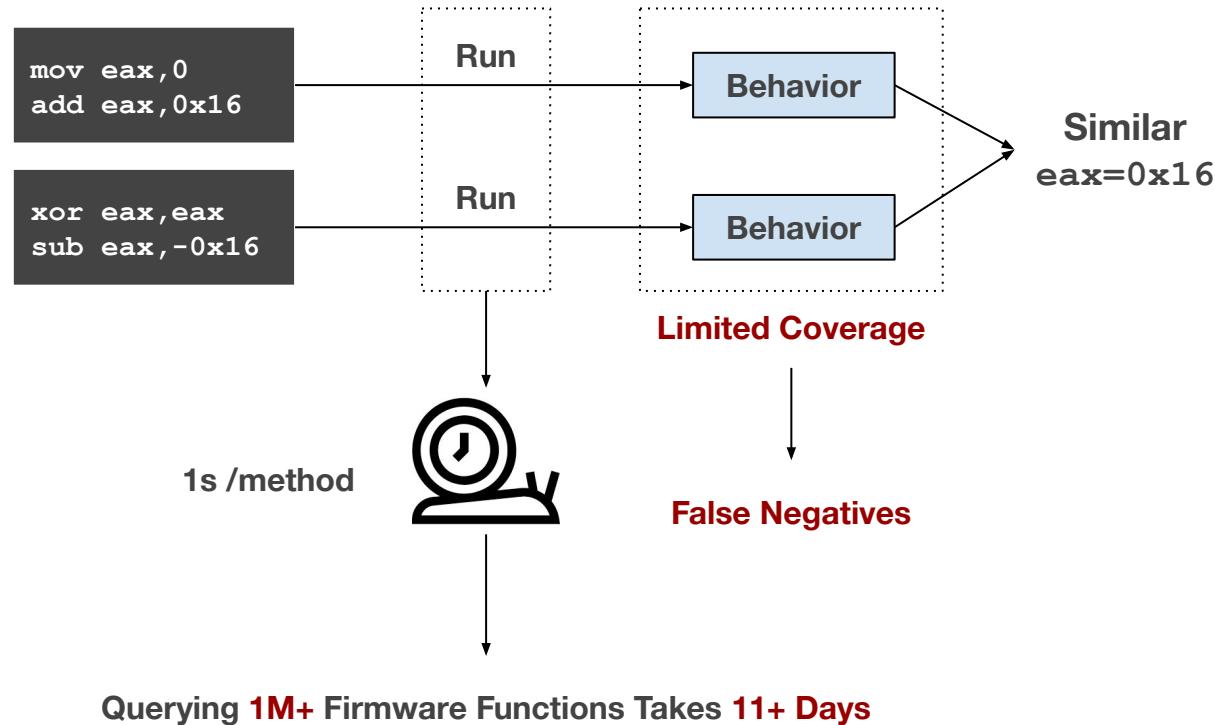


is it similar to
xor eax, eax
sub eax, -0x16?

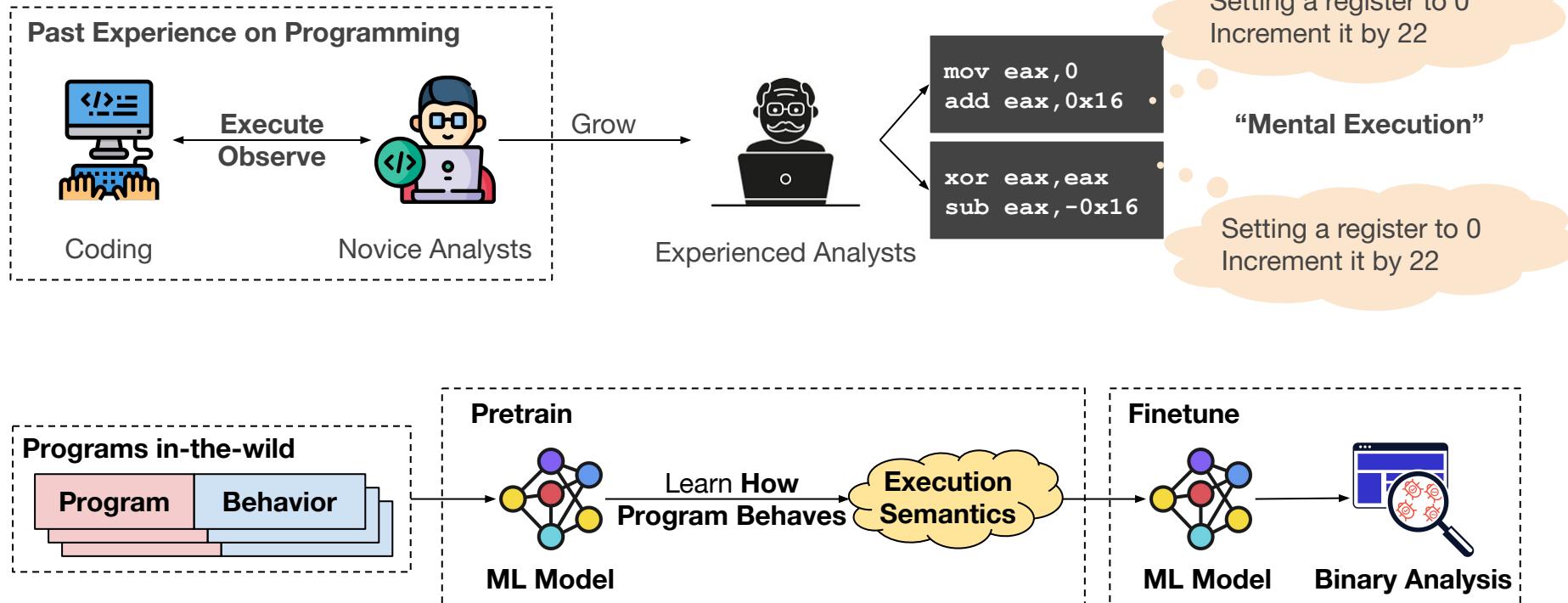


No...the second set uses the "xor" and "sub" instructions to set the value of **EAX** register to **-22**...

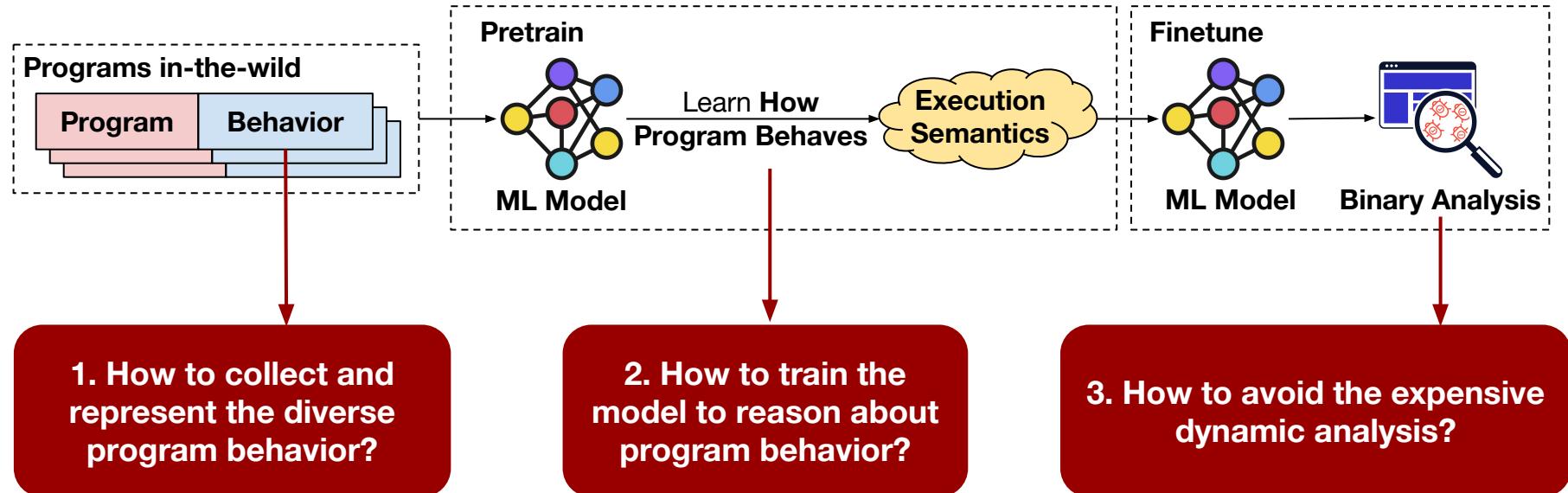
Why not dynamic analysis?



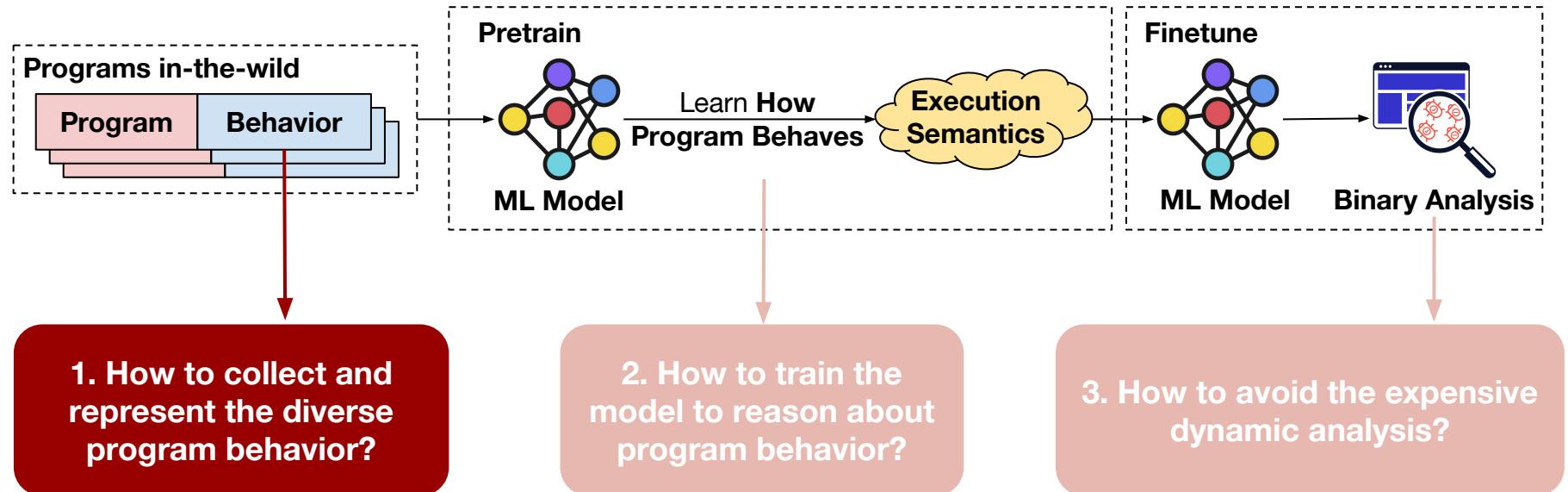
Learning Execution Semantics and Transferring it without Dynamic Analysis



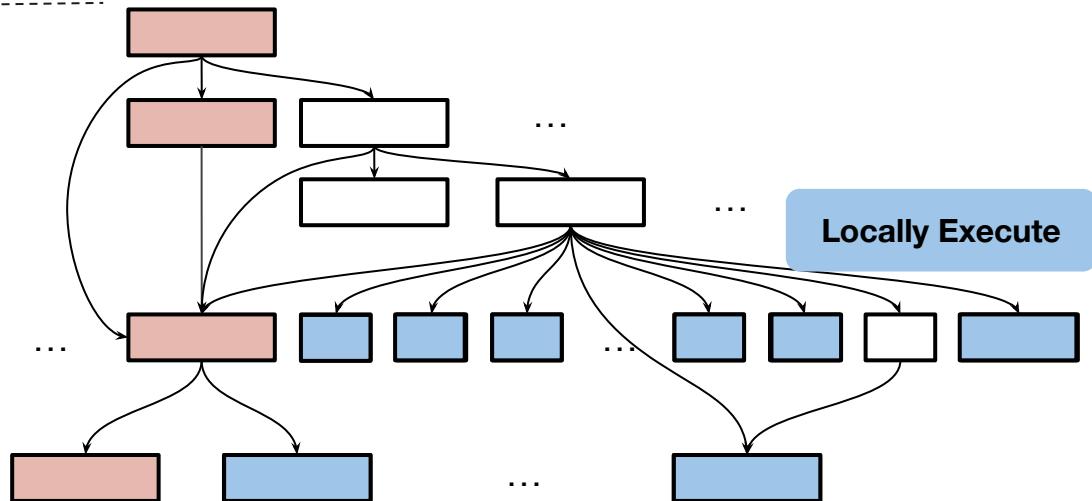
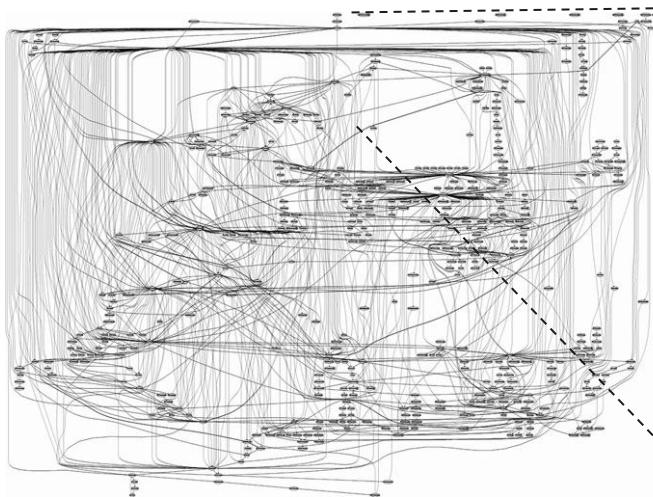
Challenges of Learning Execution Semantics



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

```
.....  
## 0xc,0x4  
## [0xc+0x8],0x3  
## [0xc+0x8],0x2  
## 0x3a  
## [0xc+0x8],0x1  
## 0x1,0x1a  
.....
```

Control flow states

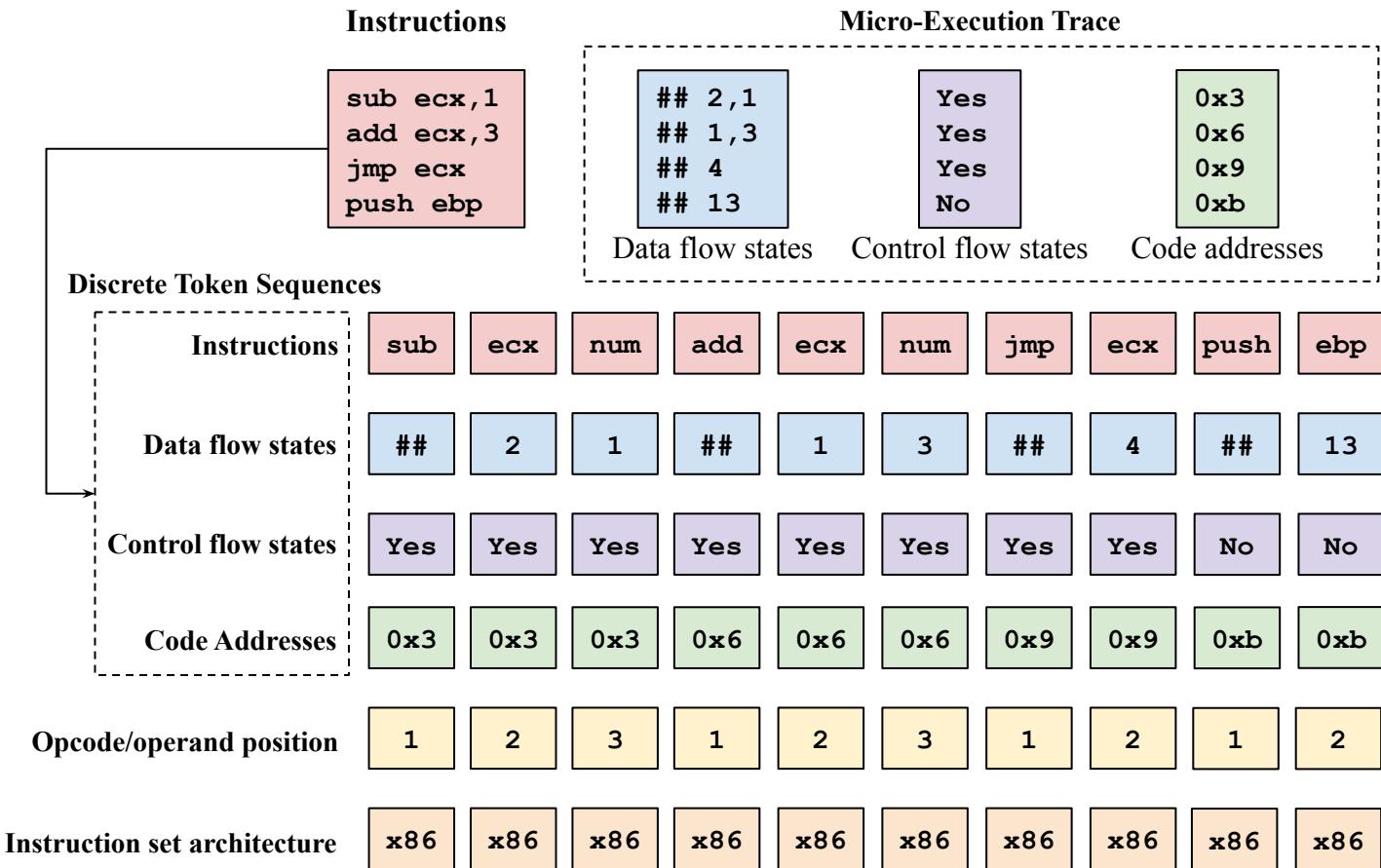
```
.....  
✓ Yes  
✓ Yes  
✓ Yes  
✓ Yes  
✗ No  
✗ No  
.....
```

Code Addresses

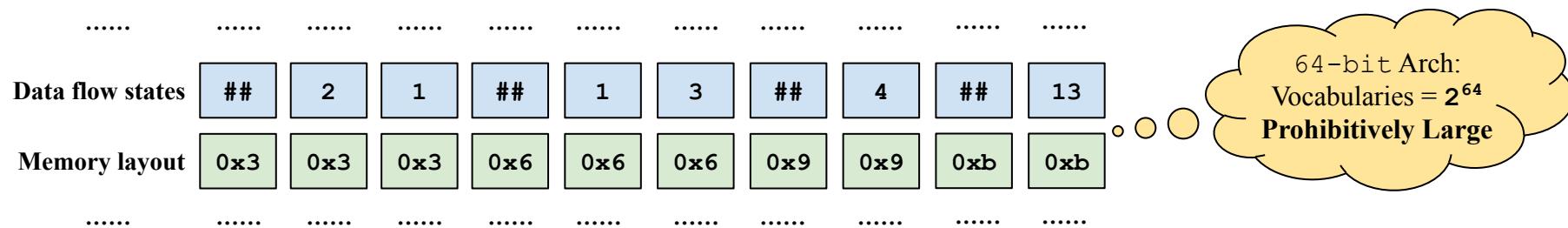
```
.....  
0x1c  
0x1f  
0x26  
0x2d  
0x33  
0x3a  
.....
```

Aligned with Program Instructions

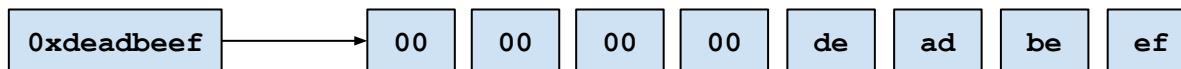
Input Representation



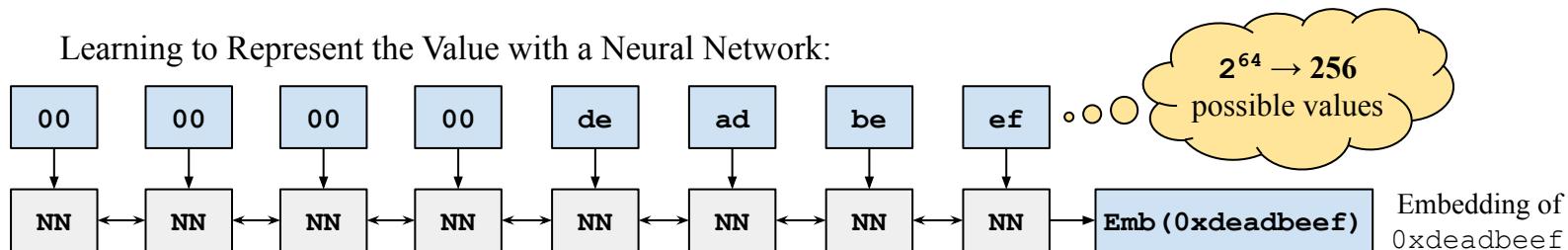
Numerical Representation



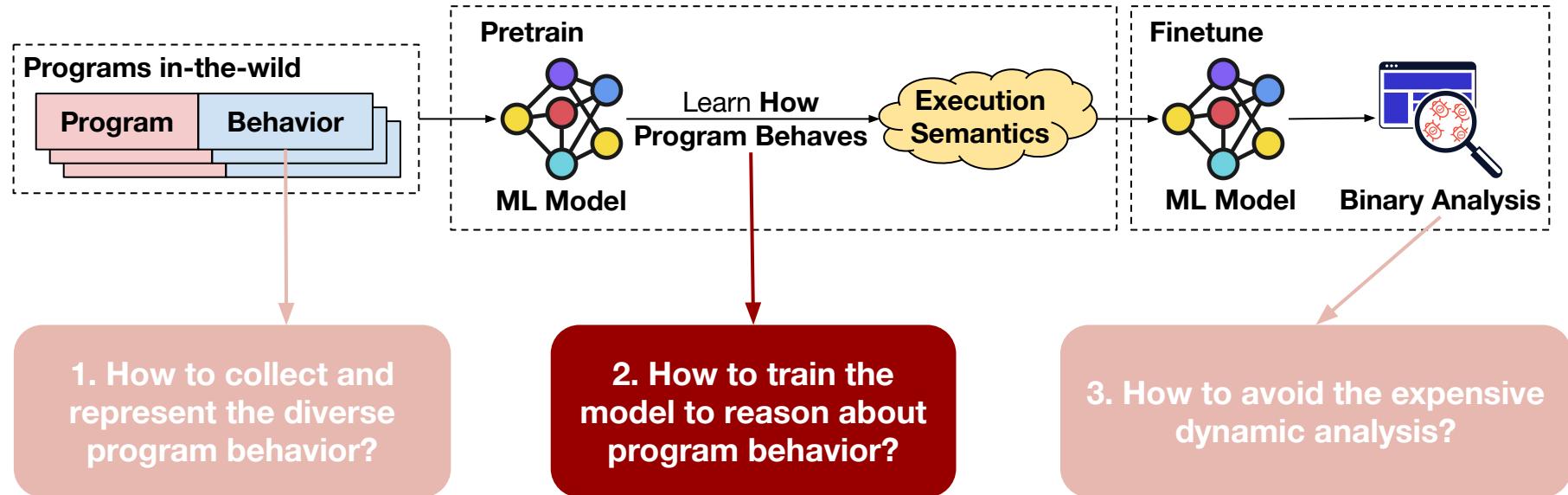
Pad Each Numeric Token as a Fixed-Length 8-Byte Sequence:



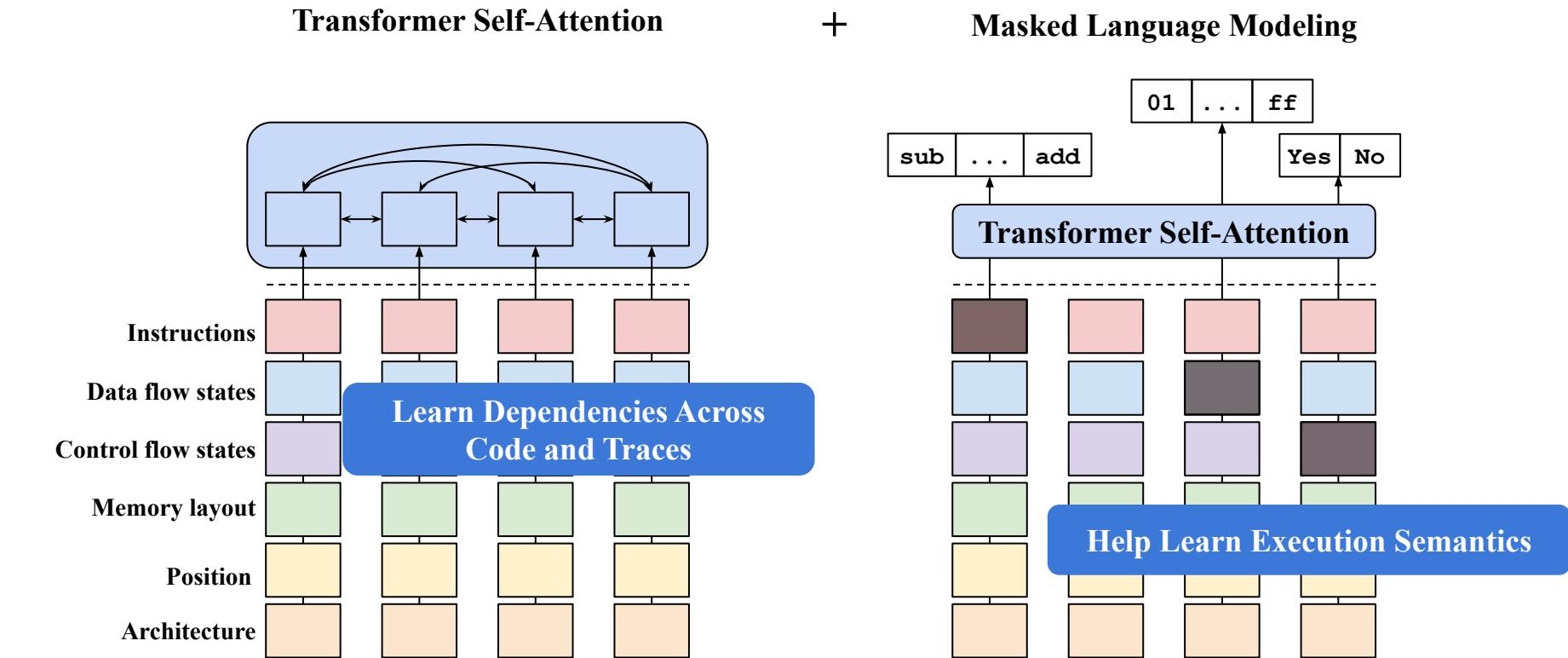
Learning to Represent the Value with a Neural Network:



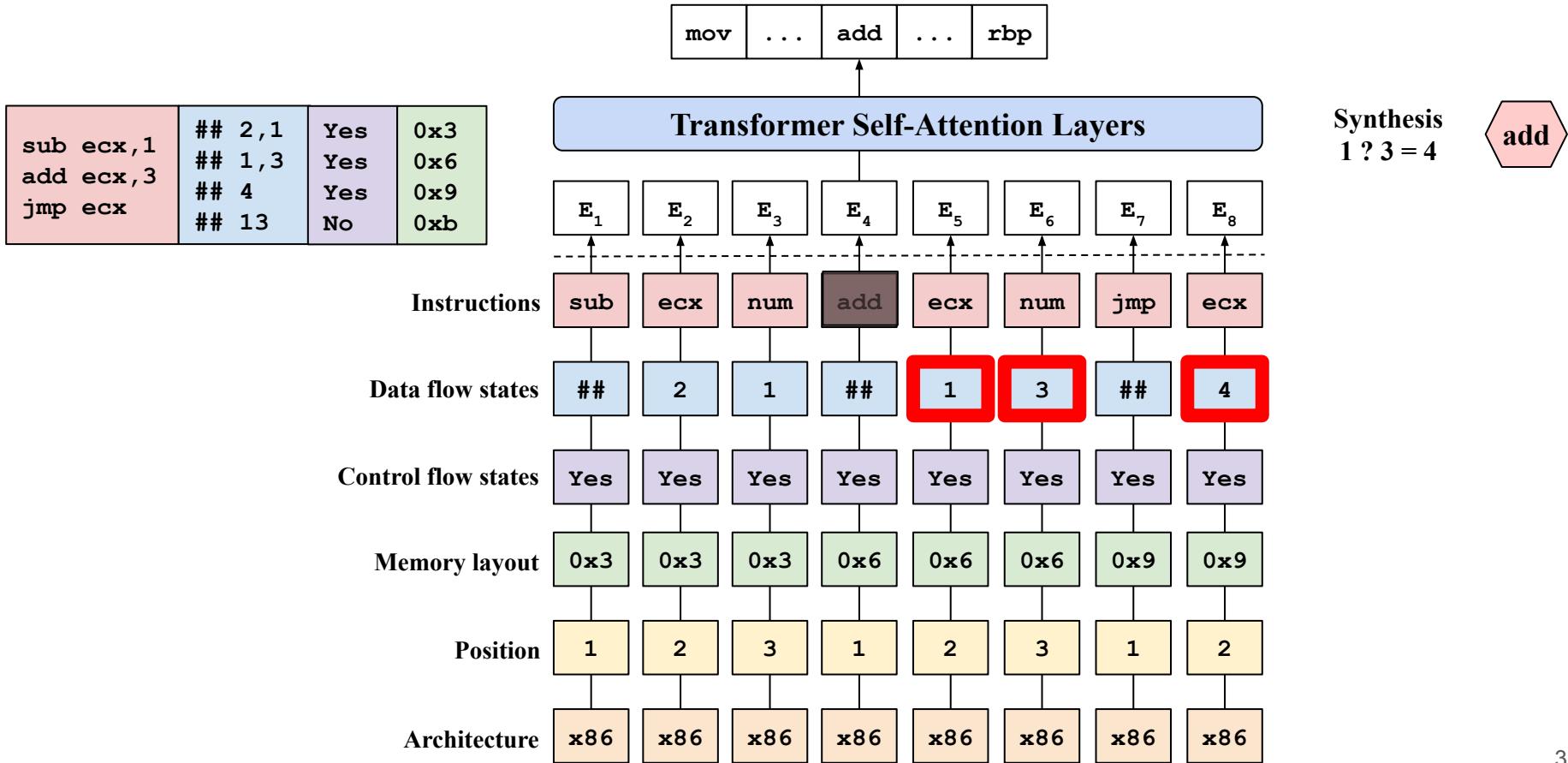
Challenges of Learning Execution Semantics



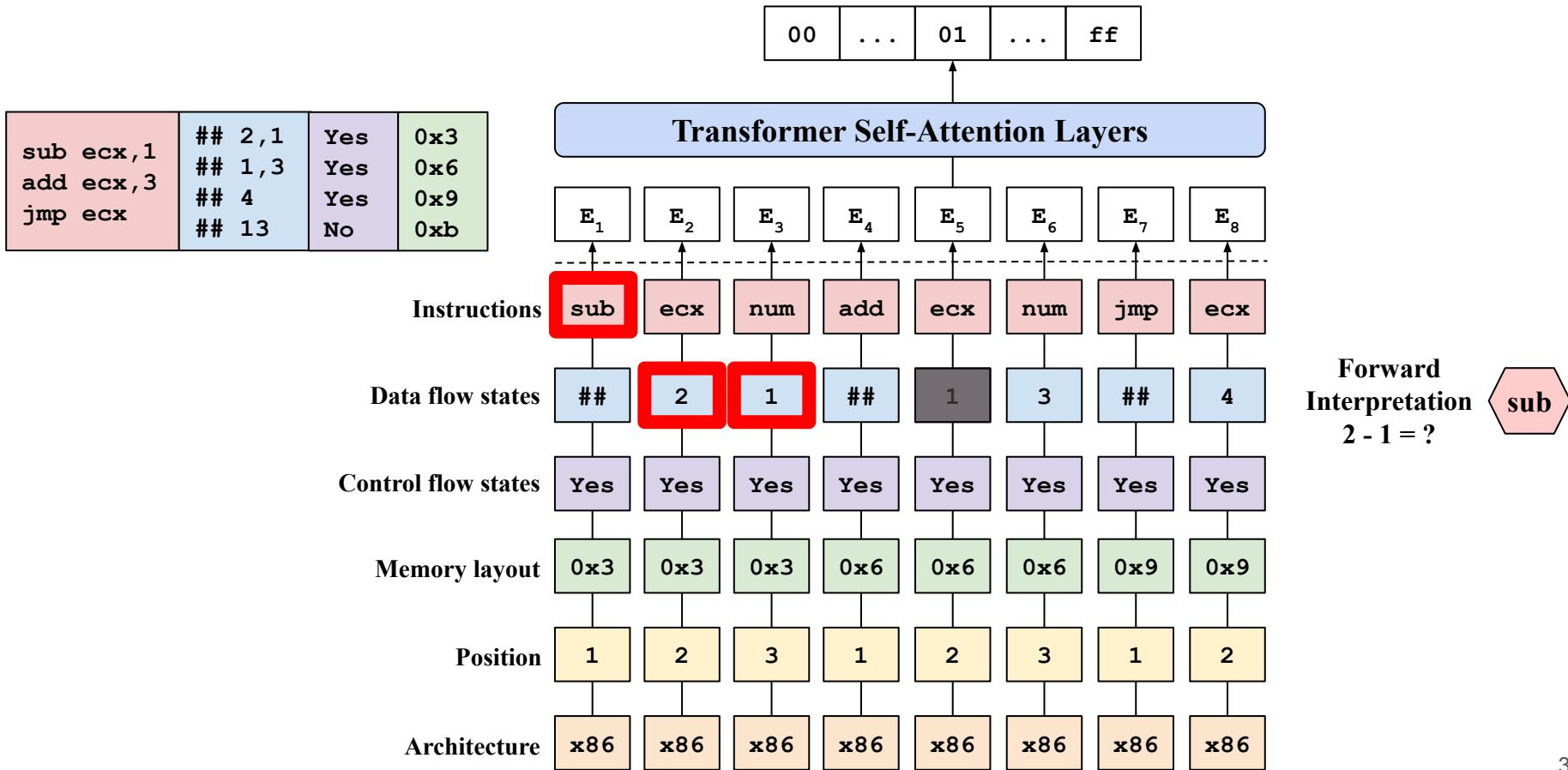
How to train the model to reason about program behavior?



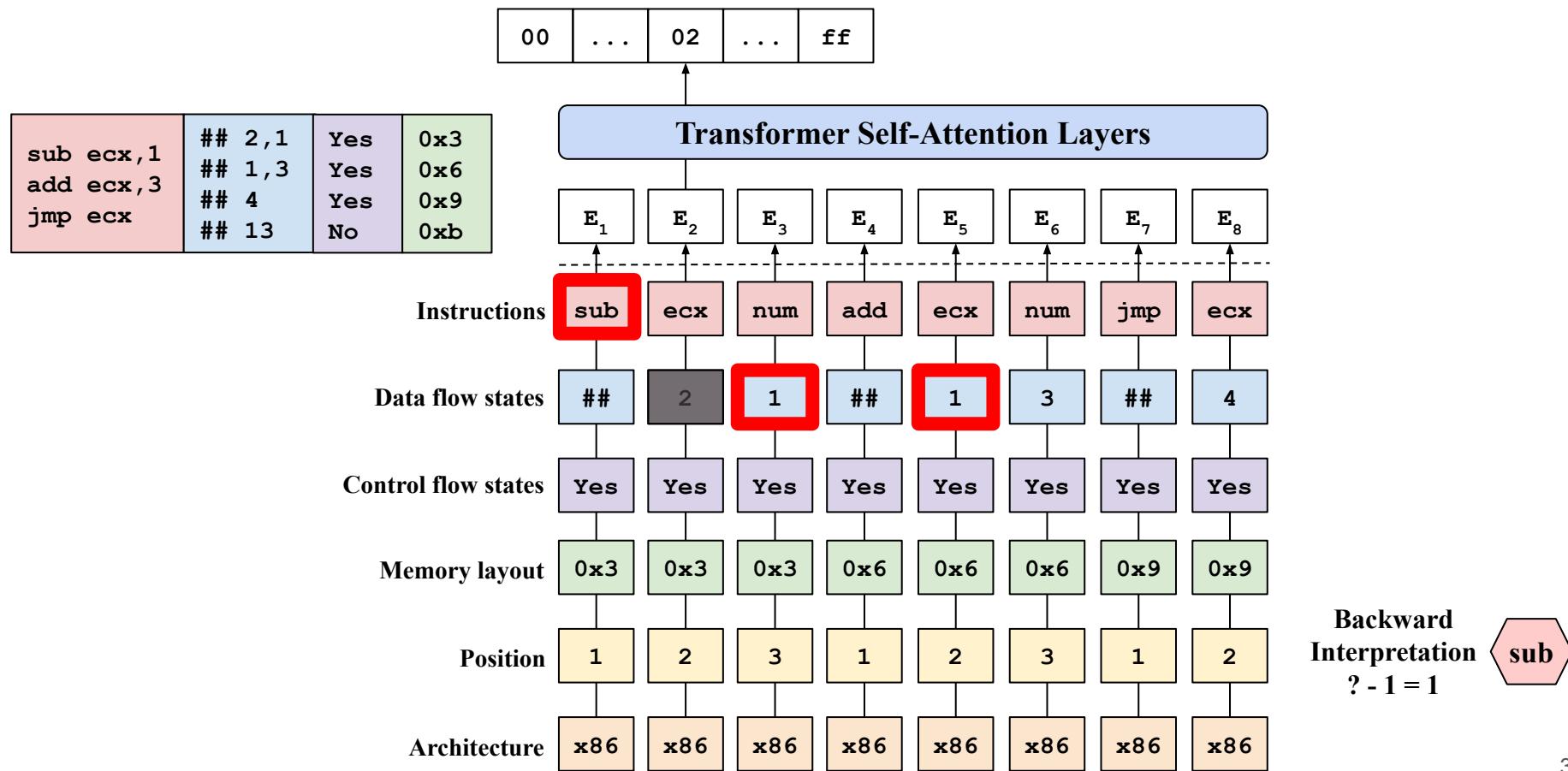
Motivating Example



Motivating Example



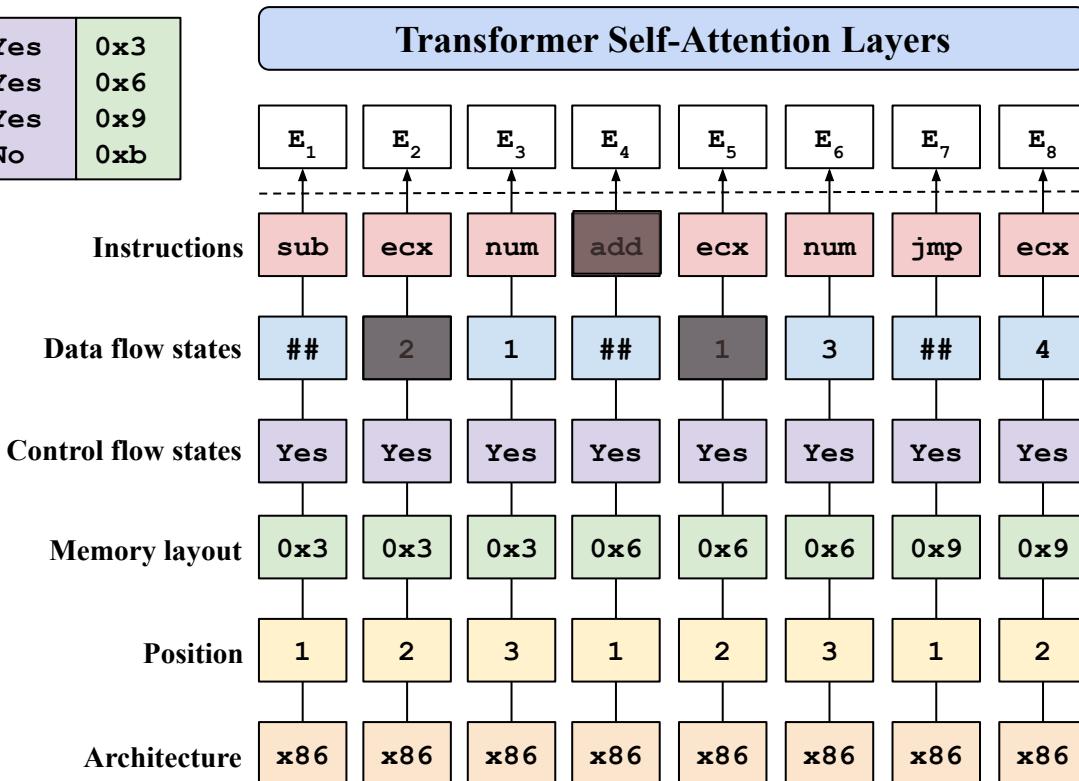
Motivating Example



Motivating Example

- Program Interpretation
- Program Synthesis

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



Synthesis
1 ? 3 = 4



Forward Interpretation
2 - 1 = ?



Backward Interpretation
? - 1 = 1

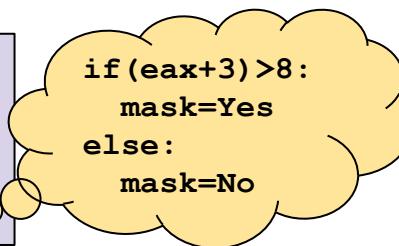


How to train the model to reason about program behavior?

How to train the model to reason about
control flow?

Mask control flow states

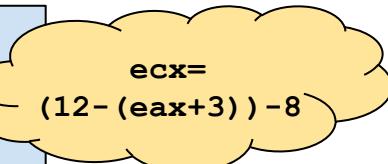
add eax, 3	Yes
cmp eax, 8	Yes
jle 0x8	Yes
mov ecx, ebx	No



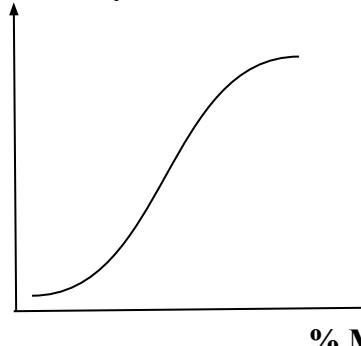
How to train the model to reason about
instruction compositions?

Mask more

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



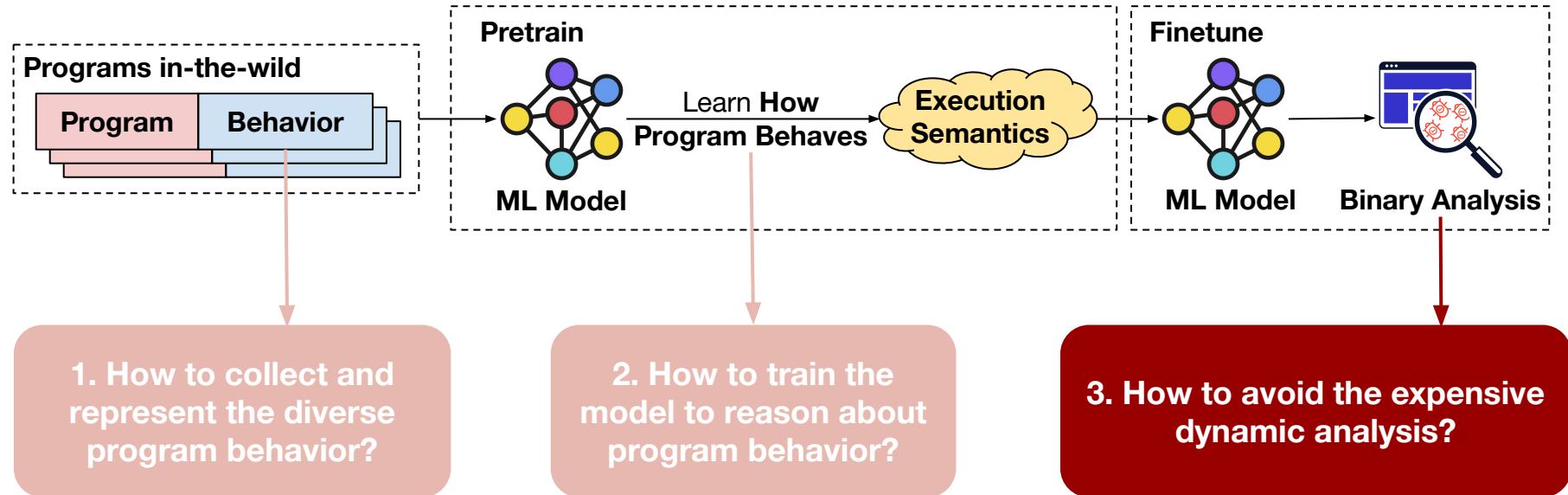
Compositionality



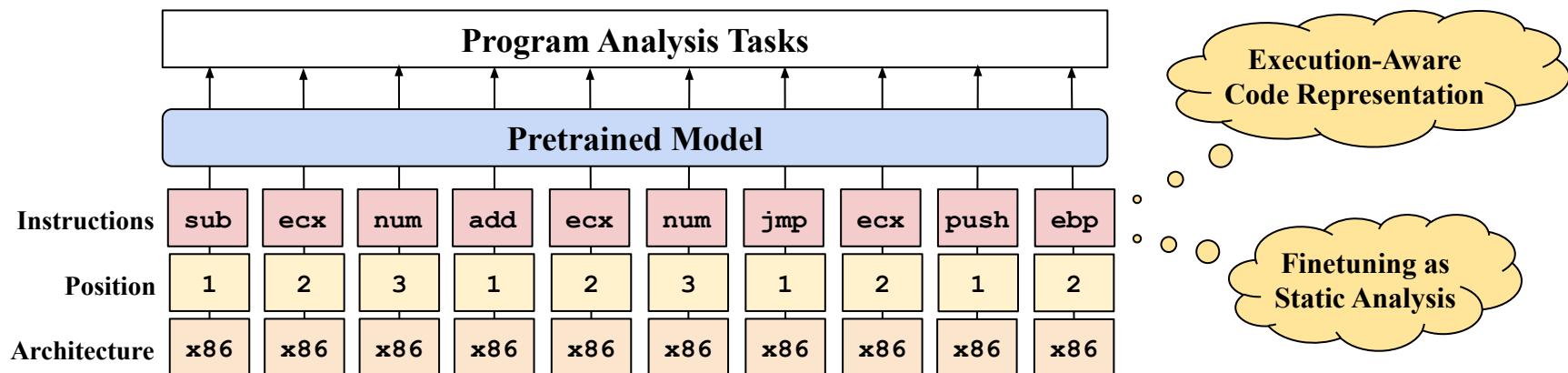
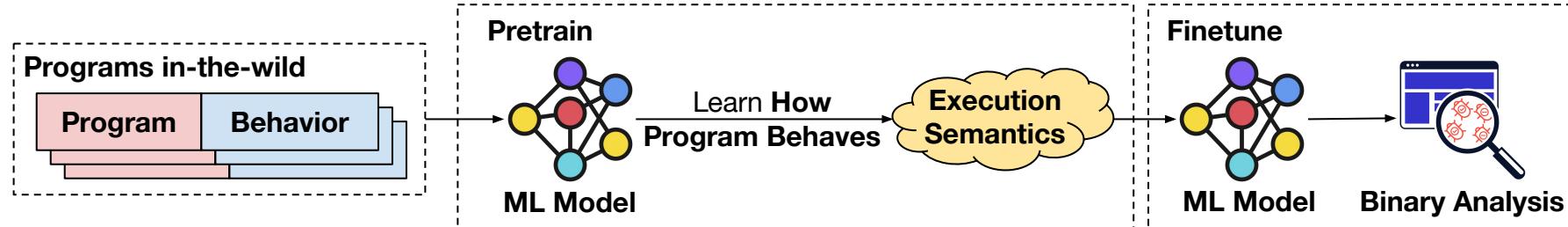
Training

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

Challenges of Learning Execution Semantics

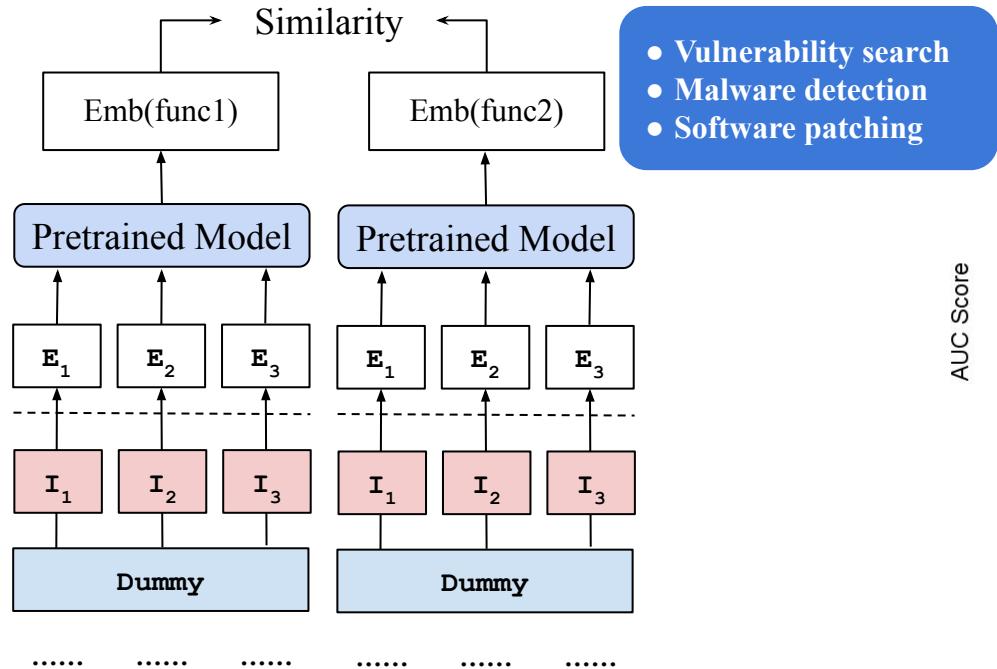


How to Avoid the Expensive Dynamic Analysis?



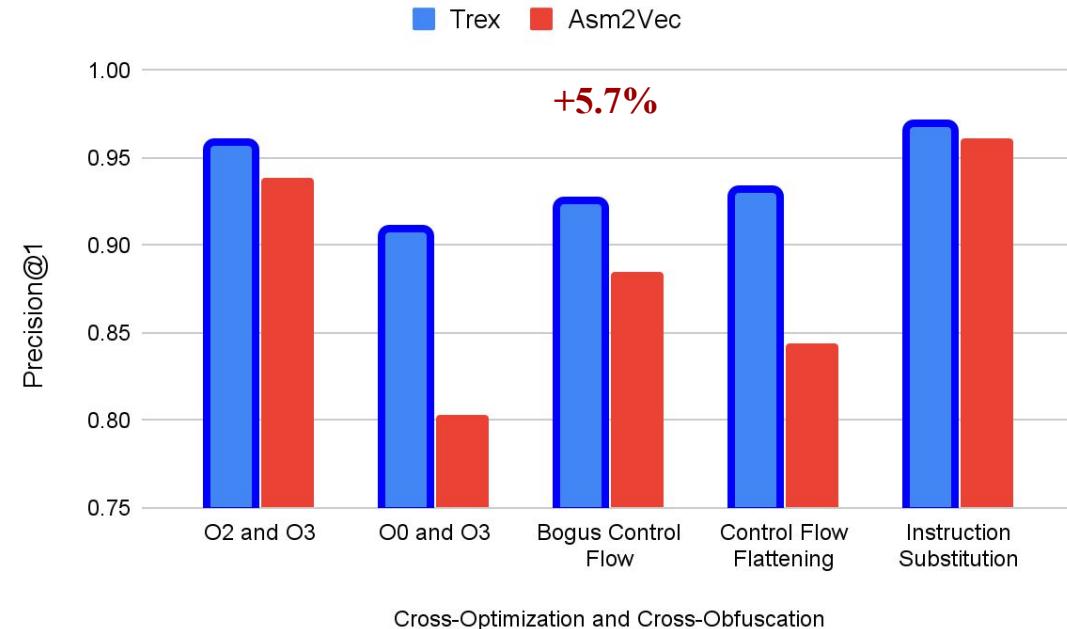
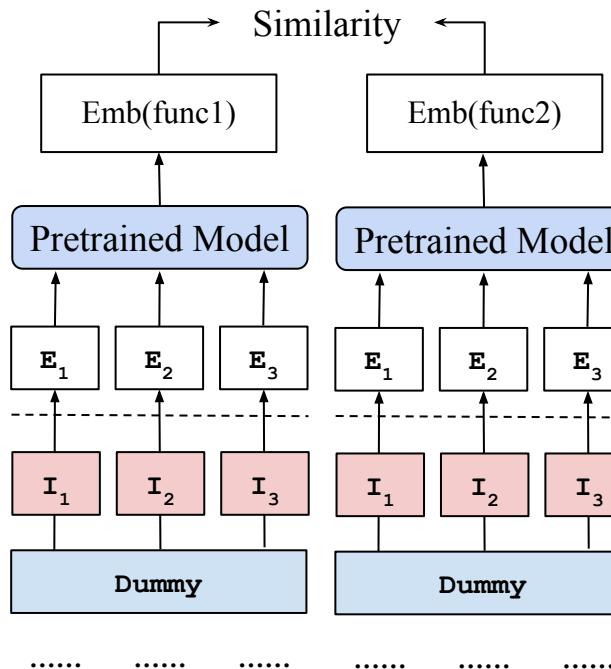
How Much do Learned Execution-Aware Program Representations Help?

Task 1: Binary Similarity



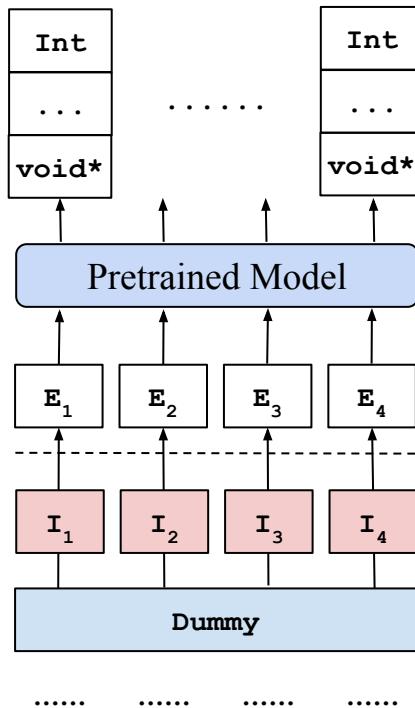
Finetuning for Matching Semantically Similar Binary Functions

Task 1: Binary Similarity



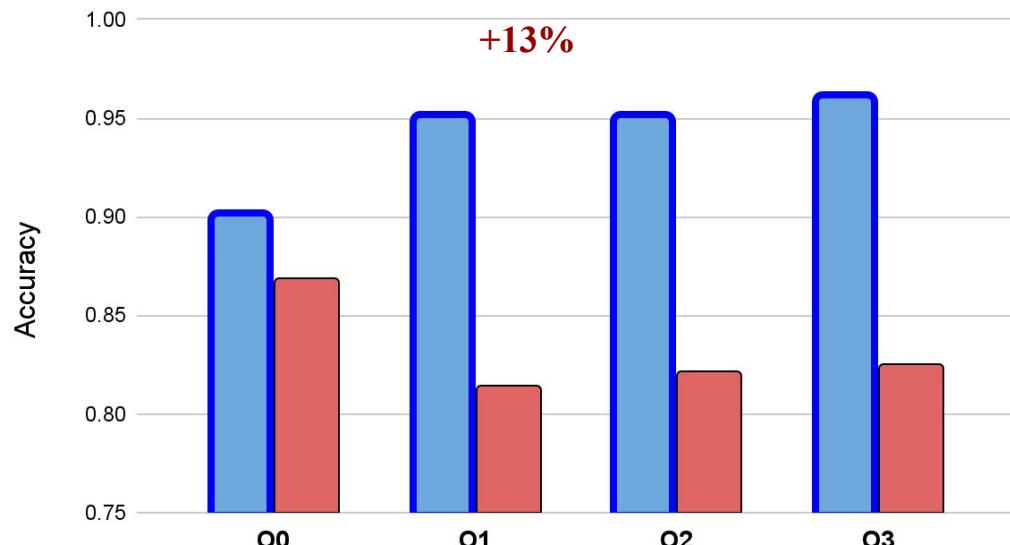
Finetuning for Predicting Function Signatures and Type Inference

Task 2: Type Inference



- Security retrofitting
- Decompilation
- Vulnerability detection

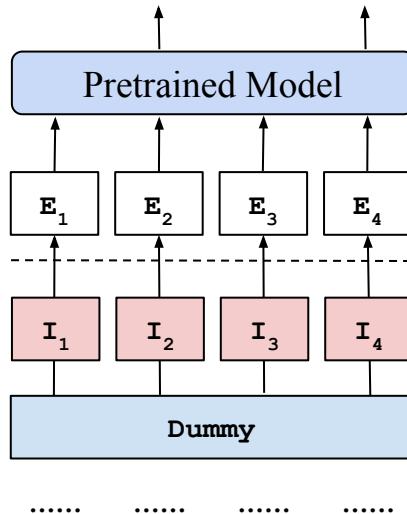
StateFormer EKLAZYA



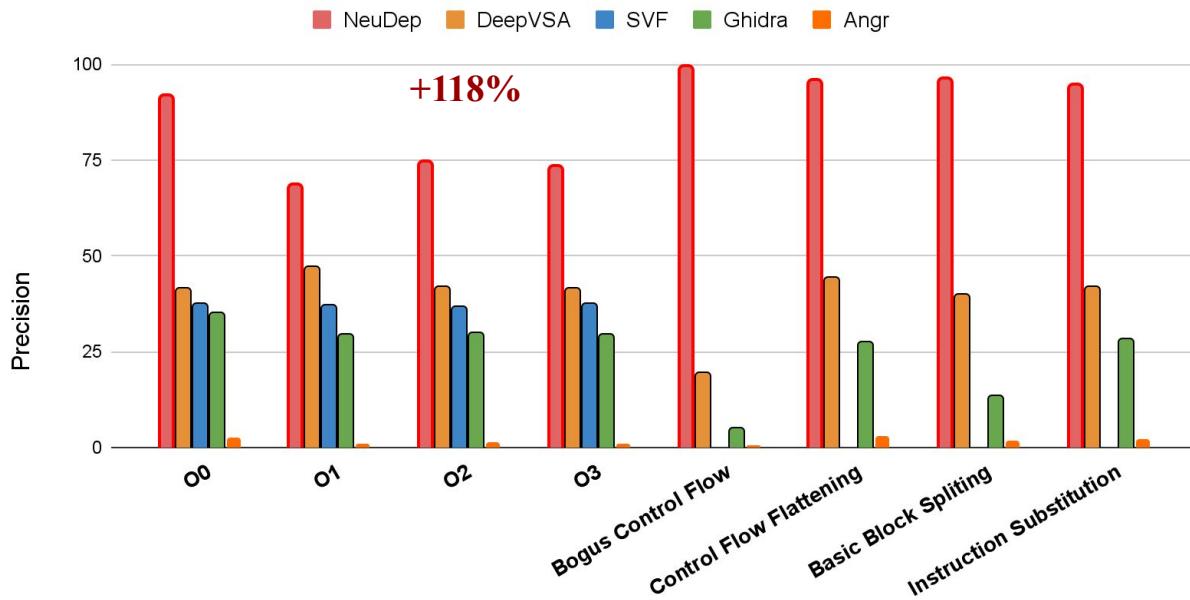
Finetuning for Analyzing Memory Dependence

Task 3: Memory Dependence Analysis

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



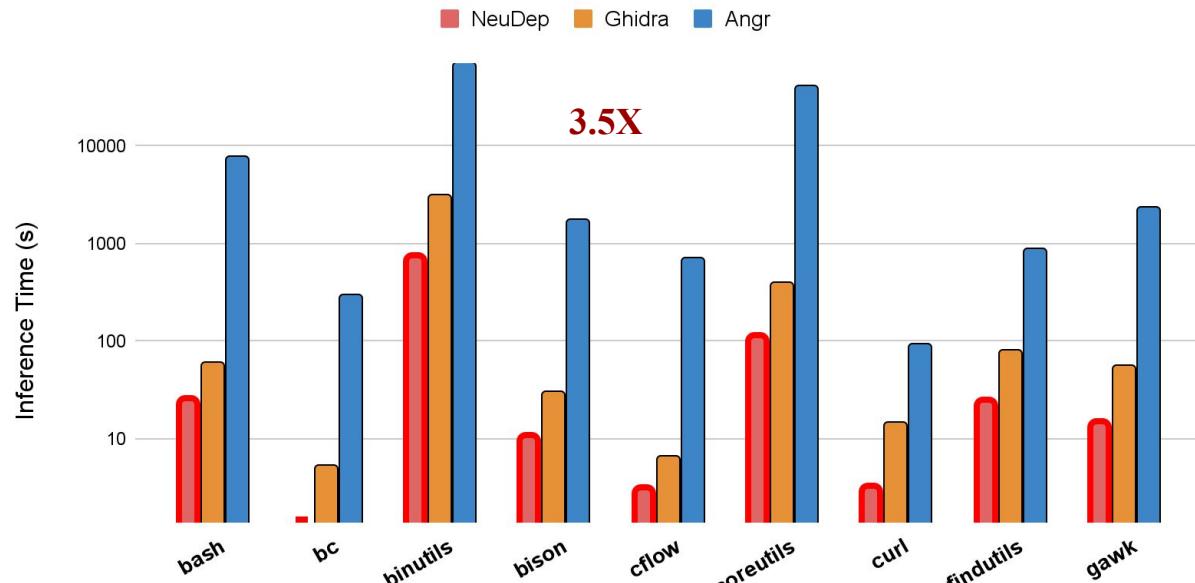
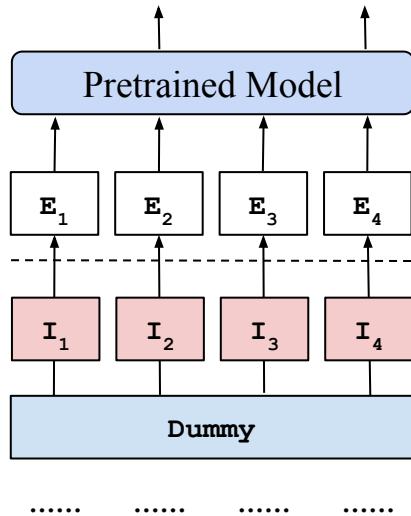
- Taint analysis
- Malware Analysis



Finetuning for Analyzing Memory Dependence: Inference Time

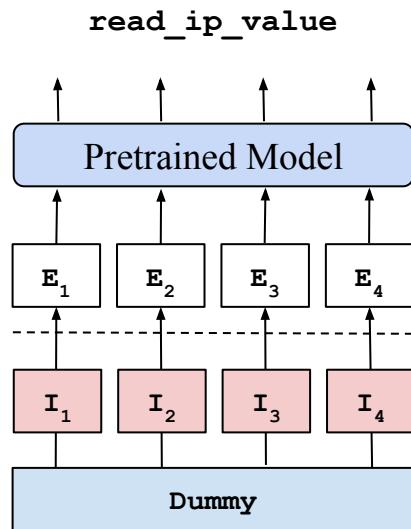
Task 3: Memory Dependence Analysis

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



Finetuning for Function Name Prediction and Memory Region Prediction

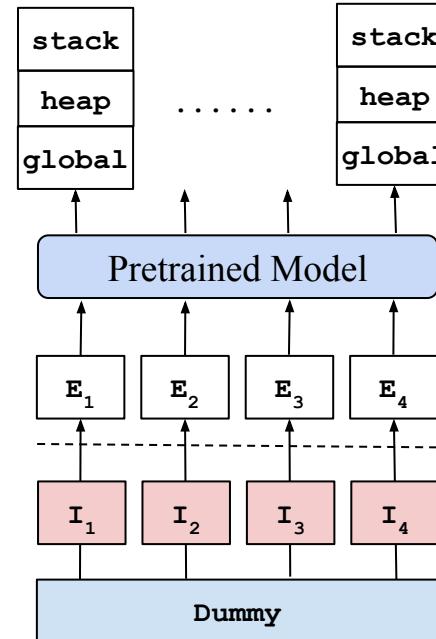
Task 4: Function Name Prediction



- 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

Firmware Images: Arm, MIPS with unknown compiler flags



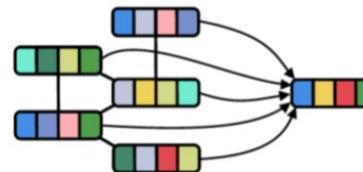
Ubiquiti sunMax TP-Link Deco-M4 NETGEAR R7000 Linksys RE7000

15 CVEs

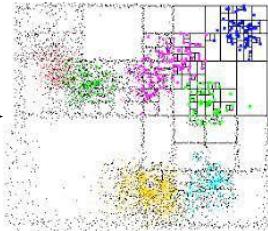
16 CVEs

7 CVEs

14 CVEs



Learned Function Embeddings



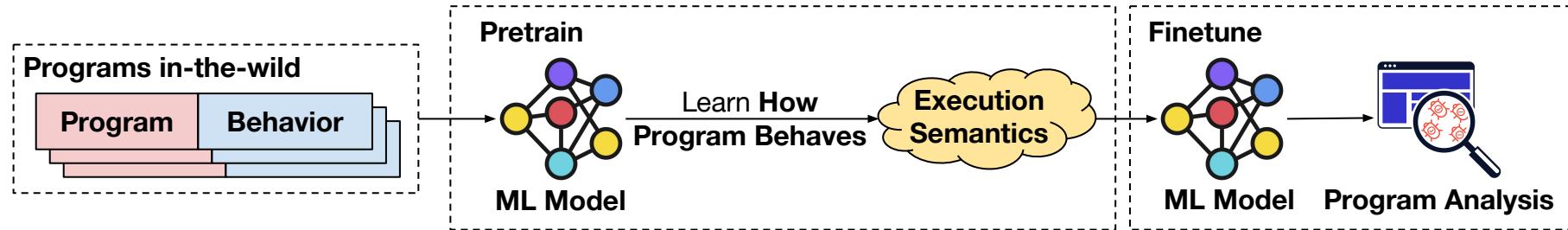
Approximate Nearest Neighbors



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

Summary: Learning Program Semantics via Execution-Aware Pre-training

Improves Program Analysis for Security Applications



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

-00 -01 -02
-03 -0d 0x



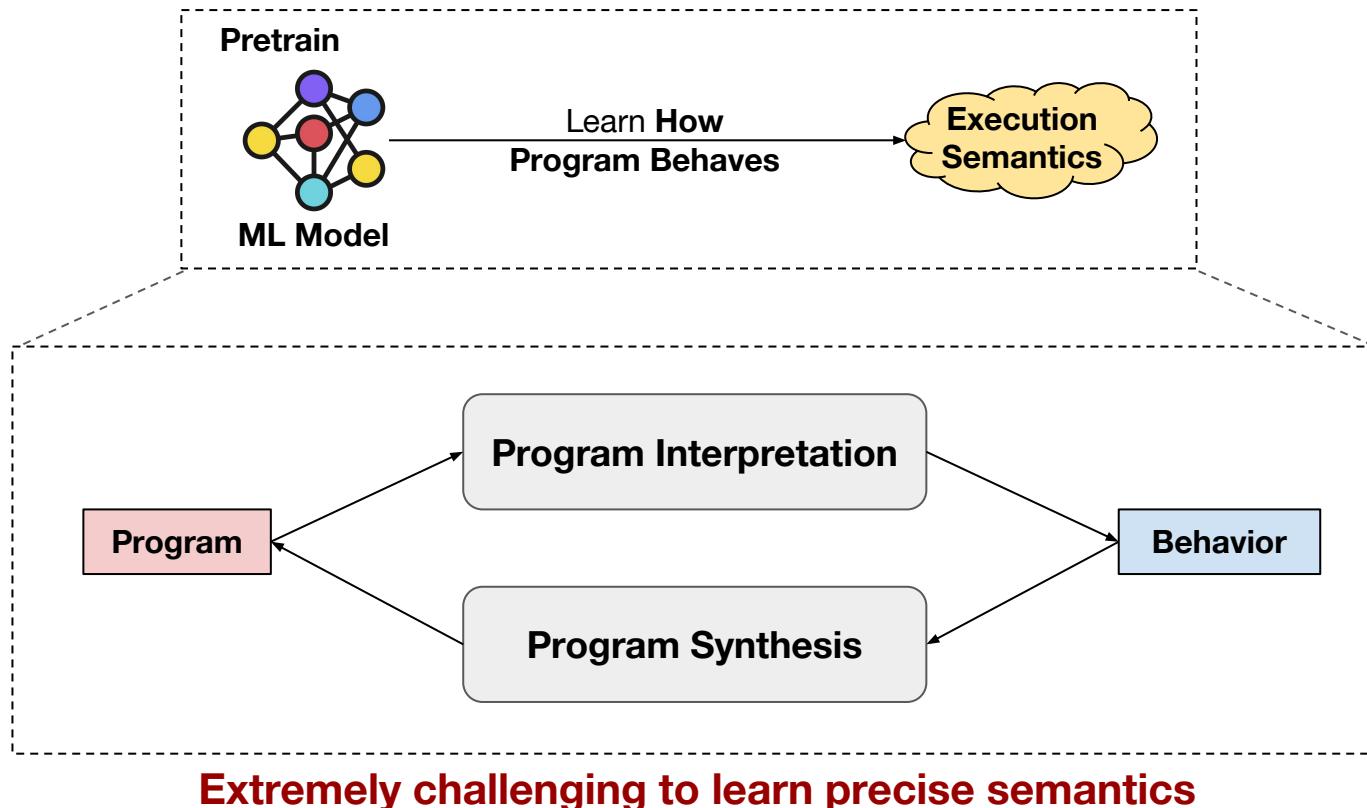
Optimizations



Obfuscations



Limitation: Learning Execution-Aware Program Representations is Challenging



Limitation: What the Model has Learned during Pretraining?

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

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

Machine Learning Security

