

Learning Program Dependencies with ML

Dongdong She

Kexin Pei

Abhishek Shah

What's happening?

```
input_byte = read("foo.txt")
if input_byte[0] == '1':
    // divide by zero error
else:
    // exit gracefully
```

Problem

- Programs have complex dependencies
 - Control flow
 - Data flow
- Program Analysis
 - Examples:
 - Taint Analysis
 - Symbolic Execution
 - Do not scale well

Solution

- Use ML to automatically learn dependencies
 - ML can excel at finding relationships in data
 - NLP: tagging parts-of-speech in a sentence
- Finding them is useful
 - Code Coverage
 - Debugging
 - Vulnerability Discovery

Outline

- Learn Dependencies
 - Examine 1 program dynamically
 - Neuro-symbolic Execution
 - Learn path dependencies
 - Examine many programs statically
 - Idea of “Big Code”
 - Binaries
 - Learn dependency between variable and registers
 - Source Code
 - Learn dependency between variable use and definition

Neuro-symbolic Execution: **Augmenting** Symbolic Execution with **Neural Constraints**

Dongdong She

Problem

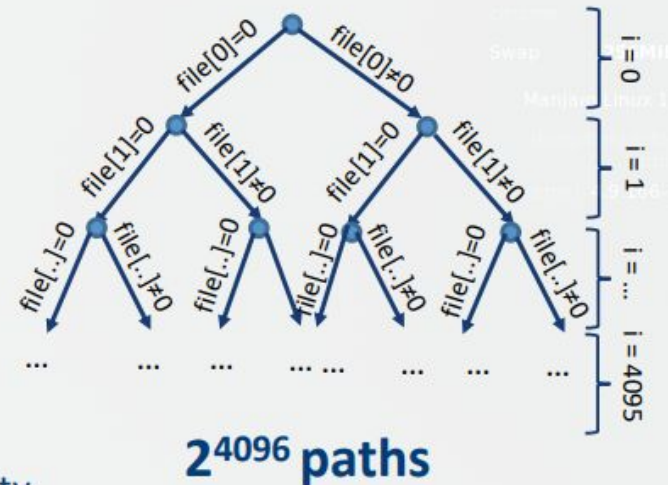
Symbolic execution has many limitations.

- Poor scalability by path explosion.
- Language-specific implementation.
- Failure to model complex dependencies.
- Limited expressiveness of satisfying theories.

A simple example

```
1 int main (...) {
2   if (strlen(filename)>1 && filename[0]=='-')
3     exit(1)
4   copy_data(...);
5   ...
6 }
7 void copy_data(..., int *file,...) {
8   static double data[4096], value;
9   read_double value(file, ... );
10  value = fabs (data [0]);
11  for(i=0; i<4096; i++)
12    if(file[i] == 0.0) count++;
13  data[1] /= (value+count-3);
14  ...
15 }
```

- #1 Limitations of SMT Solvers
 - #2 Unmodeled Semantics
 - #3 Path Explosion
- Candidate Vulnerability Point (CVP): Divide-by-zero



Locations of interest

Candidate Vulnerability Point(CVP)

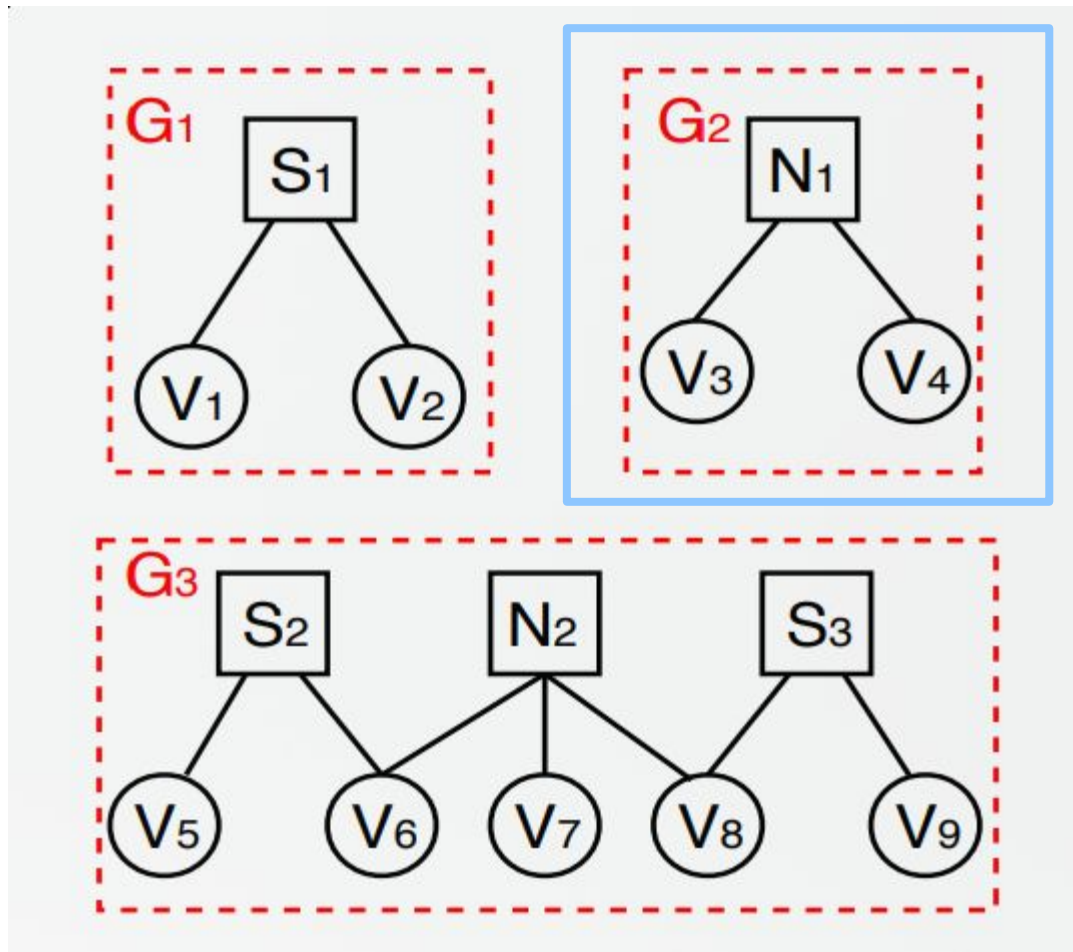
- Statically analyze program in advance.
- Identify two specific program locations.
 - Division operations (check zero division).
 - Boundary checking in buffer accesses.
- Instrument CVP to record values of denominators/ index number for further training.

Neuro-symbolic execution

- Represent most of program logic with symbolic constraints.
- **Approximate** the remaining logic that is hard to solve with **NN**.
- **Solve** the **combination** of exact constraints & approximated constraints.

$$\text{Constraints} = S_1 \wedge N_1 \wedge N_2 \wedge S_2$$

Neuro-symbolic execution



How to generate the neuro-constraints?

Neuro-constraints

Problem

Inputs:

1. Source code
2. Symbolic Variables
(e.g., filename & file)
3. Candidate Vulnerability Points (CVPs)

- Divide by zero
- Buffer overflow

Outputs:

Validated Exploits

```
1 int main (...) {
2   if (strlen filename >1 && filename[0]=='-')
3     exit(1)
4   copy_data(...);
5   ...
6 }
7 void copy_data(..., int *file ...) {
8   static double data[4096], value;
9   read_double_value(file, ... );
10  value = fabs (data [0]);
11  for(i=0; i<4096; i++)
12    if(file[i] == 0.0) count++;
13    data[1] /= (value+count-3); CVP: Divide-by-zero
14  ...
15 }
```

Neuro-constraints

Key Insights

Values of Symbolic
Variables



Values of Vulnerable
Variables in CVP

Learn an approximation
with small number of I/O
examples

100,000 samples

```
1 int main (...) {
2   if (strlen filename >1 && filename[0]=='-')
3     exit(1)
4   copy_data(...);
5   ...
6 }
7 void copy_data(..., int *file ...) {
8   [redacted]
9   [redacted]
10  [redacted]
11  [redacted]
12  [redacted]
13  data[1] /= (value+count-3); CVP: Divide-by-zero
14  ...
15 }
```


CPU
RAM
Swap 254MiB
Manjaro Linux 3.5
4.9.166.1

Neuro-constraints

1. Neural nets can represent a large category of functions (universal approximation theorem).

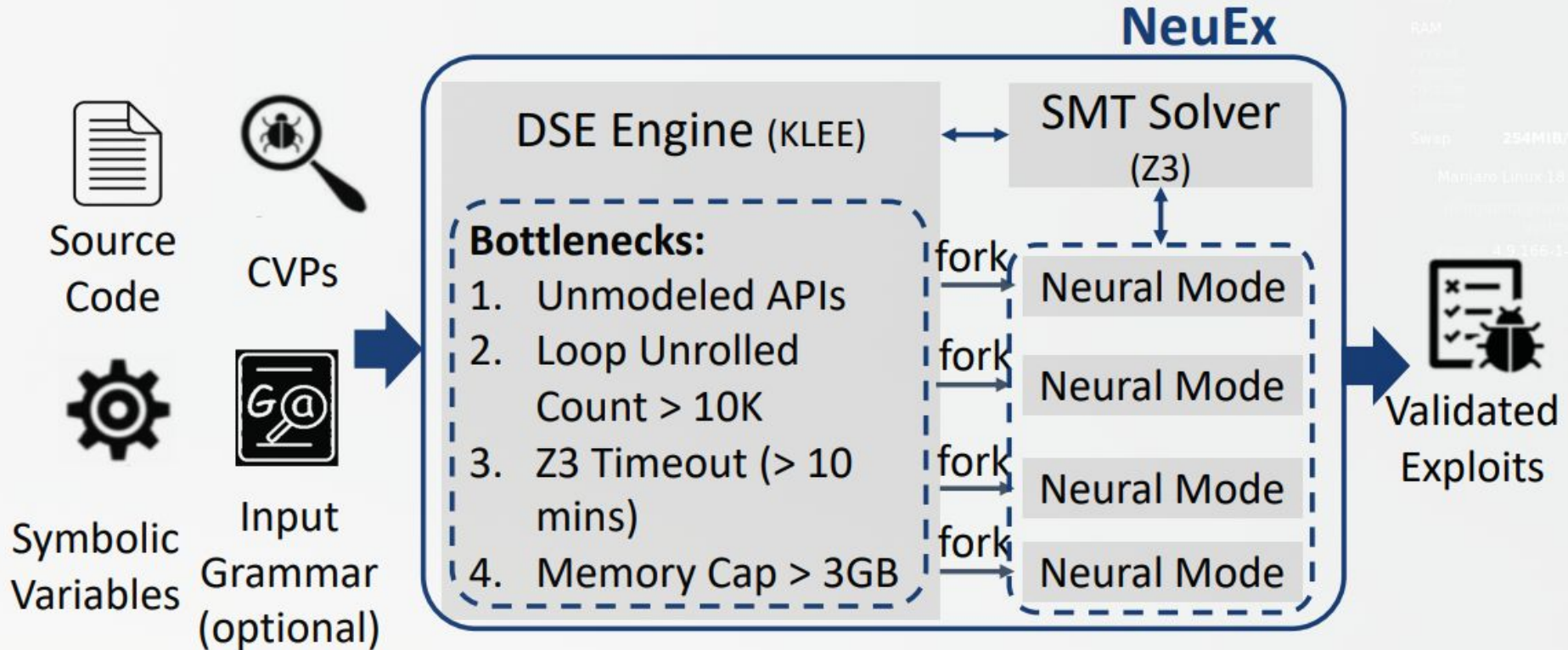
2. Multiple applications show that neural nets are learnable for many practical functions.

Approach

```
1 int main (...) {
2   if (strlen filename >1 && filename[0]=='-')
3     exit(1)
4   copy_data(...);
5   ...
6 }
7 void copy_data(..., int *file ...) {
8   Approximate Constraint (as a neural net):
9
10  file →  → count & value
11
12  data[1] /= (value+count-3); CVP: Divide-by-zero
13  ...
14  ...
15 }
```

Overview

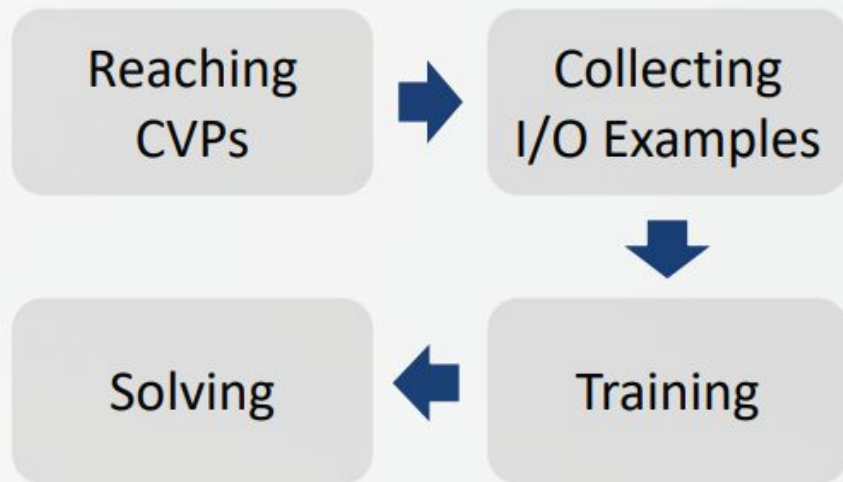
NeuEx Overview



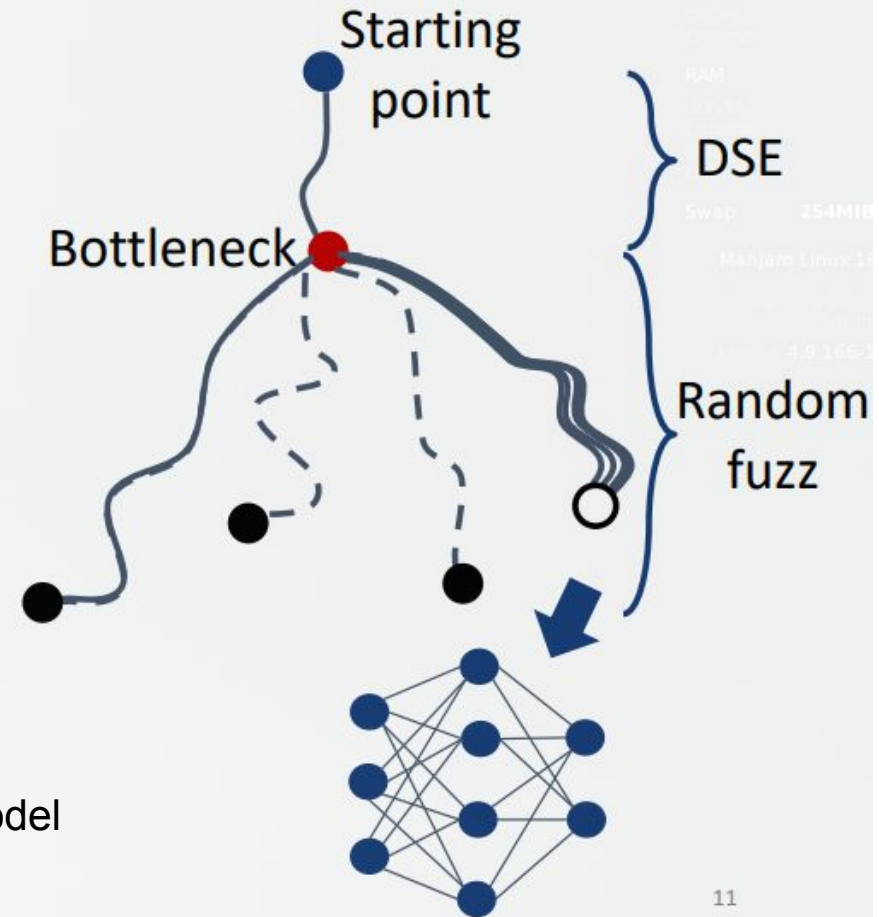
Hybrid mode design (symbolic mode and neural mode)

Neural Mode

Neural Mode



- MLP + ReLU
- Simple regression model
- 100,000 samples



Constraints

Generated Constraints

1. Reachability constraints:

$$\begin{aligned} & \text{strlen}(\text{filename}) \leq 1 \\ & \forall \text{filename} \neq \text{'-' } \end{aligned}$$

\wedge

$N: \text{infile} \rightarrow (\text{value}, \text{count})$

2. Vulnerability condition:

$$\text{value} + \text{count} - 3 == 0$$

```
1 int main (...) {
2   if (strlen(filename)>1 && filename[0]=='-')
3     exit(1)
4   copy_data(...);
5   ...
6 }
7 void copy_data(..., int *file,...) {
8   static double data[4096], value;
9   read_double_value(file, ... );
10  value = fabs (data [0]);
11  for(i=0; i<4096; i++)
12    if(file[i] == 0.0) count++;
13  data[1] /= (value+count-3);
14  ...
15 }
```

CVP: Divide-by-zero

Constraints

Constraint Solving

1. Reachability constraints:

$$\begin{aligned} \text{strlen}(\text{filename}) &\leq 1 \\ \forall \text{filename} &\neq \text{'-' } \end{aligned}$$

\wedge

2. Vulnerability condition:

$$\text{value} + \text{count} - 3 == 0$$

Purely symbolic constraints:

No variable shared with neural constraints

➔ SMT solver

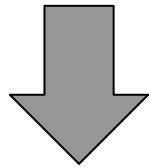
Mixed constraints:

Including both neural constraints and symbolic constraints with shared variables



How to solve mixed constraints

Symbolic constraints



Optimization objectives of the neural net

Encoding Symbolic Constraints as an Optimization Objective

$N: infile \rightarrow (value, count) \wedge value + count - 3 == 0$

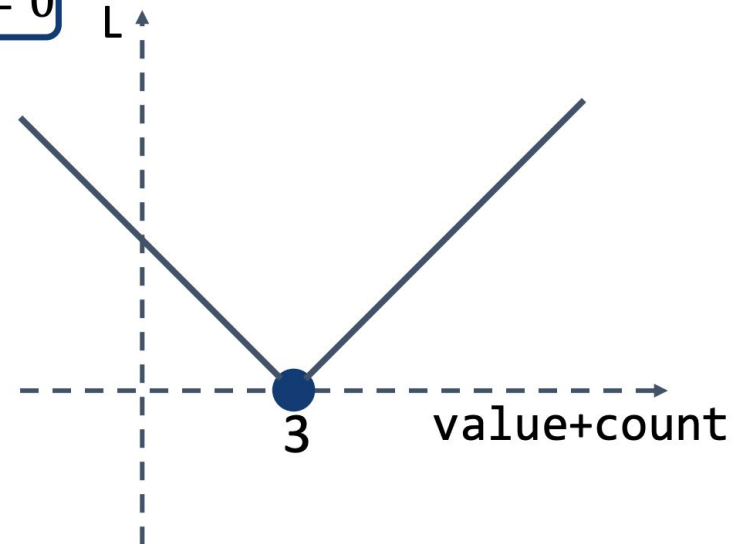
Symbolic constraint

Criterion for crafting the loss function:
The minimum point of the loss function satisfies the symbolic constraints.



One possible encoding:

$$L = \text{abs}(value + count - 3)$$

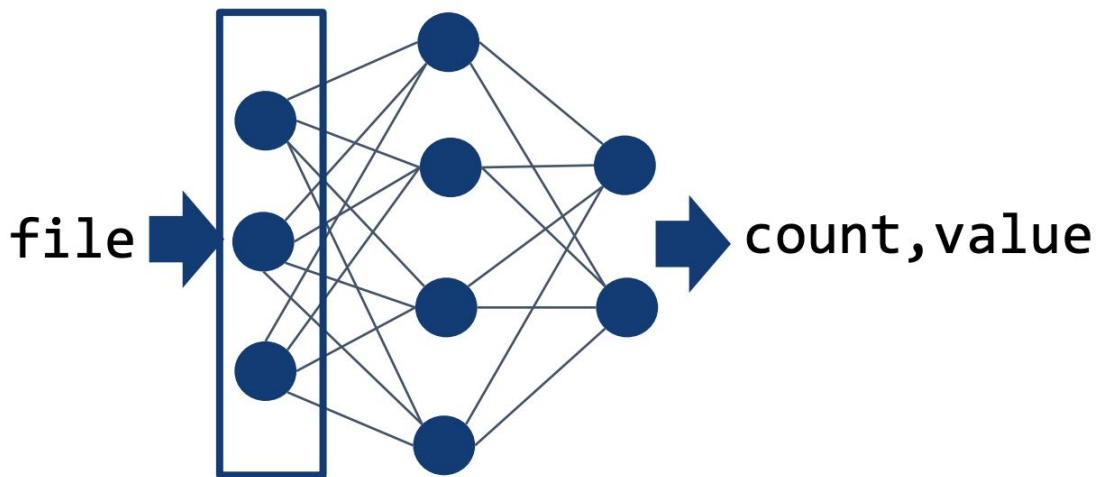


Constraints \Rightarrow Loss

Symbolic Constraint	Loss Function (L)
$S_1 ::= a < b$	$L = \max(a - b + \alpha, 0)$
$S_1 ::= a > b$	$L = \max(b - a + \alpha, 0)$
$S_1 ::= a \leq b$	$L = \max(a - b, 0)$
$S_1 ::= a \geq b$	$L = \max(b - a, 0)$
$S_1 ::= a = b$	$L = \text{abs}(a - b)$
$S_1 ::= a \neq b$	$L = \max(-1, -\text{abs}(a - b + \beta))$
$S_1 \wedge S_2$	$L = L_{S_1} + L_{S_2}$
$S_1 \vee S_2$	$L = \min(L_{S_1}, L_{S_2})$

Solving Mixed Constraints via Gradient Descent

Gradient: $\nabla_{file} L$



count	value	loss
257	20	274
1	20	18
0	20	17
...
0	3	0

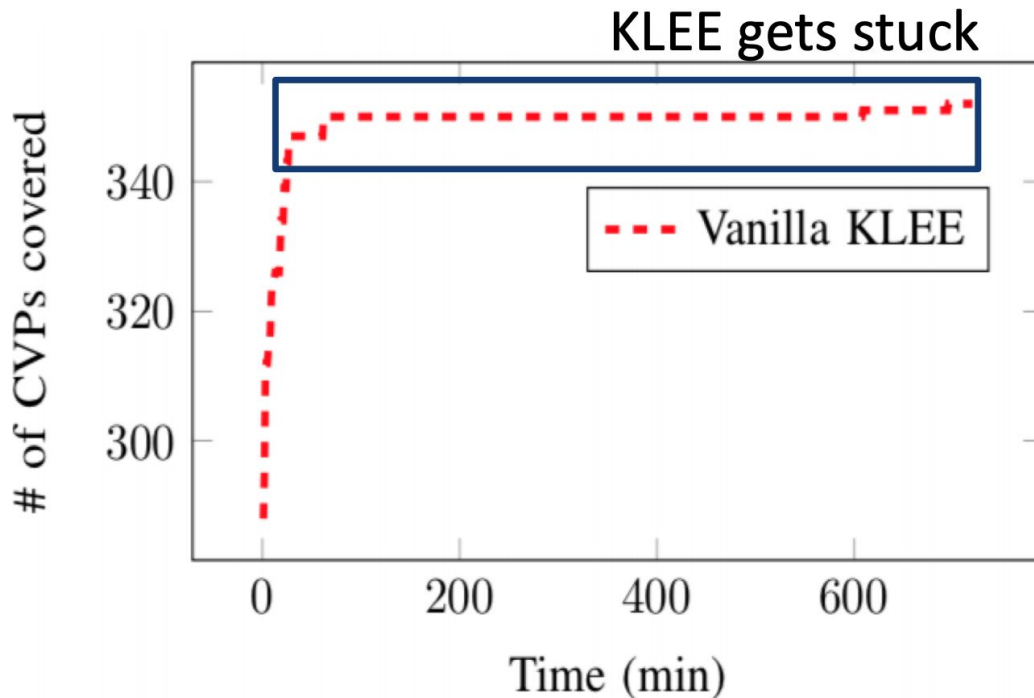
Concretely validate the exploit

file: 000...

Evaluation

- **Recall:** Neural mode is only triggered when DSE encounters bottlenecks
 - **Benchmarks:** 7 Programs known to be difficult for classic DSE
 - 4 Real programs
 - cURL: Data transferring
 - SQLite: Database
 - libTIFF: Image processing
 - libsndfile: Audio processing
 - LESE benchmarks
 - BIND, Sendmail, and WuFTP
- Include:
1. Complex loops
 2. Floating-point variables
 3. Unmodeled APIs

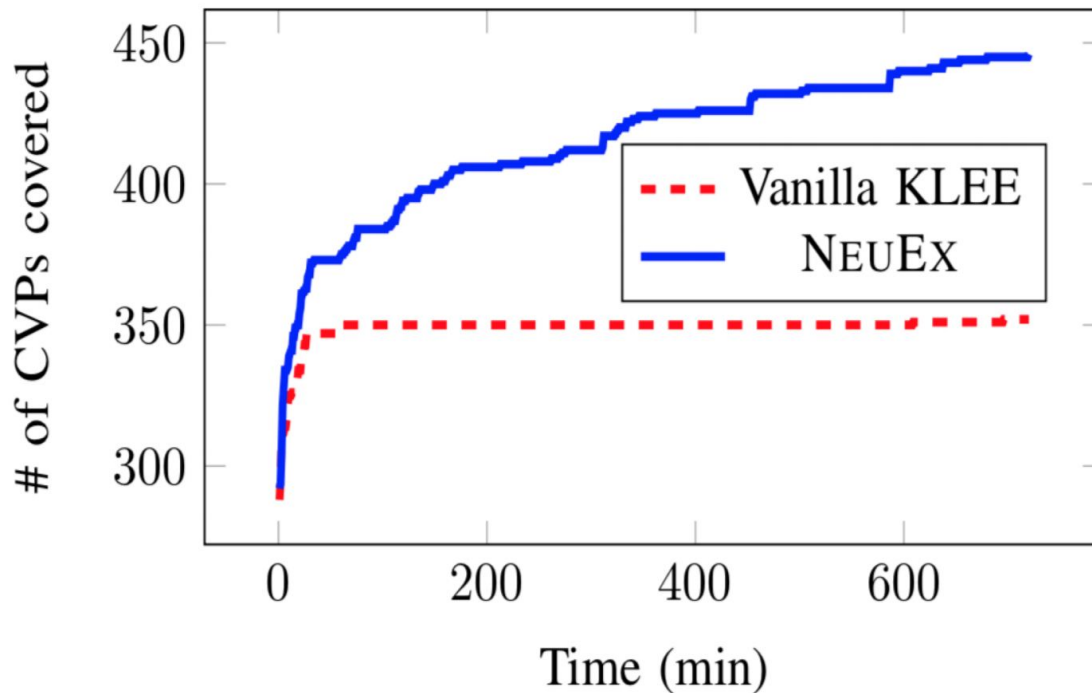
CVP Coverage & Bottlenecks for DSE



of bottlenecks: 61

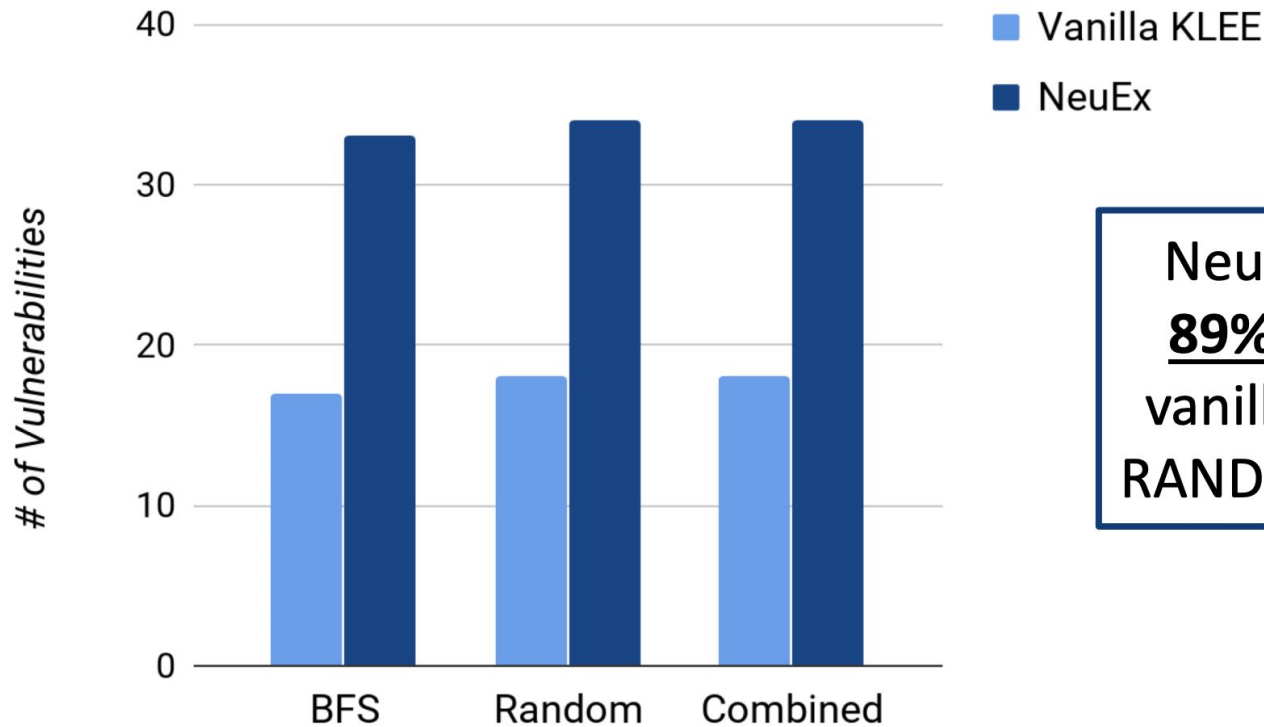
- Unmodeled APIs (6)
- Complicated loops (53)
- Z3 timeout (1)
- Memory exhaustion (1)

CVP Coverage of NeuEx vs KLEE



The number of CVPs reached or covered by NeuEx is **25%** higher than vanilla KLEE.

of Bugs Found by NeuEx vs KLEE



NeuEx finds **94%** and **89%** more bugs than vanilla KLEE in BFS and RAND mode in 12 hours.

Debin: Recovering Stripped Info from Binaries

Kexin Pei

Binaries with debug symbols

x86 malware samples from VirusShare

Assembly

```
80534BA:  
push %ebp  
push %edi  
push %esi ...
```

Debug symbols

```
80534BA  rfc1035_init  int  
8053DB1  fopen64      int  
8063320  num_entries  int  
      ⋮
```

Binary with debug symbols

Hex-rays

Descriptive names for
functions and variables

```
int rfc1035_init() {  
    ...  
    if ( num_entries <= 0 ) {  
        v0 = ("/etc/resolv.conf", 'r');  
        if ( v0 || (v1 =  
            fopen64("resolv.conf"))){  
            // code to read and  
            // manipulate DNS settings  
        }  
        ...  
    }  
}
```

Decompiled code

Stripped Binaries

Assembly

```
80534BA:  
push %ebp  
push %edi  
push %esi ...
```

Debug symbols



Hex-rays

Non-descriptive names

```
int sub_80534BA() {  
    ...  
    if ( dword_8063320 <= 0 ) {  
        v0 = ("/etc/resolv.conf", 'r');  
        if ( v0 || (v1 =  
            sub_8053B1("resolv.conf"))){  
            ...  
            ...  
        }  
    }  
}
```

Can we recover the debug symbols?

Stripped binary

Challenges

1. No mapping from registers and memory offsets to semantic variables

Computes
 $1 + 2 + \dots + n$

```
<sum> start:  
  mov 4(%esp), %ecx  
  mov $0, %eax  
  mov $1, %edx  
  add %edx, %eax  
  add $1, %edx  
  cmp %ecx, %edx  
  jne 8048400  
  repz ret  
<sum> end
```

Stores the value of
a semantic variable

Stores intermediate
(non-semantic) value

Challenges

2. No names and types

```
<sum> start:  
  mov 4(%esp), %ecx  
  mov $0, %eax  
  mov $1, %edx  
  add %edx, %eax  
  add $1, %edx  
  cmp %ecx, %edx  
  jne 8048400  
  repz ret  
<sum> end
```

Store the values of
the unsigned integer
variable n

Stores the result in an
integer variable res

DeBIN: Recovering debug info

Assembly

```
<sum> start:
  mov  4(%esp), %ecx
  mov  $0, %eax
  mov  $1, %edx
  add  %edx, %eax
  add  $1, %edx
  cmp  %ecx, %edx
  jne  8048400
  repz ret
<sum> end
```

Debug information



Assembly

```
<sum> start:
  mov  4(%esp), %ecx
  mov  $0, %eax
  mov  $1, %edx
  add  %edx, %eax
  add  $1, %edx
  cmp  %ecx, %edx
  jne  8048400
  repz ret
<sum> end
```

Debug information

Location	Name	Type
	sum	int
	n	uint
	i	uint
	res	int

Design Choice

How will you do this?

Assembly

```
<sum> start:  
  mov  4(%esp), %ecx  
  mov  $0, %eax  
  mov  $1, %edx  
  add  %edx, %eax  
  add  $1, %edx  
  cmp  %ecx, %edx  
  jne  8048400  
  repz ret  
<sum> end
```

Debug information



Assembly

```
<sum> start:  
  mov  4(%esp), %ecx  
  mov  $0, %eax  
  mov  $1, %edx  
  add  %edx, %eax  
  add  $1, %edx  
  cmp  %ecx, %edx  
  jne  8048400  
  repz ret  
<sum> end
```

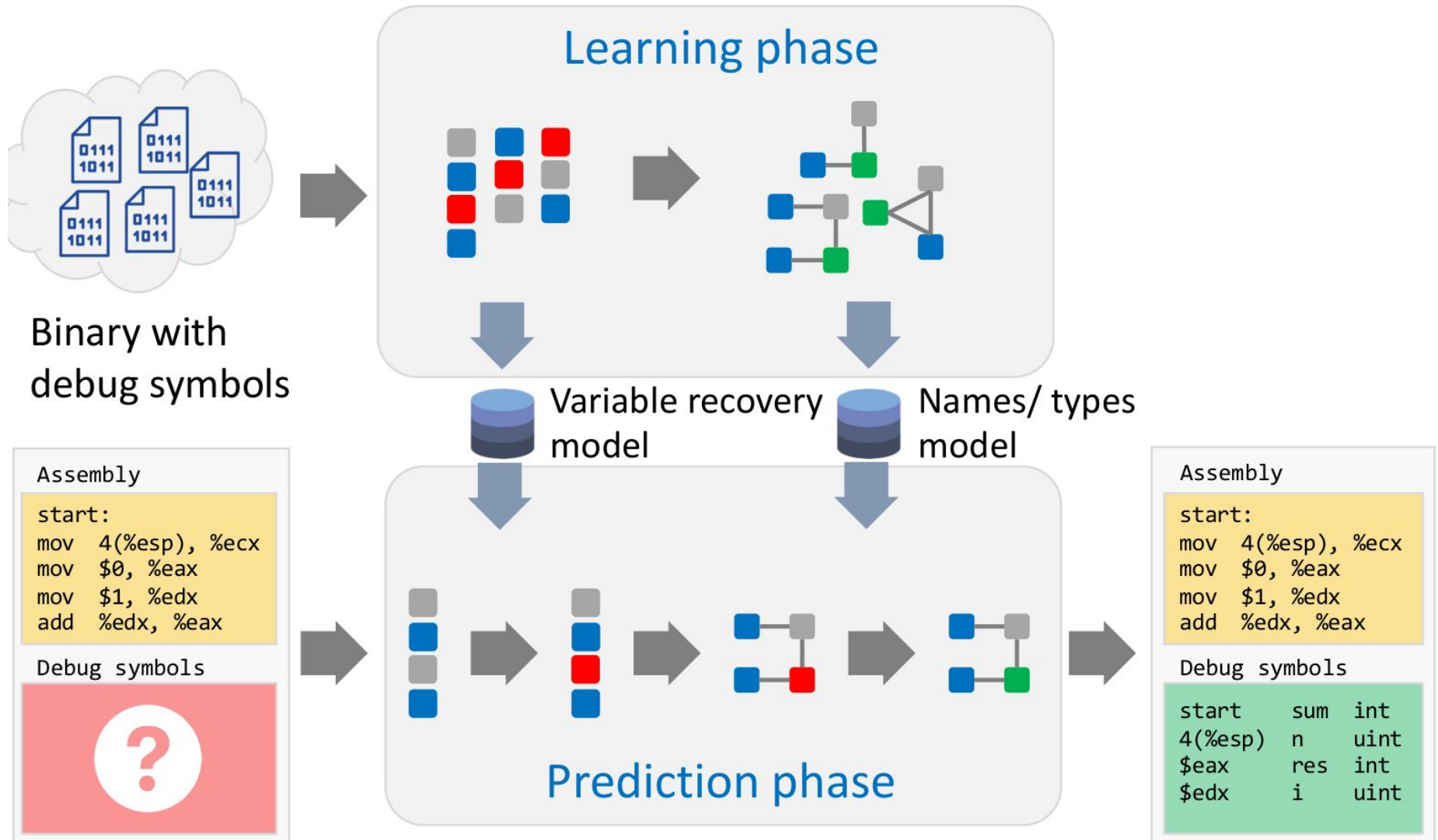
Debug information

Location	Name	Type
	sum	int
	n	uint
	i	uint
	res	int

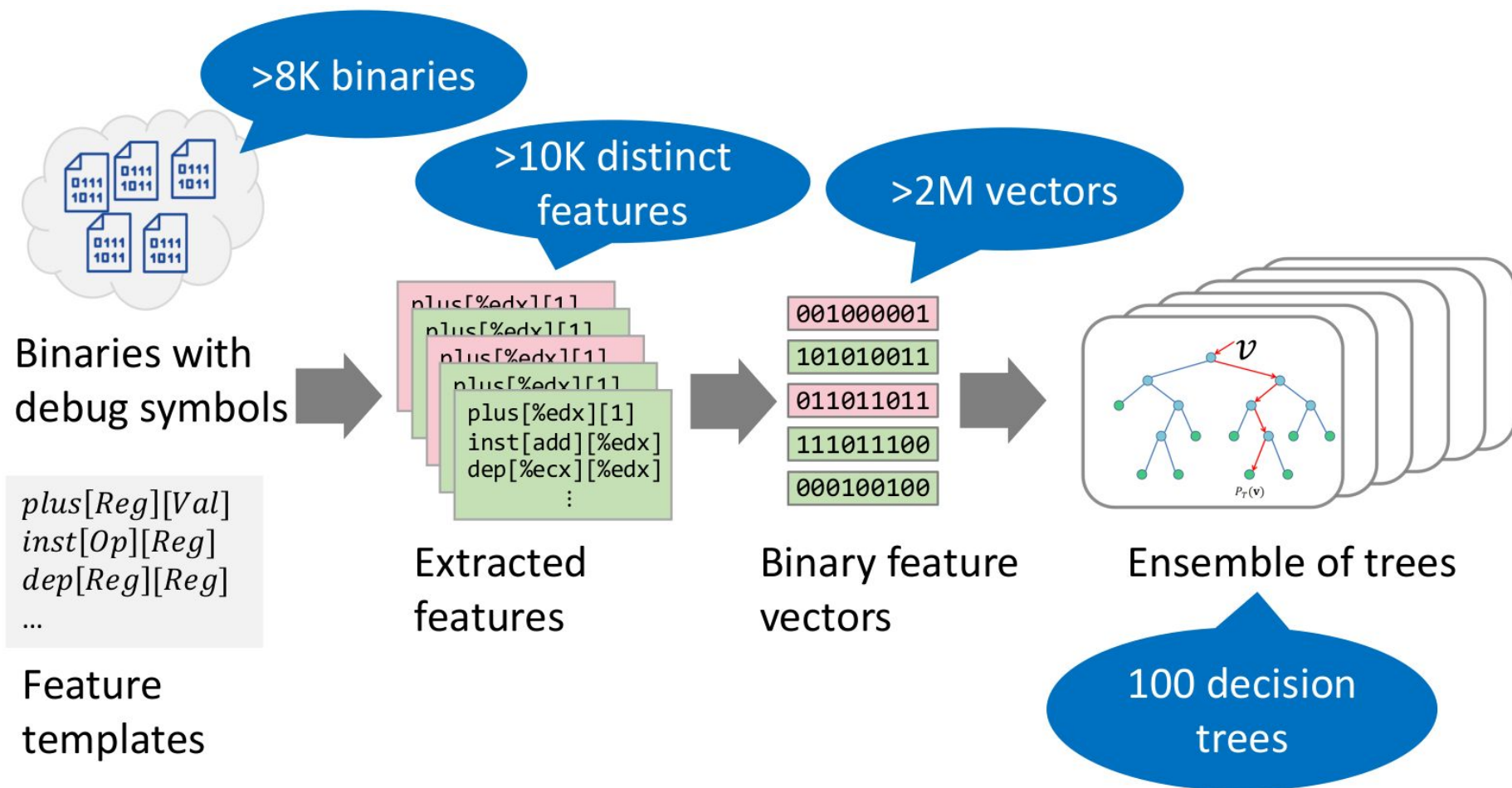
Recap: importance of dependency

1. Naive way of doing this?
 - a. Feature template
 - b. Individual classification
2. Smarter way of doing this?
 - a. RNN/LSTM
 - b. Sequential dependency
3. More advanced (best result):
 - a. PGM(CRF,MRF,Bayesian Network)/TreeLSTM/GNN/GCN/GGNN...
 - b. **Structured** learning

How does Debin work?

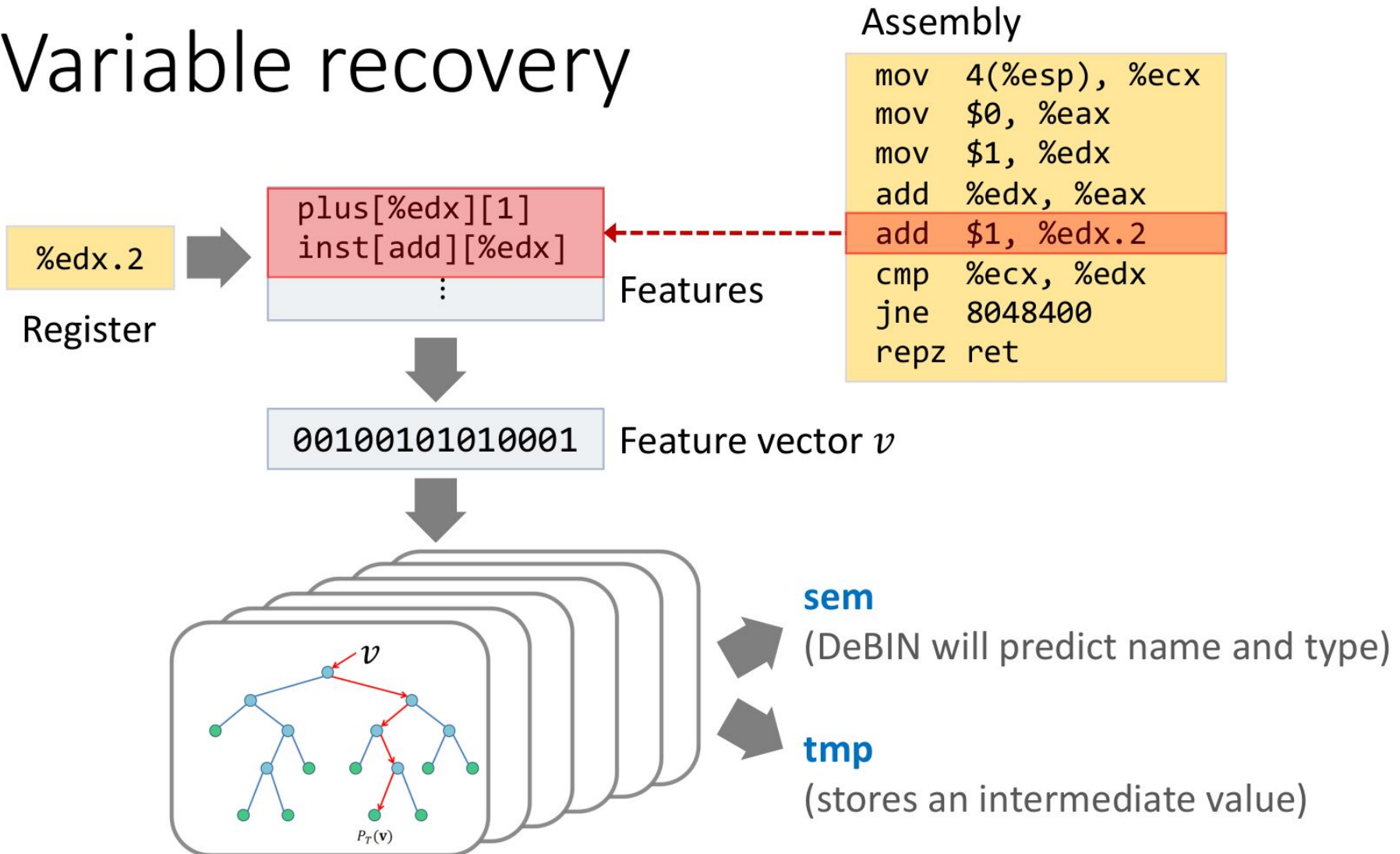


Step 1: Recovering Variables



Step 1: Recovering Variables

Variable recovery



Extremely randomized trees

Decision tree:

- One dataset
- All features

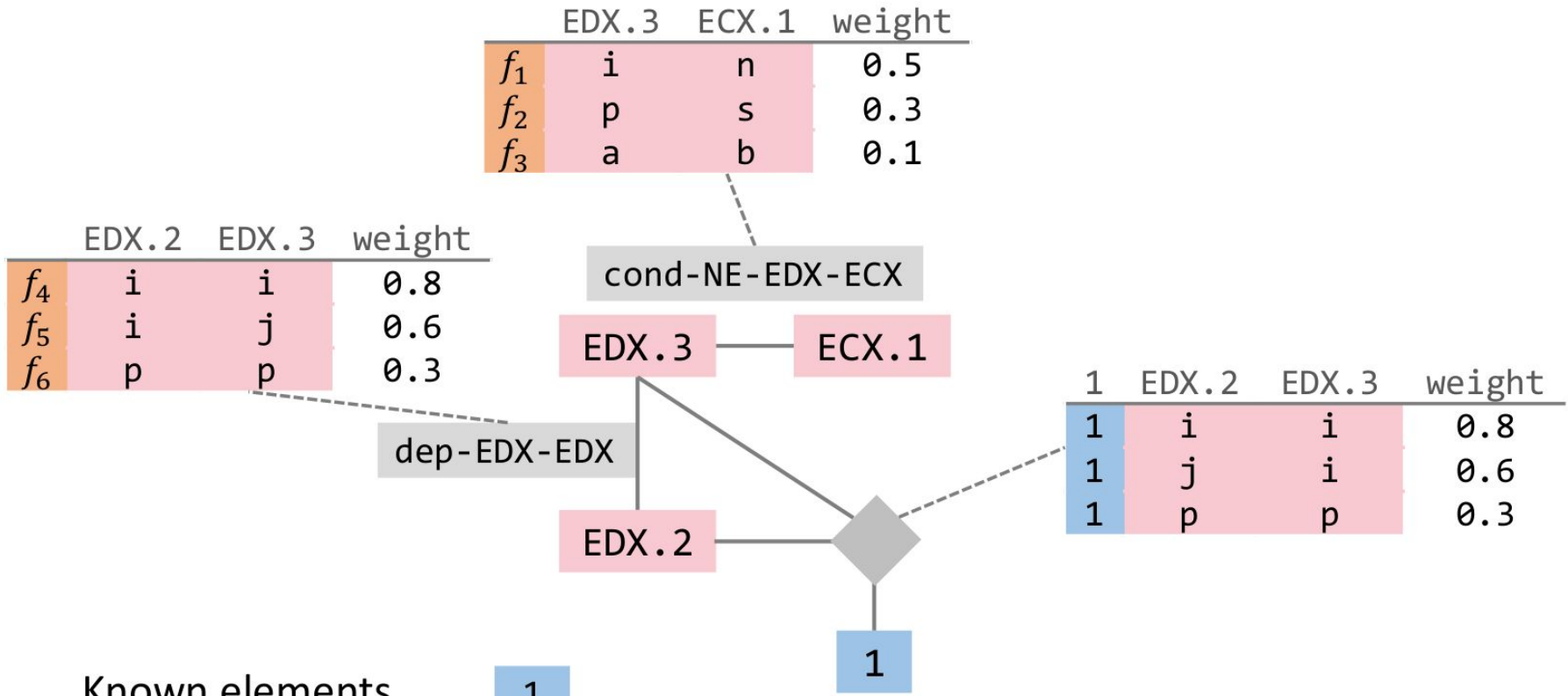
Random forest:

- Multiple sampled sub-dataset
- Sampled set of features

Extremely Randomized trees:

- Randomized division of feature values

Step 2: Predicting names and types



Known elements

1

Unknown elements

ECX.1, ...

Binary features

f_1, f_2, \dots

Factors



Pairwise Feature functions

Relationship	Template	Condition for adding an edge
<i>Function Relationships</i>		
Element used in Function	$(f, v, \text{func-loc}(v))$	variable v is accessed inside the scope of function f
	$(f, a, \text{arg-loc}(a))$	variable a is an argument of function f by calling conventions
	$(f, c, \text{func-str})$	string constant c is accessed inside the scope of function f
	$(f, s, \text{func-stack})$	stack location s is allocated for function f
Function Call	(f_1, f_2, call)	function f_2 is called by function f_1
<i>Variable Relationships</i>		
Instruction	$(v, \text{insn}, \text{insn-loc}(v))$	there is an instruction insn (e.g., <code>add</code>) that operates on variable v
Location	$(v, l, \text{locates-at})$	variable v locates at location l (e.g., memory offset <code>mem[RSP+16]</code>)
Locality	$(v_1, v_2, \text{local-loc}(v_1))$	variable v_1 and v_2 are locally allocated (e.g., <code>EDX.2</code> and <code>EDX.3</code>)
Dependency	$(v_1, v_2, \text{dep-loc}(v_1)\text{-loc}(v_2))$	variable v_1 is dependent on variable v_2
Operation	$(v, op, \text{unary-loc}(v))$	unary operation op (e.g. unsigned and low cast) on variable v
	$(n_1, n_2, \text{op-loc}(n_1)\text{-loc}(n_2))$	binary operation op (e.g., <code>+</code> , left shift <code><<</code> and etc.) on node n_1 and n_2
	$(v_1, v_2, \text{phi-loc}(v_1))$	there is a ϕ expression in BAP-IR: $v_1 = \phi(\dots v_2, \dots)$
Conditional	$(v, op, \text{cond-unary})$	there is a conditional expression $op(v)$ (e.g., <code>not(EDX.2)</code>)
	$(n_1, n_2, \text{cond-op-loc}(n_1)\text{-loc}(n_2))$	there is a conditional expression $n_1 op n_2$ (e.g. <code>EDX.3 != ECX.1</code>)
Argument	$(f, a, \text{call-arg-loc}(a))$	there is a call $f(\dots, a, \dots)$ with argument a
<i>Type Relationships</i>		
Operation	$(t, op, \text{t-unary-loc}(t))$	unary operation op on type t
	$(t_1, t_2, \text{t-op-loc}(t_1)\text{-loc}(t_2))$	binary operation op on type t_1 and t_2
	$(t_1, t_2, \text{t-phi-loc}(t_1))$	there is a ϕ expression: $t_1 = \phi(\dots t_2, \dots)$
Conditional	$(t, op, \text{t-cond-unary})$	there is a unary conditional expression $op(t)$
	$(t_1, t_2, \text{t-cond-op-loc}(t_1)\text{-loc}(t_2))$	there is a binary conditional expression $t_1 op t_2$
Argument	$(f, t, \text{t-call-arg-loc}(t))$	call $f(\dots, t, \dots)$ with an argument of type t
Name & Type	$(v, t, \text{type-loc}(v))$	variable v is of type t
	$(f, t, \text{func-type})$	function f is of type t



Factor Feature functions

Factors:

- All nodes that appear in the same ϕ expression of BAP-IR
- Function node of a call and its arguments
- Elements that are accessed in the same statement

Learning to predict

> 8,000 binaries



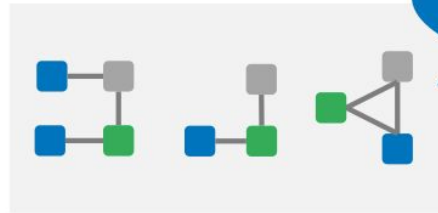
Binaries with debug symbols

(f_{unary}, Op, Var)
 $(f_{var-dep}, Var_1, Var_2)$
 ...

Feature templates

23 templates

Static analysis



Dependency graphs

binary features

f_1	i	n
f_2	p	s
f_3	a	b
f_4	i	i
f_5	i	j
f_6	p	p

Train model

	name1	name2	weight
f_1	i	n	0.4
f_2	p	s	0.5
f_3	a	b	0.2
f_4	i	i	0.3
f_5	i	j	0.6
f_6	p	p	0.4

3-factor			weight
1	i	i	0.4
1	j	i	0.2
1	p	p	0.1

4-factor				weight
1	i	i	k	0.3
1	j	i	a	0.5
1	p	p	v	0.2

Find **weights** that maximize $P(\vec{U} = \vec{u} | \vec{K} = \vec{k}_i)$ for all training samples (\vec{u}_i, \vec{k}_i)

Binary features and factors

End-to-end recovery of debug information

```
<sum> start :  
  mov  4(%esp), %ecx  
  mov  $0, %eax  
  mov  $1, %edx  
  add  %edx, %eax  
  add  $1, %edx.2  
  cmp  %ecx.1, %edx.3  
  jne  8048400  
  repz ret  
<sum> end
```

Stripped binary



Registers / mem offsets

EDX.2 EDX.3

EDX.1 ECX.1

Known elements

0 1 mov



Semantic variables

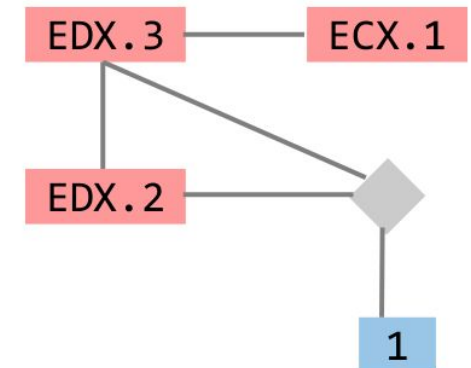
EDX.2 EDX.3 ECX.1

Temporary

EDX.1

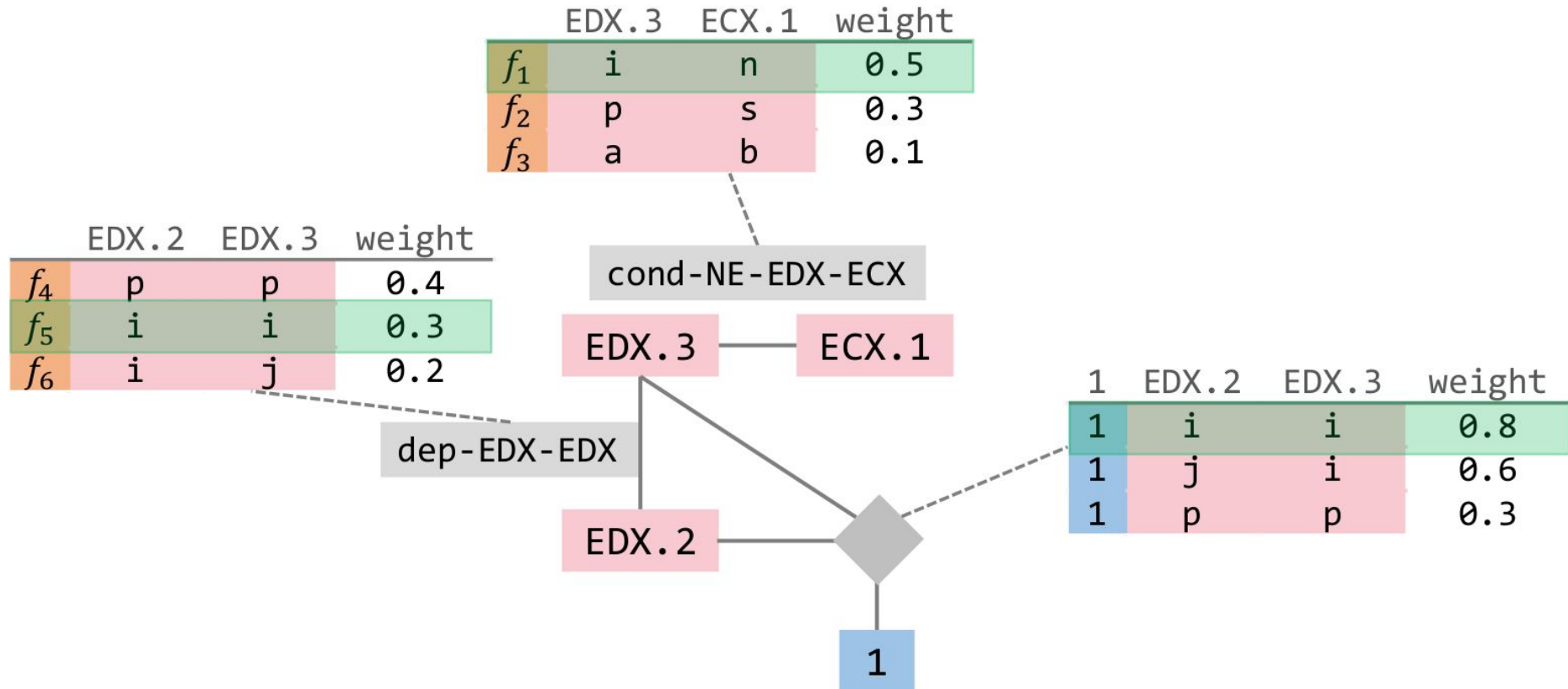
Known elements

0 1 mov



End-to-end recovery of debug information

MAP inference



End-to-end recovery of debug information

```
<sum> start :  
mov 4(%esp), %ecx  
mov $0, %eax  
mov $1, %edx  
add %edx, %eax  
add $1, %edx.2  
cmp %ecx.1, %edx.3  
jne 8048400  
repz ret  
<sum> end
```

Stripped binary

Registers / mem offsets

EDX.2 EDX.3

EDX.1 ECX.1

Known elements

0 1 mov

Semantic variables

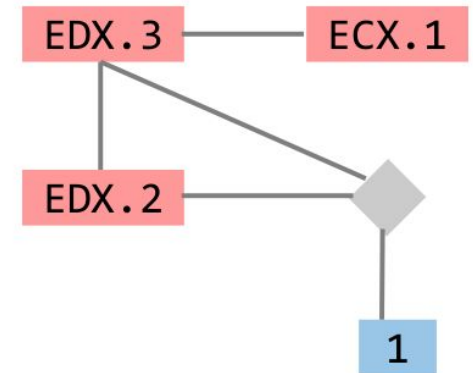
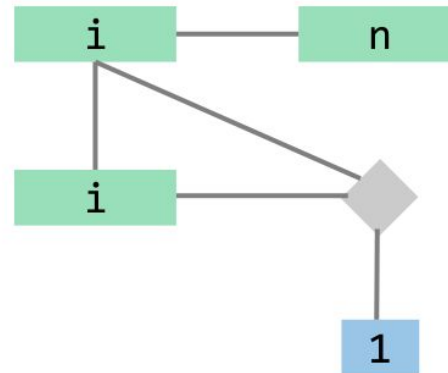
EDX.2 EDX.3 ECX.1





Temporary

EDX.1

Known elements

0 1 mov



Loc	Name	Type
	sum	int
	n	uint
	i	uint
	res	int

Debug information

Implementation

Static analysis: BAP

<https://github.com/BinaryAnalysisPlatform/bap/>

Learning and inference



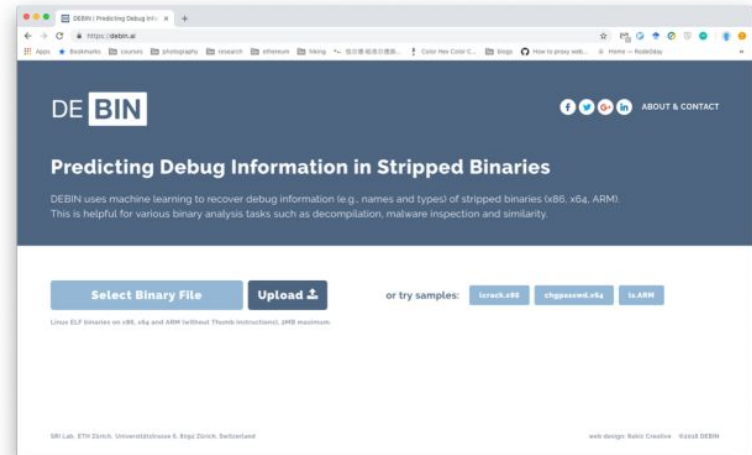
<http://scikit-learn.org>

NICE 2 Predict

<http://nice2predict.org>



830 Linux packages
x86, x64, ARM

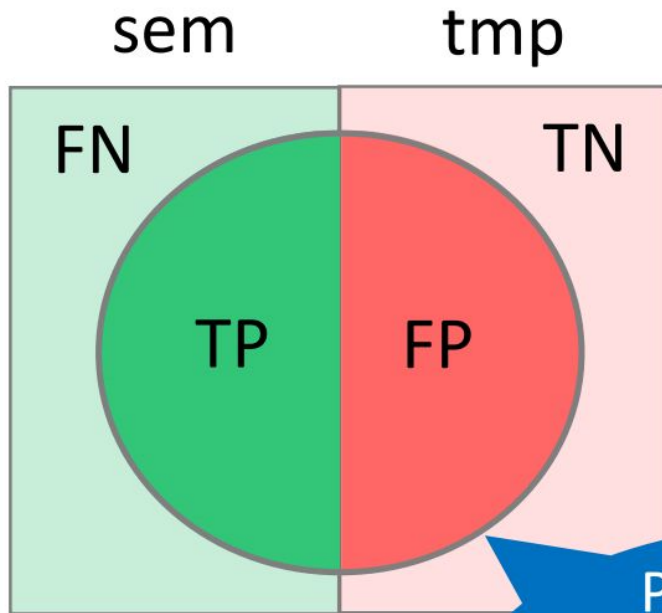


<https://debin.ai>

Evaluation

- How accurate is DeBIN's variable recovery?
- How accurate is DeBIN's name and type prediction?
- Is DeBIN useful for malware inspection?

Variable recovery accuracy



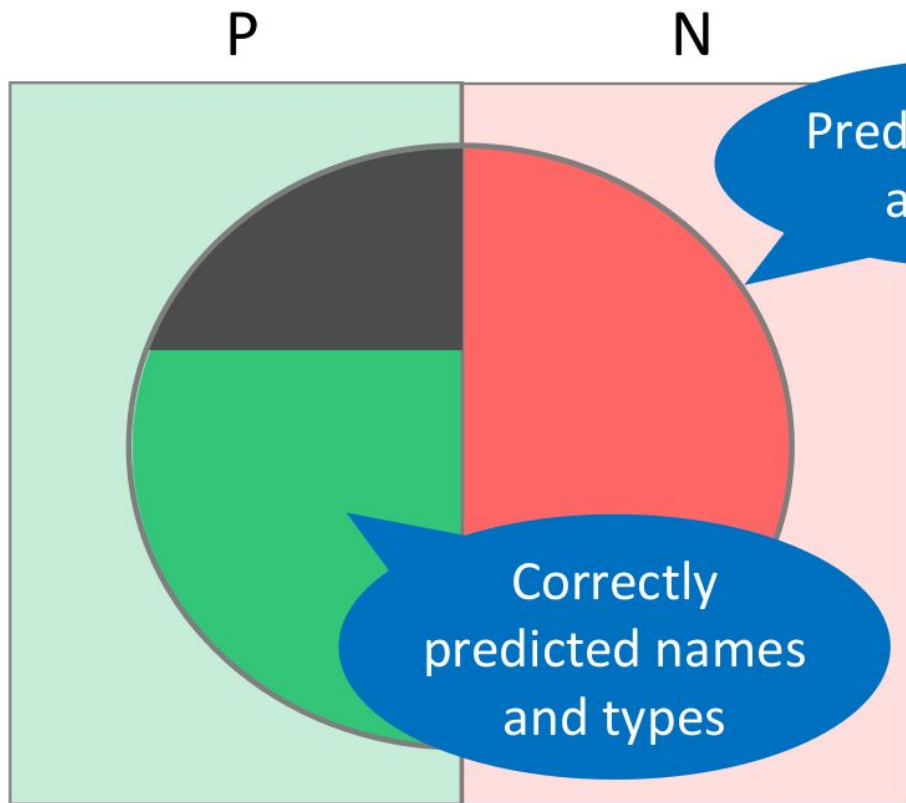
$$\text{Accuracy} = \frac{|\text{TP}| + |\text{TN}|}{|\text{sem}| + |\text{tmp}|} = \frac{\text{Green Circle} + \text{Red Circle}}{\text{Green Box} + \text{Red Box}}$$

Results

Arch	Accuracy
x86	87.1%
x64	88.9%
ARM	90.6%

Predicted as semantic registers and memory offsets

Name and type prediction accuracy



Total names and types (P) =



Predicted names and types (PN) =



Correct Predictions (CP) =



$$\text{Precision} = \frac{|CP|}{|PN|} = \frac{|\text{Green Quarter}|}{|\text{Circle}|}$$

$$\text{Recall} = \frac{|CP|}{|P|} = \frac{|\text{Green Quarter}|}{|\text{Light Green Square}|}$$

$$F1 = \frac{2 * \text{Precision} * \text{Recall}}{\text{Precision} + \text{Recall}}$$

Evaluation of name and type prediction

Arch		Precision	Recall	F1
x86	Name	62.6	62.5	62.5
	Type	63.7	63.7	63.7
	Overall	63.1	63.1	63.1
x64	Name	63.5	63.1	63.3
	Type	74.1	73.4	73.8
	Overall	68.8	68.3	68.6
ARM	Name	61.6	61.3	61.5
	Type	66.8	68.0	67.4
	Overall	64.2	64.7	64.5

Malware inspection

Inspected 35 x86 malware from VirusShare

Manipulating DNS settings

```
int sub_80534BA() {  
    ...  
    if ( dword_8063320 <= 0 ) {  
        v1 = ("/etc/resolv.conf", 'r');  
        if (v1 || (v1 =  
            sub_8053B1("resolv.conf"))){  
            ...  
            ...  
        }  
    }  
}
```

DE BIN

```
int rfc1035_init_resolv() {  
    ...  
    if ( num_entries <= 0 ) {  
        v0 = ("/etc/resolv.conf", 'r');  
        if (v0 || (v1 =  
            fopen64("resolv.conf"))){  
            // code to read and  
            // manipulate DNS settings  
        }  
    }  
}
```

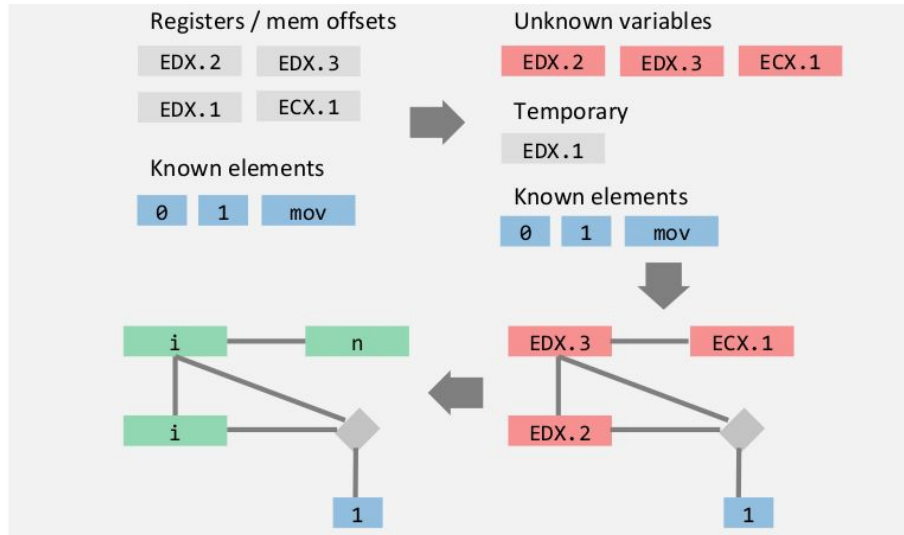
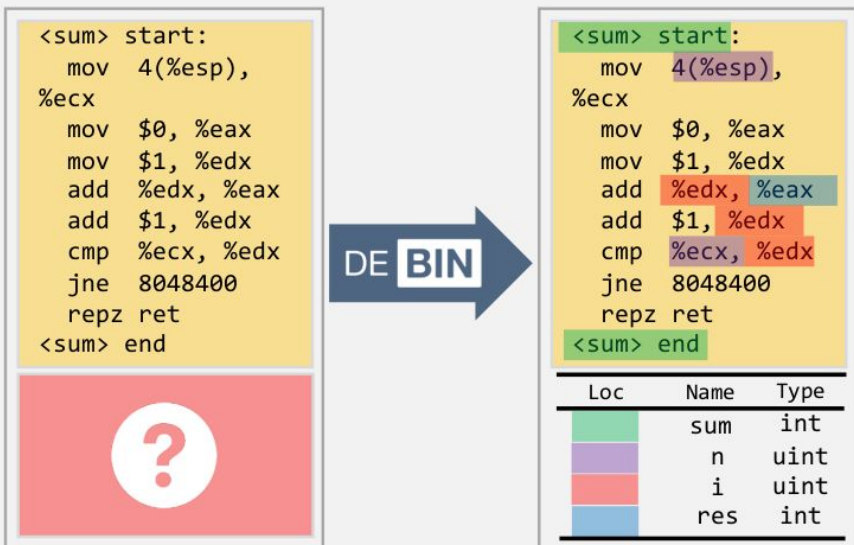
Leakage of sensitive data

```
If (sub_806d9f0(args) >= 0) {  
    ...  
    sub_80522B0(args);  
    ...  
}
```

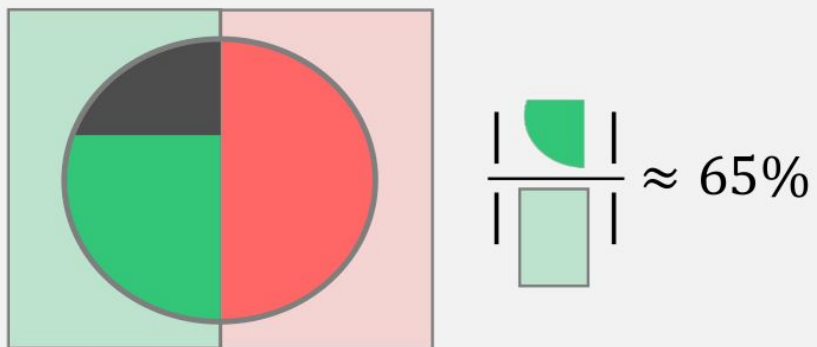
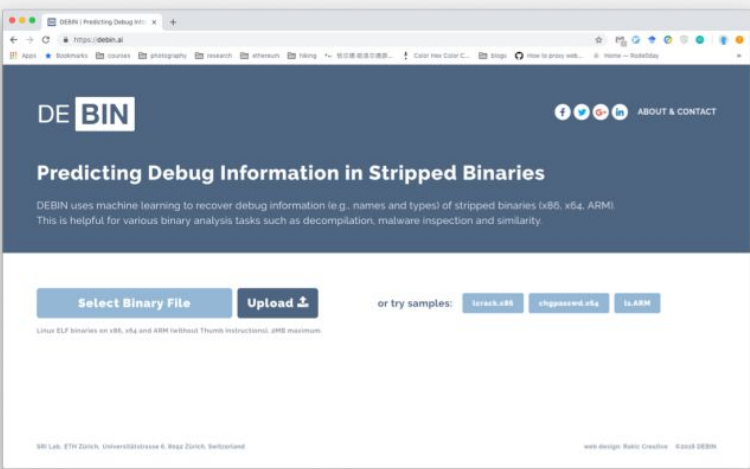
DE BIN

```
If (setsockopt(args) >= 0) {  
    ...  
    sendto(args);  
    ...  
}
```

Summary



Two-stage prediction process



Try online: <https://debin.ai>

High precision and accuracy

How can we improve?

Learning To Represent Programs with Graphs

Abhishek Shah

Problem

Neural Networks have understood:

- Images
- Speech
- Language
- *Source Code ?*

Problem

```
float getHeight { return this.width; }
```

Question: what's the bug?

Problem

Do what I **want**, not what I **wrote**

```
float getHeight { return this.width; }
```

Question: what's the bug?

Solution - Learning from “Big Code”

How to feed programs into Neural Networks?

- Sequence of Tokens (Hindle et al., 2012)
- Parse Tree (Bielik et al., 2016)

Key Insight:

- Expose semantics to NN via a Graph
 - Avoid shallow, textual structure by using data flow and type information

Outline

- Primer on Graph Neural Networks
- Converting Programs to Graphs
- Learning Representations with Graph NNs
- Downstream Tasks
- Evaluation

Primer on Graph NNs

- Why use Graphs?
 - Graphs describe a system and the complex dependencies within them
- Use Cases
 - Node Classification → is a node malicious?
 - Link Detection → are these two transactions linked in the blockchain?

Primer on Graph NNs

- Modern DL Techniques
 - CNNs → *fixed-size* images with *spatial locality*
 - RNNs → *ordered* sequences

Primer on Graph NNs

- Modern DL Techniques
 - CNNs → ~~fixed-size~~ images with ~~spatial locality~~
 - RNNs → ~~ordered~~ sequences
- Properties of Graphs
 - No obvious ordering
 - Not fixed sizes
 - Non-obvious or non-existent spatial locality

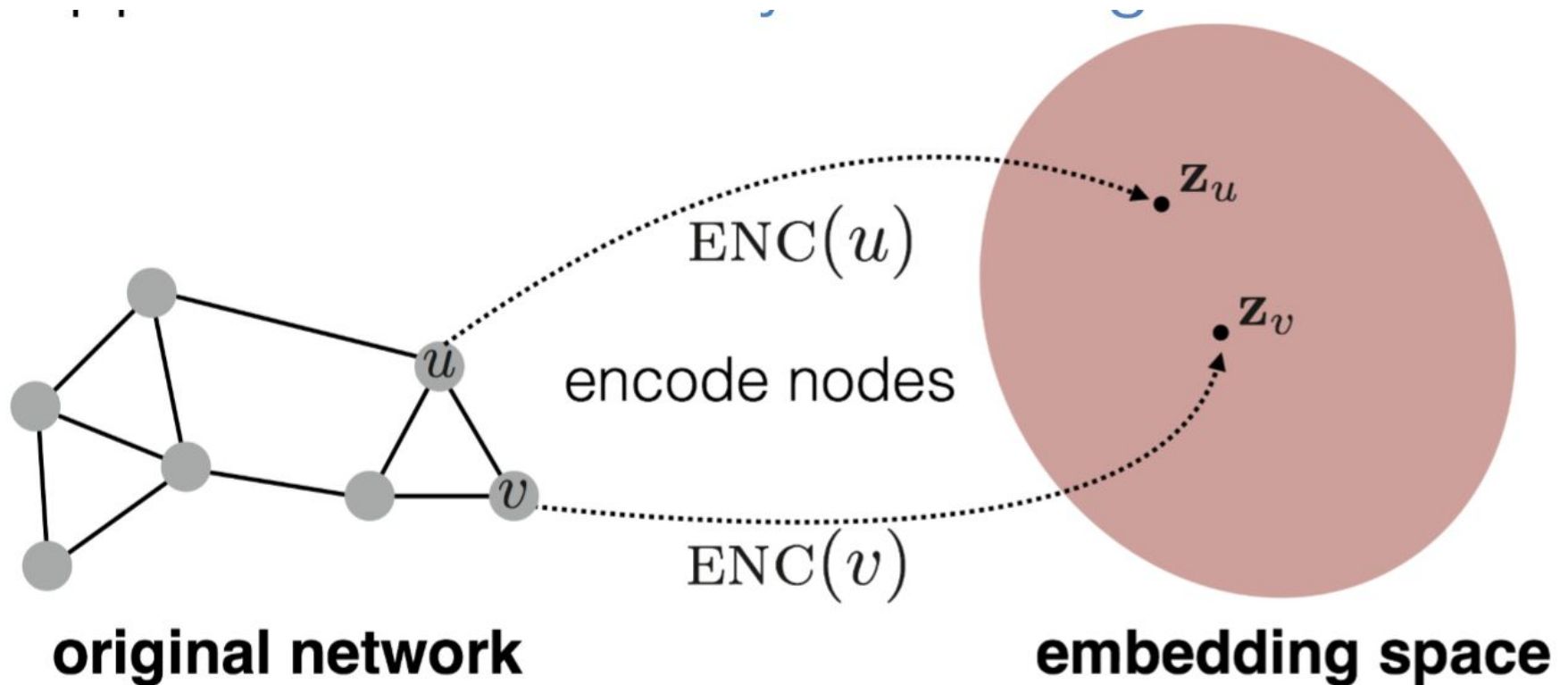
Primer on Graph NNs

- Building a Graph NN (focus on embedding)
 - Need an encoder
 - Such that similarity in original graph is preserved in embedded space

Primer on Graph NNs

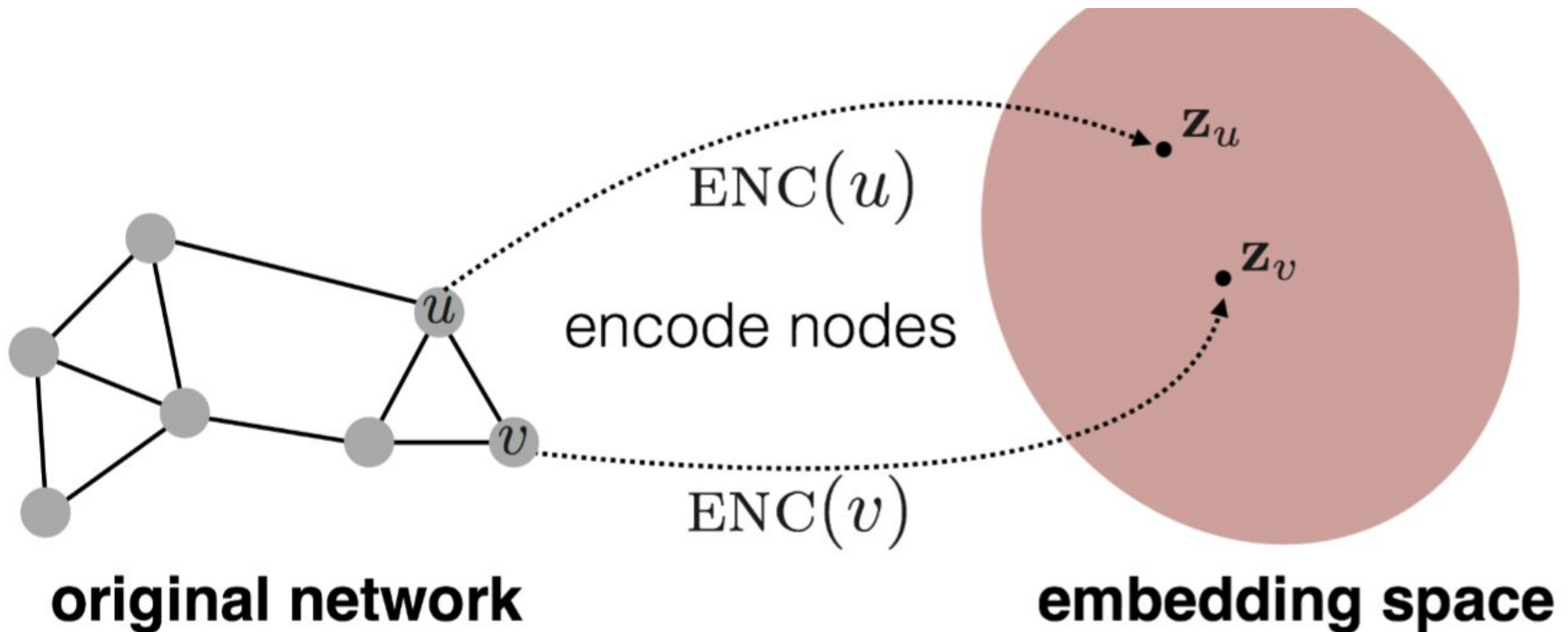
- Building a Graph NN (focus on embedding)
 - Need an encoder
 - Such that similarity in original graph is preserved in embedded space
 - Need a similarity metric
 - Learning → minimizing the distance between similar nodes

Primer on Graph NNs



Primer on Graph NNs

$$\text{similarity}(u, v) \approx \mathbf{z}_v^\top \mathbf{z}_u$$



Primer on Graph NNs

- For now, shallow encoding
 - Each node has a unique vector (“embedding-lookup”)

Primer on Graph NNs

- For now, shallow encoding
 - Each node has a unique vector (“embedding-lookup”)
- Similarity
 - Connected? or Share Neighbors?

Primer on Graph NNs

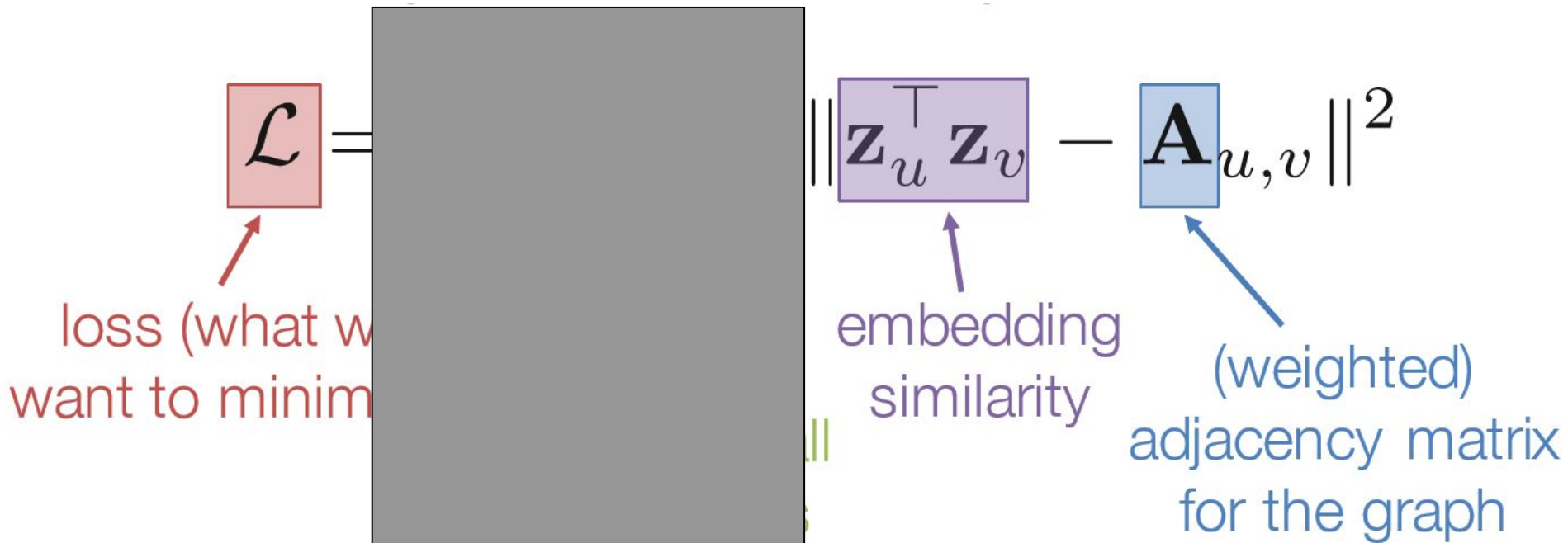
- For now, shallow encoding
 - Each node has a unique vector (“embedding-lookup”)
- Similarity
 - Connected? or Share Neighbors?
 - One Idea: dot products between node embeddings \sim edge existence

Primer on Graph NNs

- For now, shallow encoding
 - Each node has a unique vector (“embedding-lookup”)
- Similarity
 - Connected? or Share Neighbors?
 - One Idea: dot products between node embeddings \sim edge existence
 - Adjacency Matrix defines ground truth for edge existence
 - Take the difference between the two

Primer on Graph NNs

- Similarity



Primer on Graph NNs

- Similarity

$$\mathcal{L} = \sum_{(u,v) \in V \times V} \| \mathbf{z}_u^\top \mathbf{z}_v - \mathbf{A}_{u,v} \|^2$$

loss (what we want to minimize)

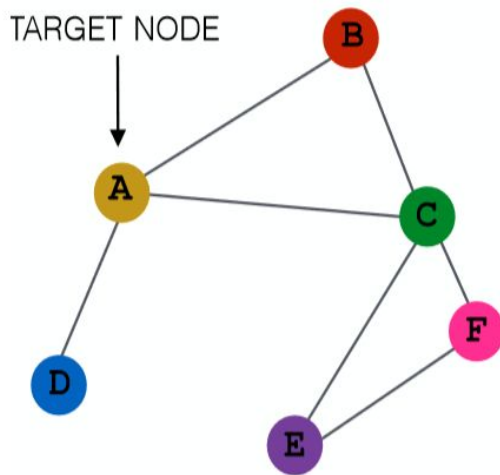
sum over all node pairs

embedding similarity

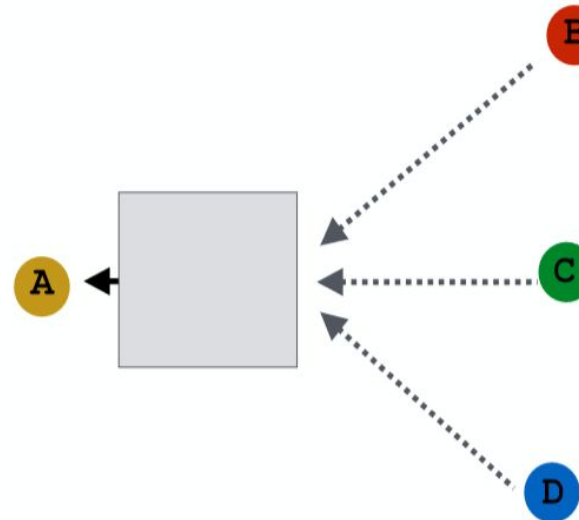
(weighted) adjacency matrix for the graph

Primer on Graph NNs

- Encoder
 - Main insight: generate node embeddings based on local neighborhoods

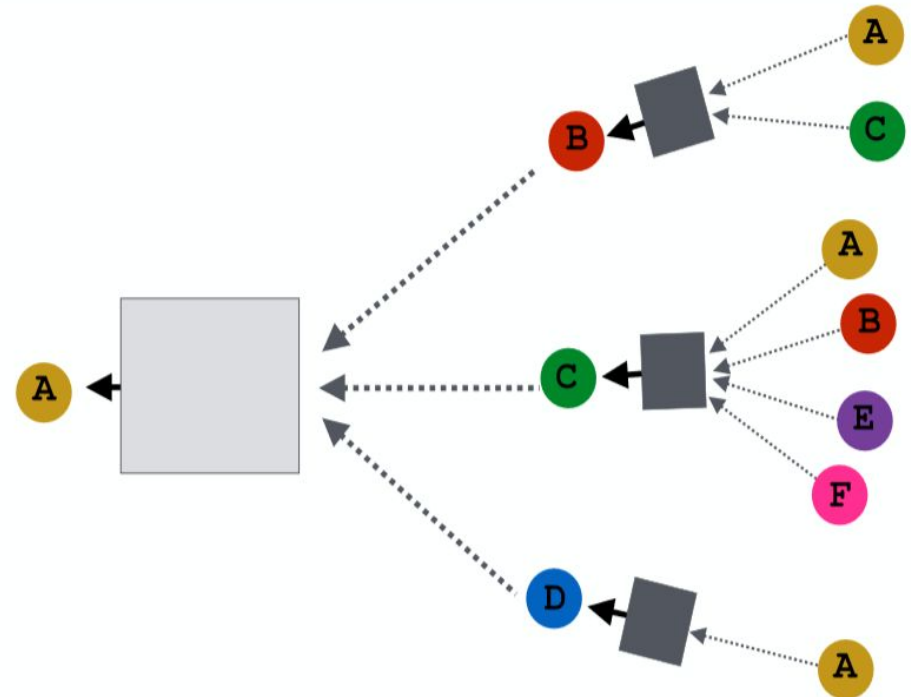
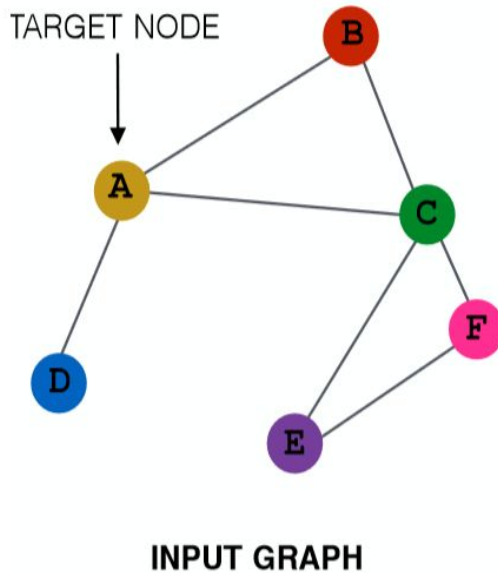


INPUT GRAPH



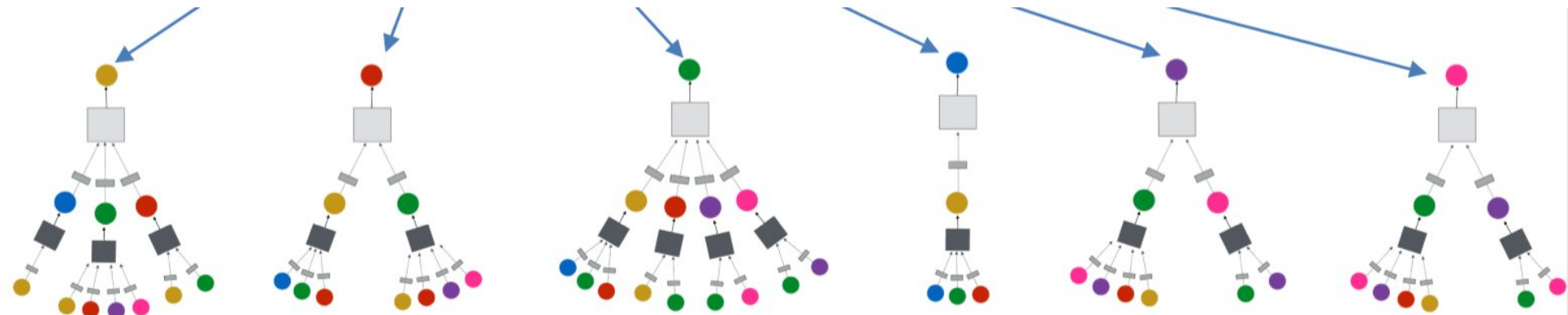
Primer on Graph NNs

- Encoder
 - Main insight: generate node embeddings based on local neighborhoods
 - NNs to aggregate information (msg) per layer



Primer on Graph NNs

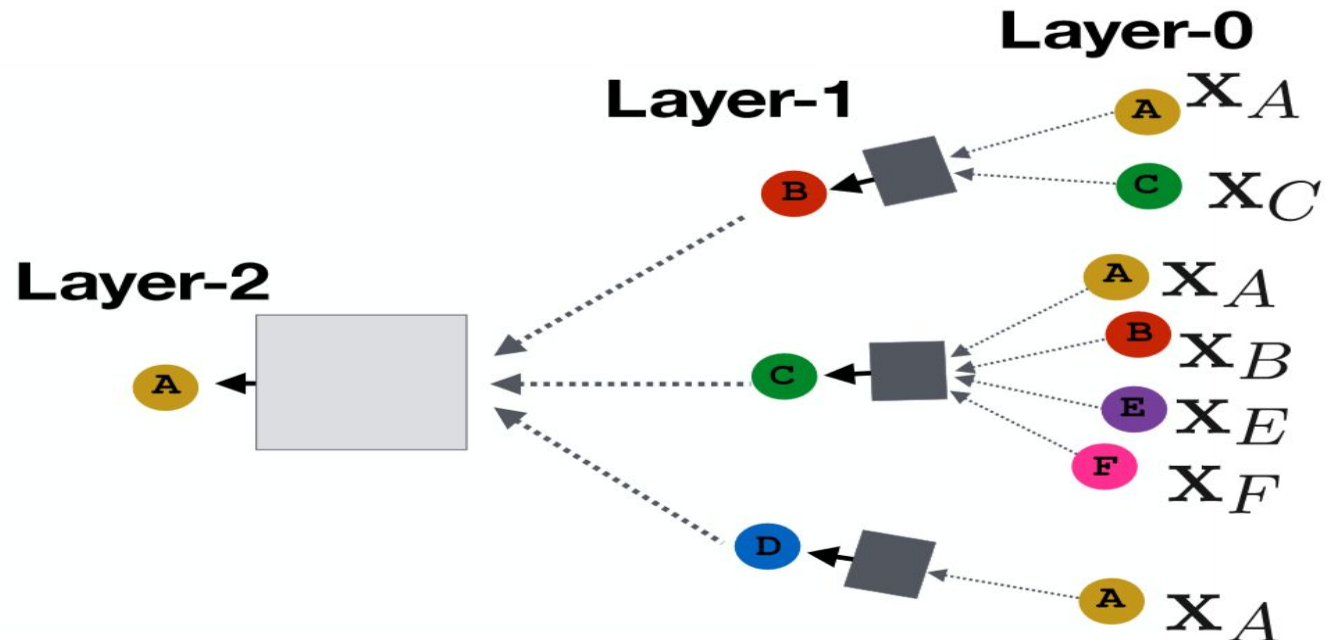
- “Deep” Encoder
 - Main insight: generate node embeddings based on local neighborhoods
 - NNs to aggregate information (msg) per layer
 - Each node has unique computation graph



Primer on Graph NNs

- Setup

- Graph $G = (V, A, X)$
 - $V \rightarrow$ Vertex Set
 - $A \rightarrow$ Adjacency Matrix
 - $X \rightarrow$ matrix of node features
 - Name, id, relationship status
- Layer 0 embedding \rightarrow input feature vector



Primer on Graph NNs

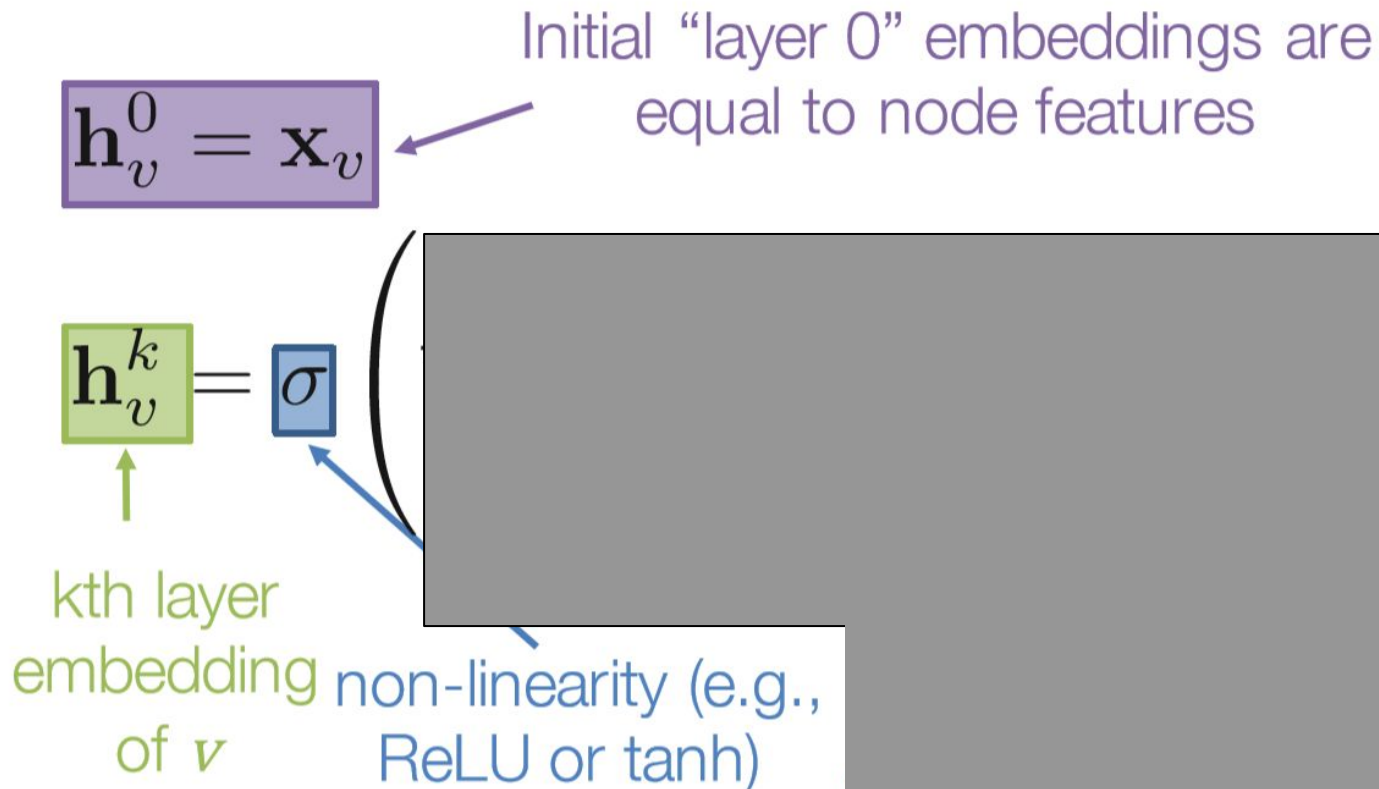
- **Basic approach:** Average neighbor messages and apply a neural network.

$$\mathbf{h}_v^0 = \mathbf{x}_v$$

Initial “layer 0” embeddings are equal to node features

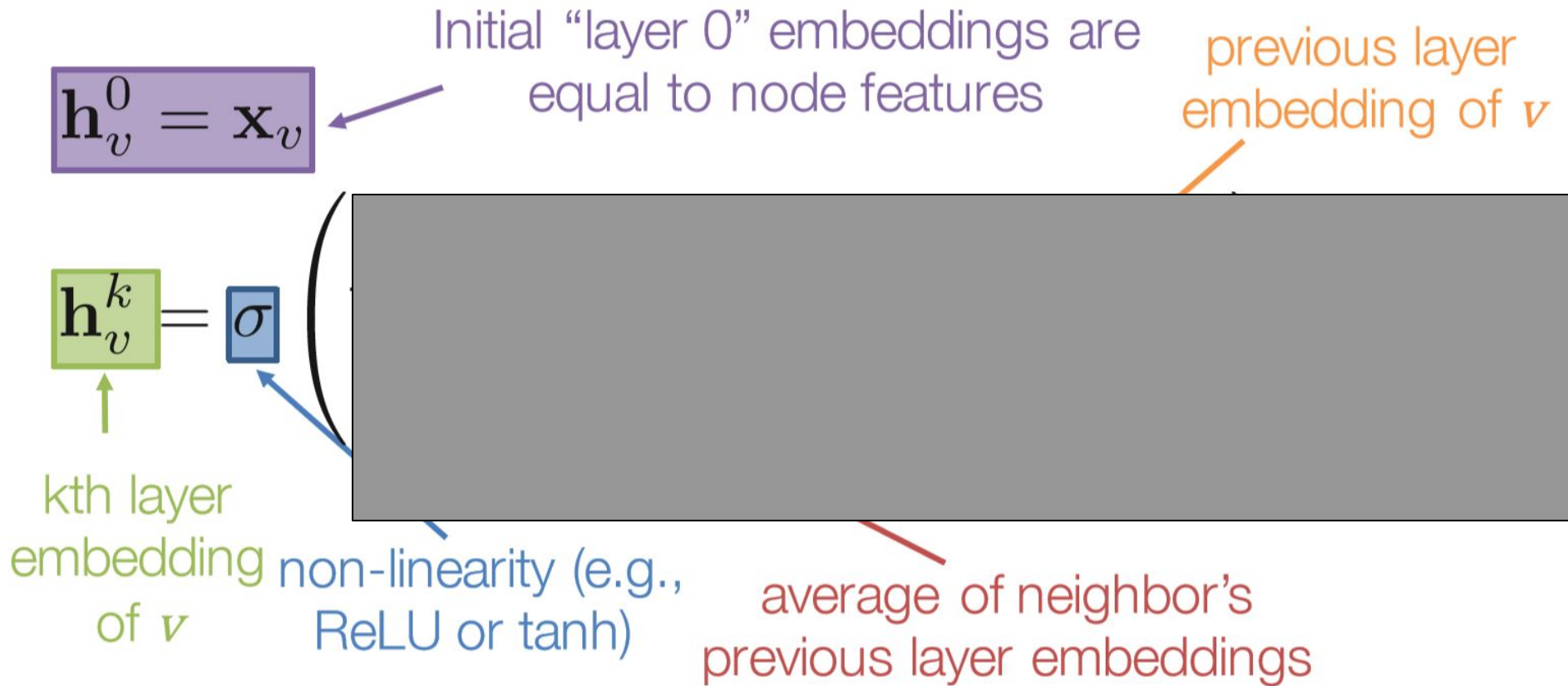
Primer on Graph NNs

- **Basic approach:** Average neighbor messages and apply a neural network.



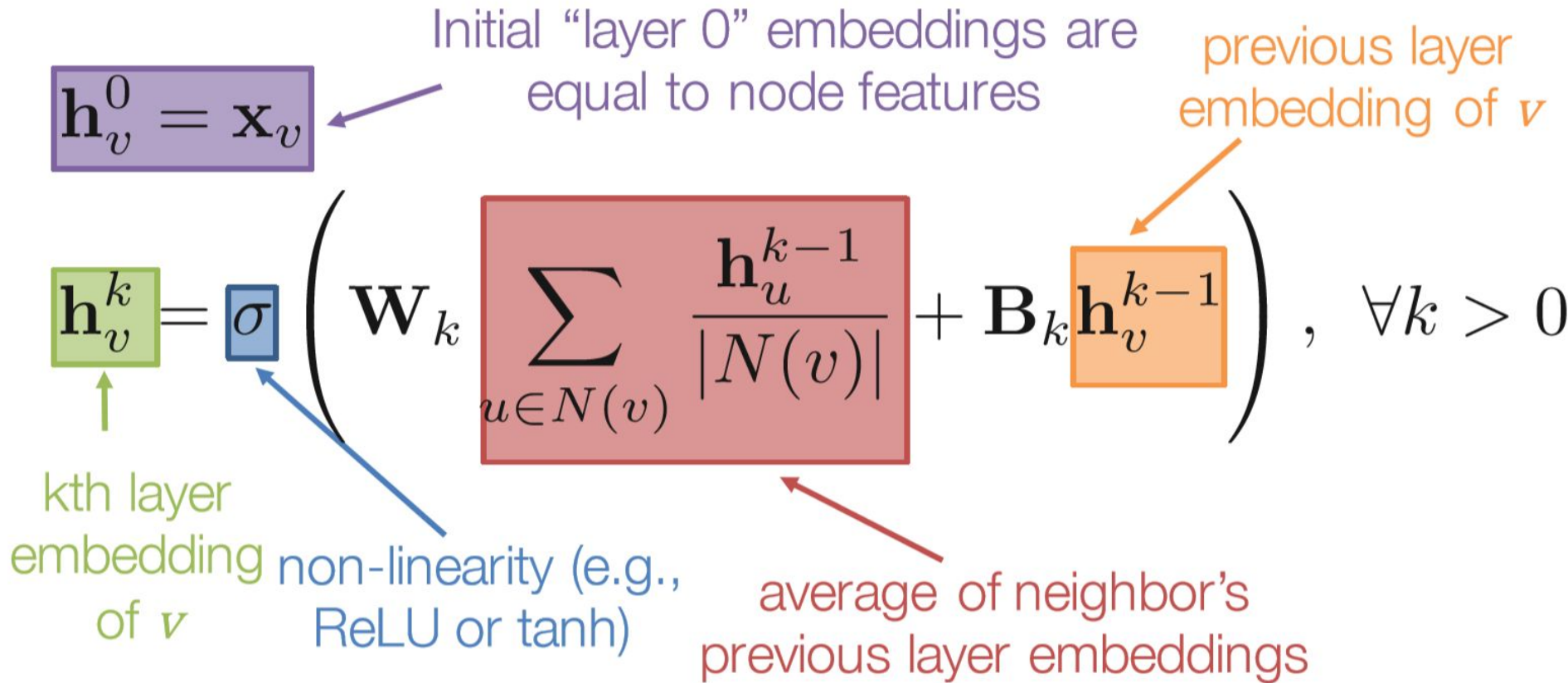
Primer on Graph NNs

- **Basic approach:** Average neighbor messages and apply a neural network.



Primer on Graph NNs


- **Basic approach:** Average neighbor messages and apply a neural network.



Primer on Graph NNs

$\mathbf{h}_v^0 = \mathbf{x}_v$

trainable matrices
(i.e., what we learn)

$$\mathbf{h}_v^k = \sigma \left(\mathbf{W}_k \sum_{u \in N(v)} \frac{\mathbf{h}_u^{k-1}}{|N(v)|} + \mathbf{B}_k \mathbf{h}_v^{k-1} \right), \quad \forall k \in \{1, \dots, K\}$$


Primer on Graph NNs

trainable matrices
(i.e., what we learn)

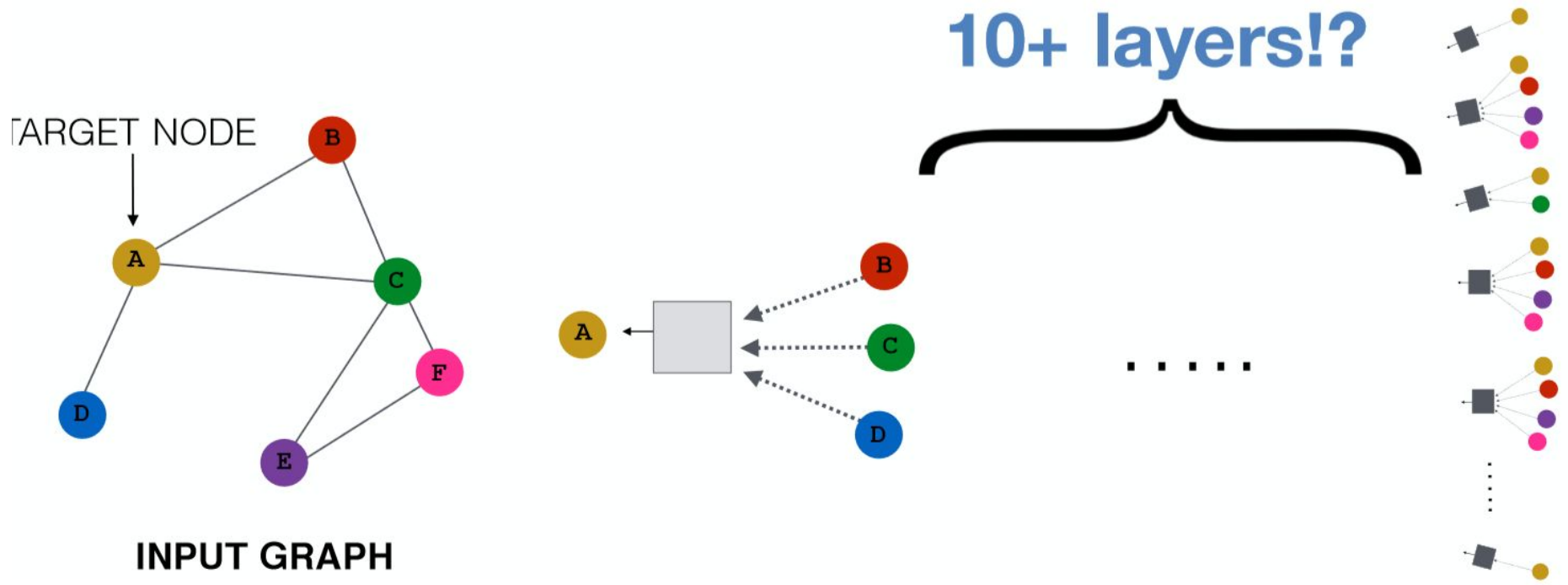
$$\mathbf{h}_v^0 = \mathbf{x}_v$$
$$\mathbf{h}_v^k = \sigma \left(\mathbf{W}_k \sum_{u \in N(v)} \frac{\mathbf{h}_u^{k-1}}{|N(v)|} + \mathbf{B}_k \mathbf{h}_v^{k-1} \right), \quad \forall k \in \{1, \dots, K\}$$

$$\mathbf{z}_v = \mathbf{h}_v^K$$

- After K -layers of neighborhood aggregation, we get output embeddings for each node.

Primer on Graph NNs

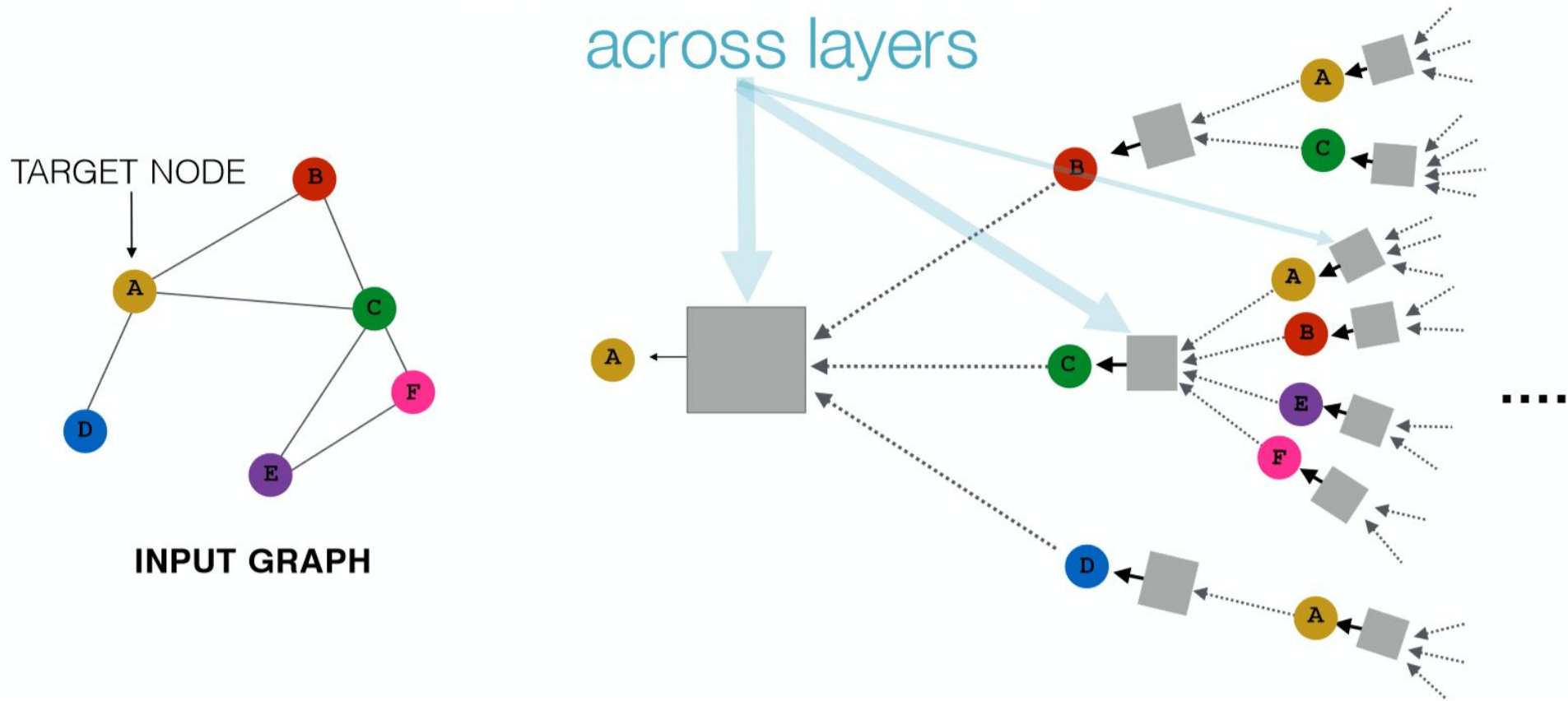
- What if we want to go deeper?
 - Overfitting from parameters



Primer on Graph NNs

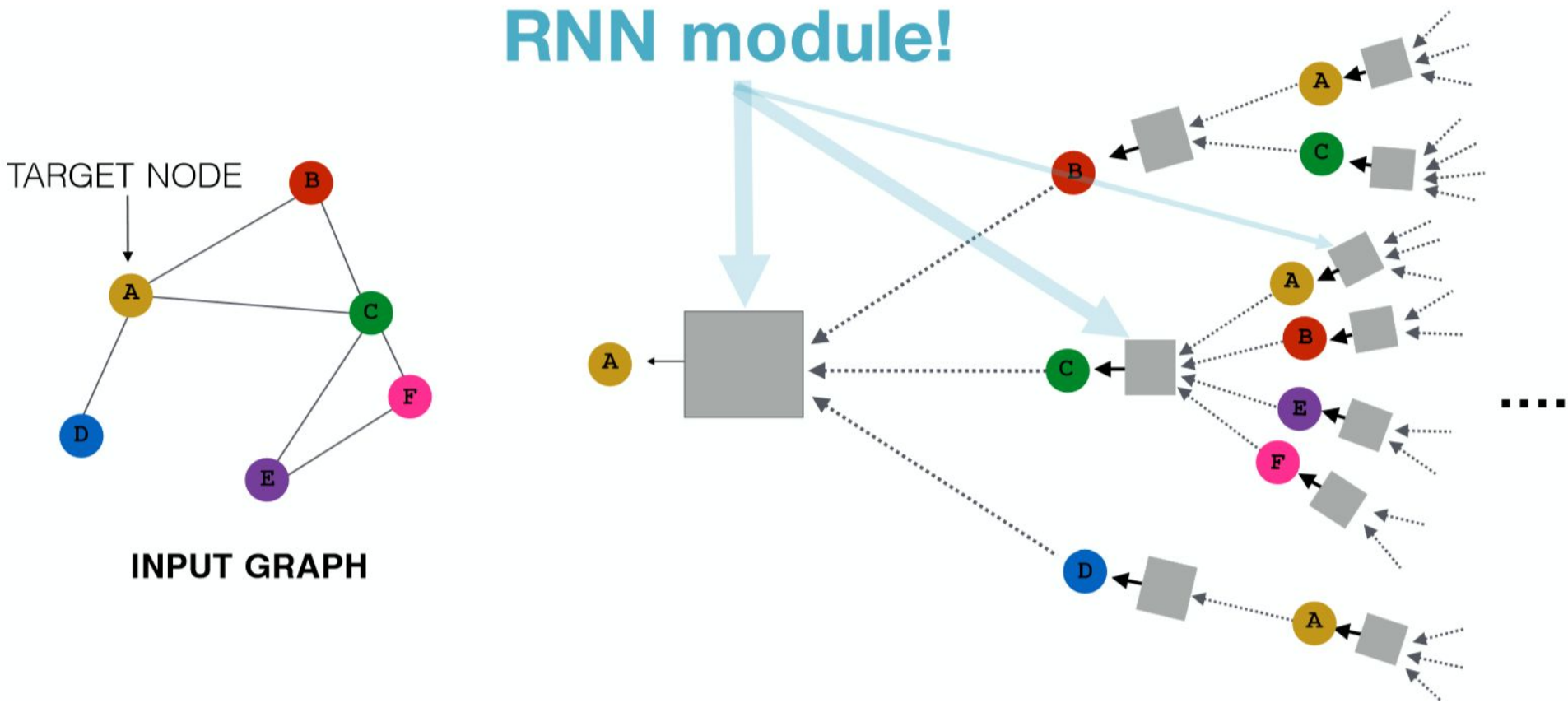
- **Idea 1:** Parameter sharing across layers.

same neural network
across layers



Primer on Graph NNs

- **Idea 2:** Recurrent state update.



Gated Graph NN

Intuition: Neighborhood aggregation with RNN state update.

1. Get “message” from neighbors at step k :

$$\mathbf{m}_v^k = \mathbf{W} \sum_{u \in N(v)} \mathbf{h}_u^{k-1}$$

← aggregation function does not depend on k

Gated Graph NN

Intuition: Neighborhood aggregation with RNN state update.

1. Get “message” from neighbors at step k :

$$\mathbf{m}_v^k = \mathbf{W} \sum_{u \in N(v)} \mathbf{h}_u^{k-1}$$

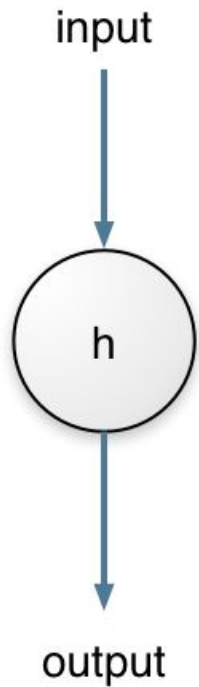
← aggregation function does not depend on k

2. Update node “state” using Gated Recurrent Unit (GRU). New node state depends on the old state and the message from neighbors:

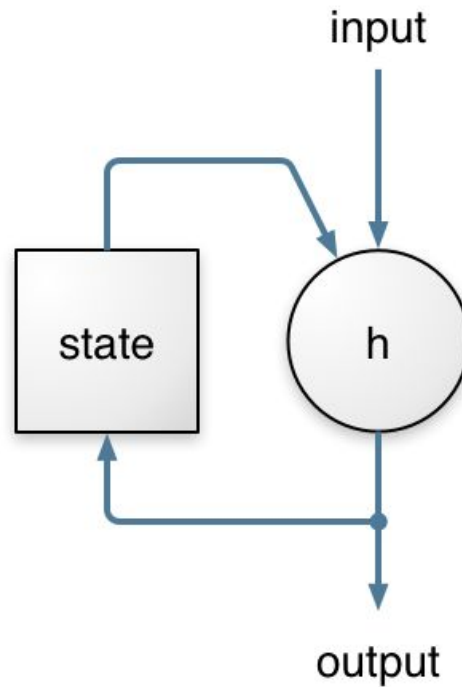
$$\mathbf{h}_v^k = \text{GRU}(\mathbf{h}_v^{k-1}, \mathbf{m}_v^k)$$

Gated Graph NN

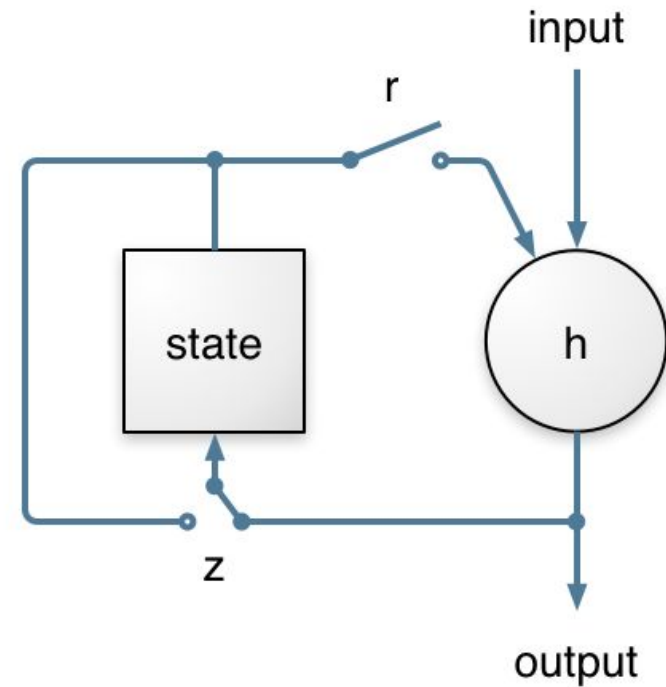
Feed-forward unit



Simple recurrent unit



Gated recurrent unit (GRU)



Outline

- Primer on Graph Neural Networks
- Converting Programs to Graphs
- Learning with Graph NNs
- Downstream Tasks
- Evaluation

Converting Programs to Graphs

Key Insight:

- Expose semantics to NN via a Graph
 - Avoid shallow, textual structure by using data flow and type information

Published as a conference paper at ICLR 2018

LEARNING TO REPRESENT PROGRAMS WITH GRAPHS

Miltiadis Allamanis
Microsoft Research
Cambridge, UK
miallama@microsoft.com

Marc Brockschmidt
Microsoft Research
Cambridge, UK
mabrocks@microsoft.com

Mahmoud Khademi*
Simon Fraser University
Burnaby, BC, Canada
mkhademi@sfu.ca

Converting Programs to Graphs

Graph: (V, E, X)

- V (AST nodes)
 - Grammar-Rule-Named Nonterminals
 - Named Program Tokens
- E
 - Syntactic
 - Semantic
- Discussion:
 - What are examples of syntactic and semantic edges?

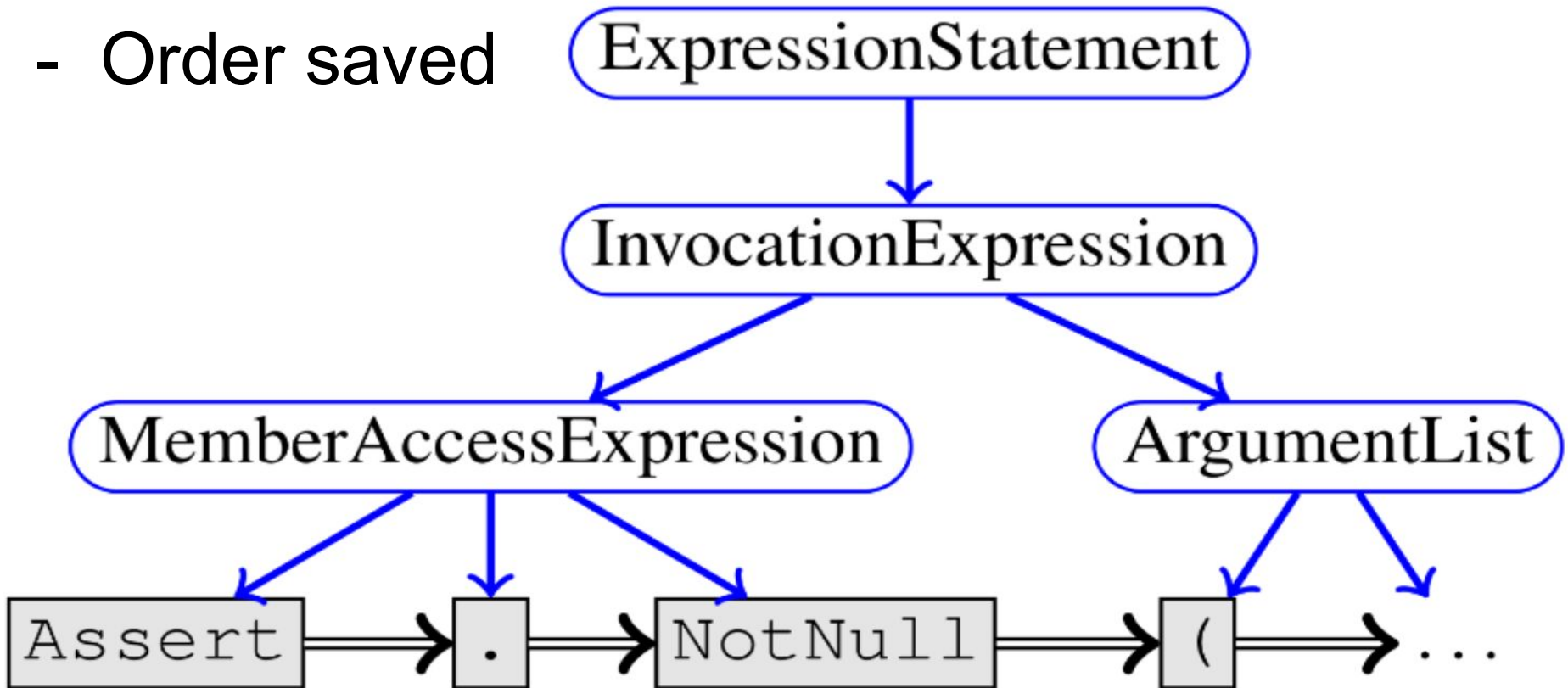
Converting Programs to Graphs

Syntactic Edges

Blue → Children

Black → NextToken

- Order saved



Converting Programs to Graphs

Semantic Edges

```
x, y = Foo();  
while (x > 0)  
    x = x + y;
```

- Let's focus on y at line 3

Converting Programs to Graphs

Semantic Edges

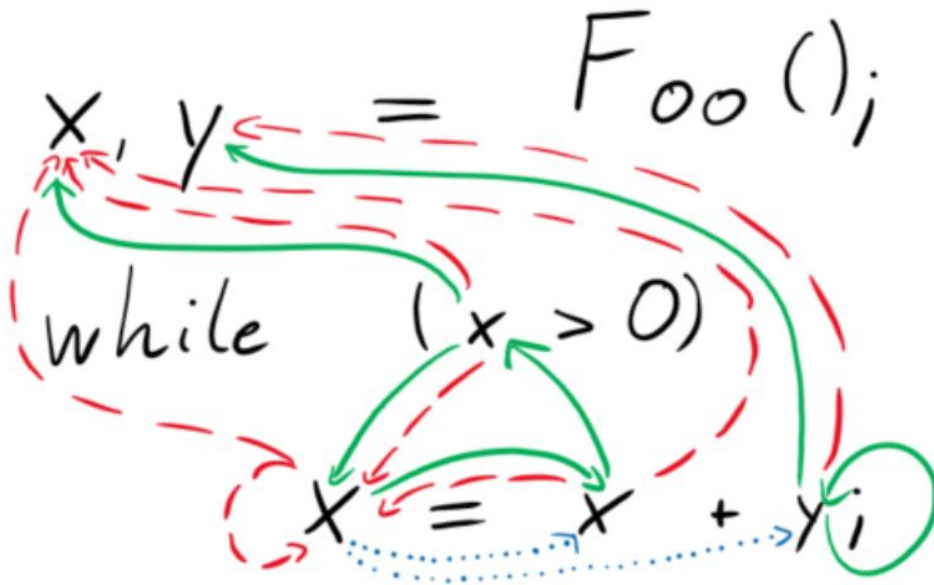
```
x, y = Foo();  
while (x > 0)  
    x = x + y;
```

- LastUse/Read(y_3) \rightarrow Line {1, 3}
 - Line 3 due to loop
- LastWrite(y_3) \rightarrow Line 1

Converting Programs to Graphs

Semantic Edges

```
x, y = Foo();  
while (x > 0)  
    x = x + y;
```



\longrightarrow Last Use
 $- - \longrightarrow$ Last Write
 $\cdots \longrightarrow$ Computed From

Converting Programs to Graphs

- Other Edges
 - Can use any other program analysis
 - Points-to analysis
 - Formal Parameter \leftrightarrow Argument Match
 - Conditional Guards
 - ReturnsTo

Converting Programs to Graphs

Variable Type Information

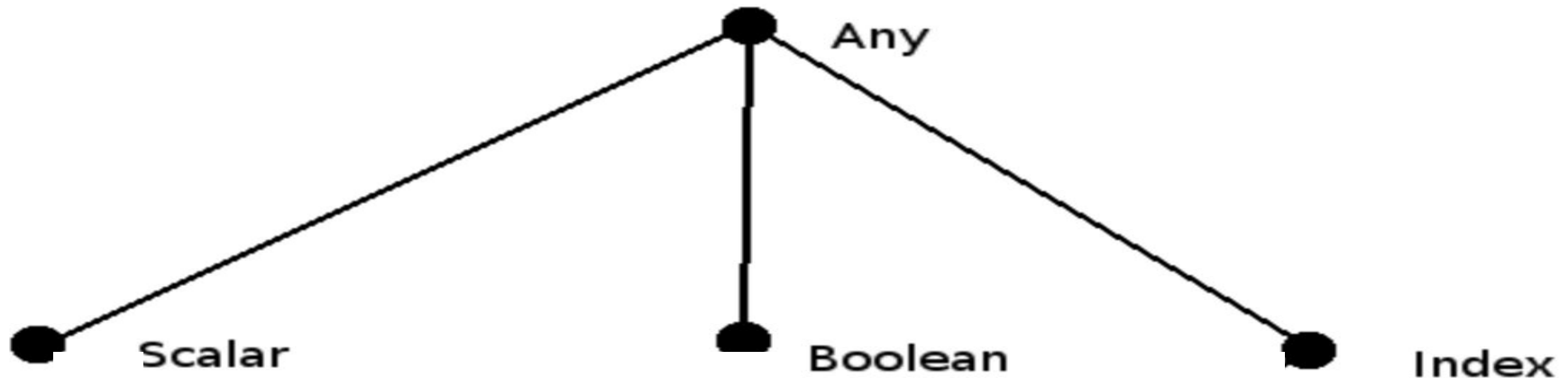
- Map variable type to max of set of supertypes
- $\text{List}\langle\text{int}\rangle \rightarrow \max(\{\text{List}\langle\text{int}\rangle, \text{List}\langle K\rangle\})$

Discussion: any flaws?

Converting Programs to Graphs

Variable Type Information

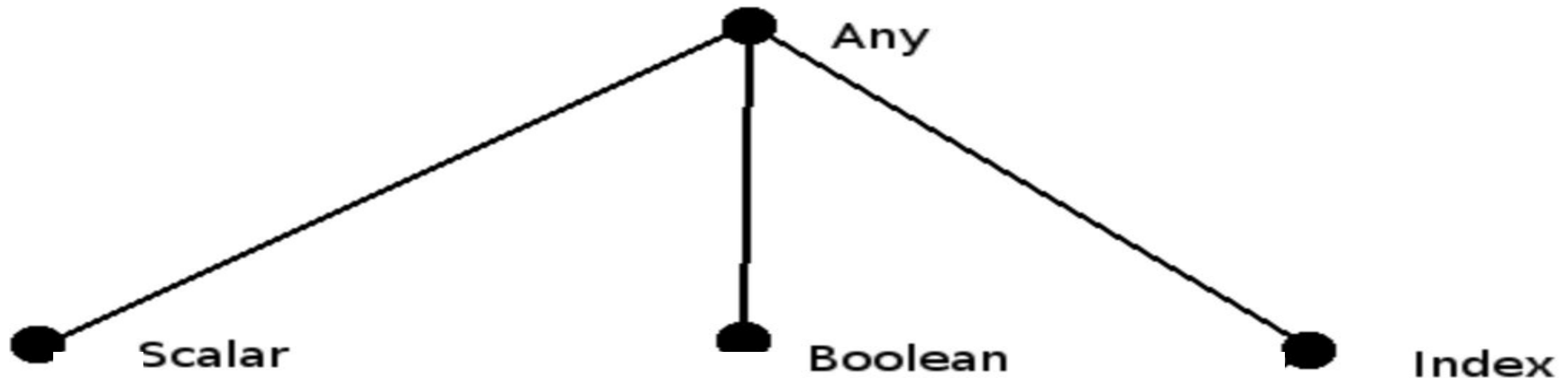
- Map variable type to max of set of supertypes
- Boolean $\rightarrow \max(\{\text{Boolean}, \text{Any}\}) \rightarrow \text{Any}$
- Scalar $\rightarrow \max(\{\text{Scalar}, \text{Any}\}) \rightarrow \text{Any}$



Converting Programs to Graphs

Variable Type Information

- Use dropout mechanisms: randomly select subset
- Boolean $\rightarrow \max(\{\text{Any}\}) \rightarrow \text{Any}$
- Scalar $\rightarrow \max(\{\text{Scalar}\}) \rightarrow \text{Scalar}$



Learning with Graph NNs

- $T = 0$ (Initial Node Representation)
 - Concatenate Name with Type string embedding
- Run Gated Graph NN propagation for 8 steps
 - 8 was experimentally determined

Downstream Tasks

- We have an embedding... now what?

Downstream Task 1 - VarNaming

```
1 private void ██████████ () {  
2     String vertexShader = "literal_1";  
3     String fragmentShader = "literal_2";  
4     shader = new ShaderProgram(vertexShader,  
5         fragmentShader);  
6     if(shader.isCompiled() == false)  
7         throw new IllegalArgumentException(  
8             "literal_3" + shader.getLog());  
9 }
```

Downstream Task 1 - VarNaming

- Goal: predict correct name at slot t
- Edit Graph
 - Insert new node at slot t (“hole”)

Downstream Task 1 - VarNaming

- Goal: predict correct name at slot t
- Edit Graph
 - Insert new node at slot t (“hole”)
 - Run Gated Graph NN for 8 steps
 - Feed representation into trained GRU to predict name as a sequence

Downstream Task 2 - VarMisuse

- Found several real-world bugs

```
var clazz=classTypes["Root"].Single() as JsonCodeGenerator.ClassType;
Assert.NotNull(clazz);

var first=classTypes["RecClass"].Single() as JsonCodeGenerator.ClassType;
Assert.NotNull( clazz );

Assert.Equal("string", first.Properties["Name"].Name);
Assert.False(clazz.Properties["Name"].IsArray);
```

Downstream Task 2 - VarMisuse

- Goal: predict correct token at slot t
 - Only type-correct tokens allowed at slot t
- Edit Graph
 - Insert new node at slot (“hole”)

Downstream Task 2 - VarMisuse

- Goal: predict correct token at slot t
 - Only type-correct tokens allowed at slot t
- Edit Graph
 - Insert new node at slot (“hole”)
 - Connect it without node v -dependent edges \rightarrow *context (i.e. $c(t)$)*

Downstream Task 2 - VarMisuse

- Goal: predict correct token at slot t
 - Only type-correct tokens allowed at slot t
- Edit Graph
 - Insert new node at slot (“hole”)
 - Connect it without node v -dependent edges \rightarrow *context (i.e. $c(t)$)*
 - Connect it with node v -dependent edges \rightarrow *usage representation (i.e. $u(t, v)$)*
 - Edges include LastUse and LastWrite
 - Add usage node per type-correct variable

Downstream Task 2 - VarMisuse

- Goal: predict correct token at slot t
 - Only type-correct tokens allowed at slot t
- Edit Graph
 - Insert new node at slot (“hole”)
 - Connect it without node v -dependent edges \rightarrow *context (i.e. $c(t)$)*
 - Connect it with node v -dependent edges \rightarrow *usage representation (i.e. $u(t, v)$)*
 - Edges include LastUse and LastWrite
 - Add usage node per type-correct variable
 - Run Gated Graph NN for 8 steps
- Correct Variable Usage
 - Node v that maximizes trained $W(c(t), u(t, v))$

Evaluation

- Dataset
 - 29 C# projects (~3 million lines of code)
 - Graphs on average: ~2300 nodes, ~8400 edges
- Baseline
 - VarMisuse (predict variable usage)
 - LOC → 2 layer bidirectional GRU
 - AVGB1RNN → LOC + simple variable usage dataflow

Evaluation

- Dataset
 - 29 C# projects (~3 million lines of code)
 - Graphs on average: ~2300 nodes, ~8400 edges
- Baseline
 - VarMisuse (predict variable usage)
 - LOC → 2 layer bidirectional GRU
 - AVGB1RNN → LOC + simple variable usage dataflow
 - VarNaming (predict name)
 - AVGLBL → Log-bilinear model (NLP-inspired)
 - AVGB1RNN (birectional RNN)

Evaluation

- LOC → captures little information
- AVGLBL/AVGB1RNN → captures some info
- Generalization --> unknown types/vocabulary

Table 1: Evaluation of models. UNSEENPROJTEST refers to projects that have no files in the training data, SEENPROJTEST refers to the test set containing projects that have files in the training set.

	SEENPROJTEST				UNSEENPROJTEST			
	Loc	AVGLBL	AVGB1RNN	GGNN	Loc	AVGLBL	AVGB1RNN	GGNN
VARMISUSE								
Accuracy (%)	15.8	—	73.5	82.1	13.8	—	59.7	68.6
PR AUC	0.363	—	0.931	0.963	0.363	—	0.891	0.909

Evaluation

- LOC → captures little information
- AVGLBL/AVGB1RNN → captures some info
- Generalization --> unknown types/vocabulary

Table 1: Evaluation of models. UNSEENPROJTEST refers to projects that have no files in the training data, SEENPROJTEST refers to the test set containing projects that have files in the training set.

	SEENPROJTEST				UNSEENPROJTEST			
	LOC	AVGLBL	AVGB1RNN	GGNN	LOC	AVGLBL	AVGB1RNN	GGNN
VARMISUSE								
Accuracy (%)	15.8	—	73.5	82.1	13.8	—	59.7	68.6
PR AUC	0.363	—	0.931	0.963	0.363	—	0.891	0.909
VARNAMING								
Accuracy (%)	—	22.0	25.5	30.7	—	15.3	15.9	19.4
F1 (%)	—	36.1	42.9	54.6	—	22.7	23.4	30.5

Evaluation

- Lacking semantic info hurts both
- Lacking syntactic info hurts VarMisuse

Table 2: Ablation study for the GGNN model on SEENPROJTEST for the two tasks.

Ablation Description	VARMISUSE Accuracy (%)	VARNAMING Accuracy (%)
Standard Model (reported in Table 1)	82.1	30.7
Only NextToken, Child, LastUse, LastWrite edges	79.0	15.4
Only semantic edges (all but NextToken, Child)	74.3	29.7
Only syntax edges (NextToken, Child)	49.6	20.5

Evaluation

- Only syntactic info impacts both
- Only semantic info impacts VarMisuse
- Node initial labeling impacts VarNaming

Table 2: Ablation study for the GGNN model on SEENPROJTEST for the two tasks.

Ablation Description	VARMISUSE Accuracy (%)	VARNAMING Accuracy (%)
Standard Model (reported in Table 1)	82.1	30.7
Only NextToken, Child, LastUse, LastWrite edges	79.0	15.4
Only semantic edges (all but NextToken, Child)	74.3	29.7
Only syntax edges (NextToken, Child)	49.6	20.5
Node Labels: Tokens instead of subtokens	82.1	16.8
Node Labels: Disabled	80.0	14.7

Contributions

- VarMisuse tasks and its practicality
- Learning Program Representations over Graphs

Questions/Discussion

- References

- [http://snap.stanford.edu/proj/embeddings-
www/](http://snap.stanford.edu/proj/embeddings-
www/)