SRILAB

ETHzürich

# DeBIN: Predicting Debug Information in Stripped Binaries
https://debin.ai

Jingxuan He

Pesho Ivanov

**Petar Tsankov**

Veselin Raychev

Martin Vechev

DEEPCODE

# Binaries with debug symbols

Assembly

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

Debug symbols

```
80534BA   rfc1035_init int
8053DB1   fopen64      int
8063320   num_entries  int
                ⋮
```
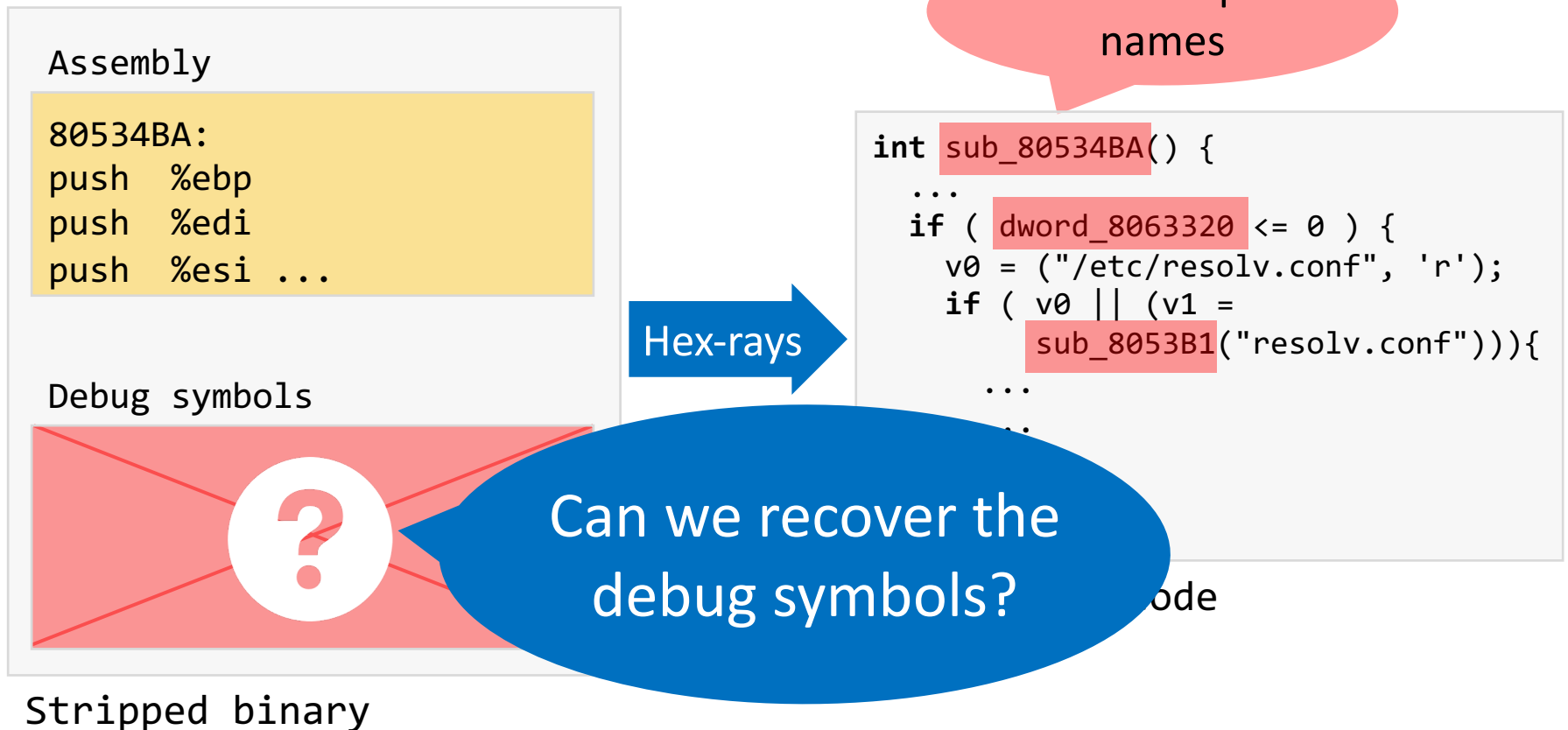
Binary with debug symbols

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

# Challenges

Stores the value of a semantic variable

Computes 1 + 2 + … + 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 intermediate (non-semantic) value

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

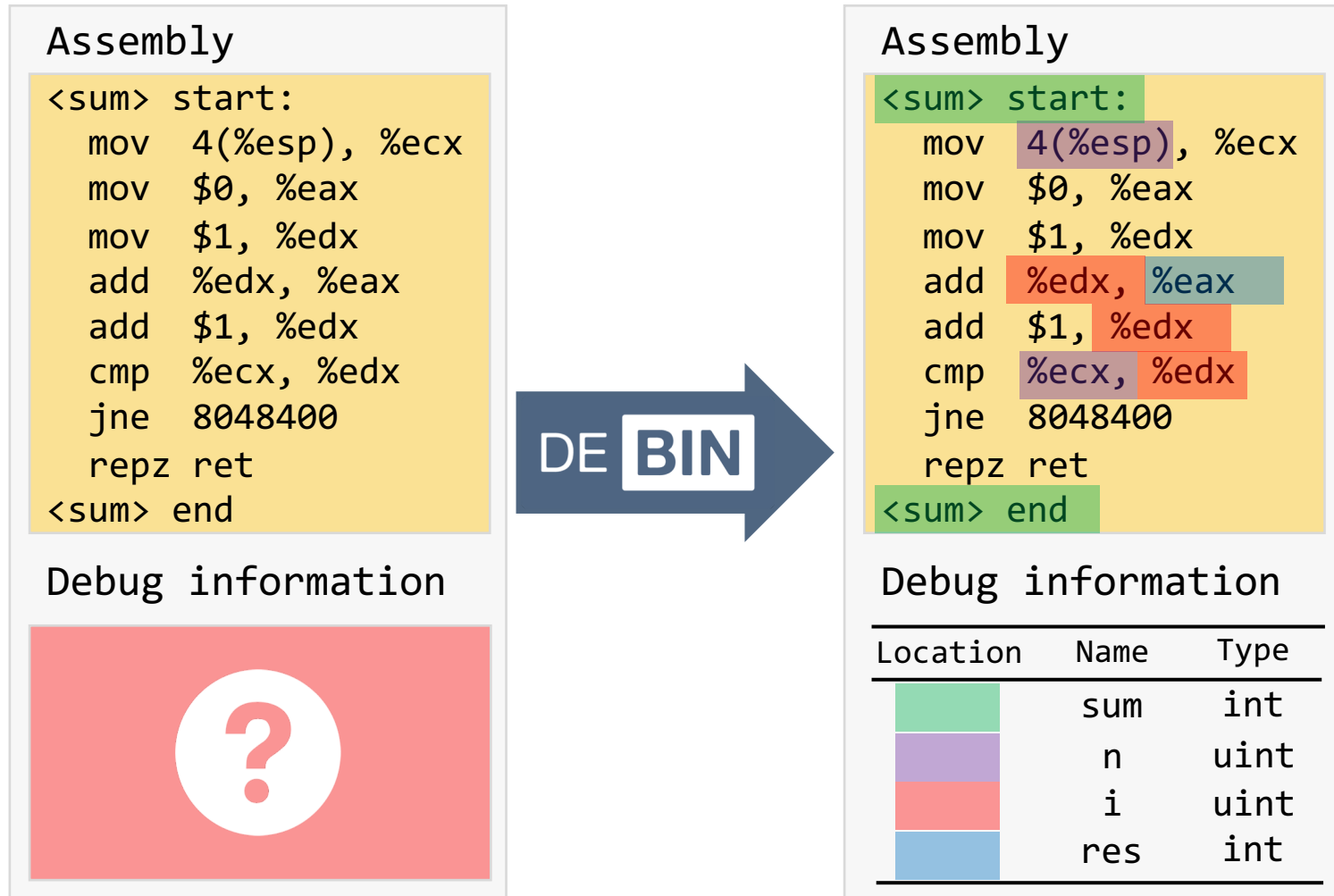# Challenges

```
<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

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

3

# DeBIN: Recovering debug information

DE BIN

ABOUT & CONTACT

# Predicting Debug Information in Stripped Binaries

DEBIN uses machine learning to recover debug information (e.g., names and types) of stripped binaries (x86, x64, ARM).
This is helpful for various binary analysis tasks like decompilation, malware inspection and similarity.

DEMO

Select Binary File    Upload ⬆

Linux ELF binaries on x86, x64 and ARM (without Thumb instructions), 2MB maximum.
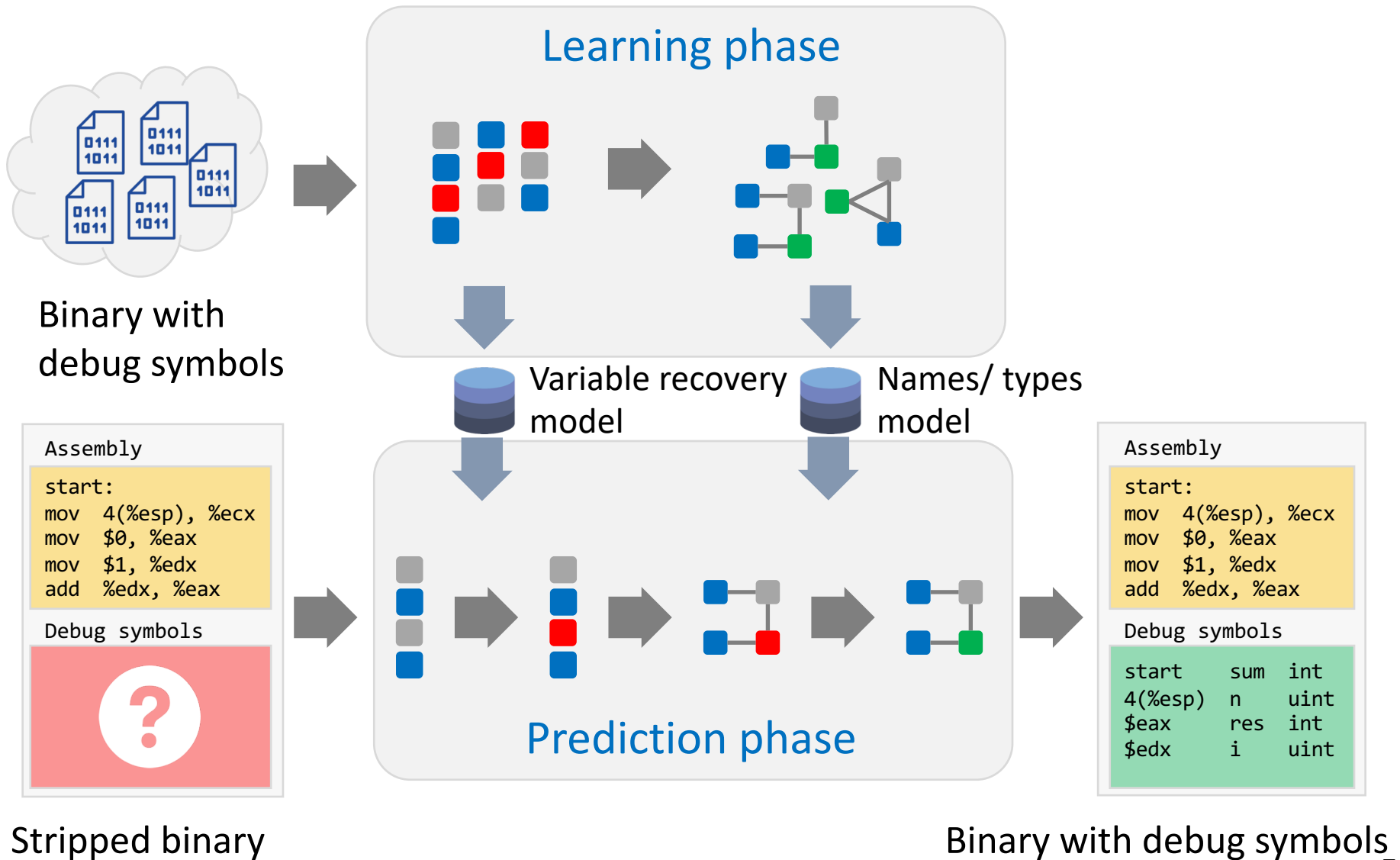
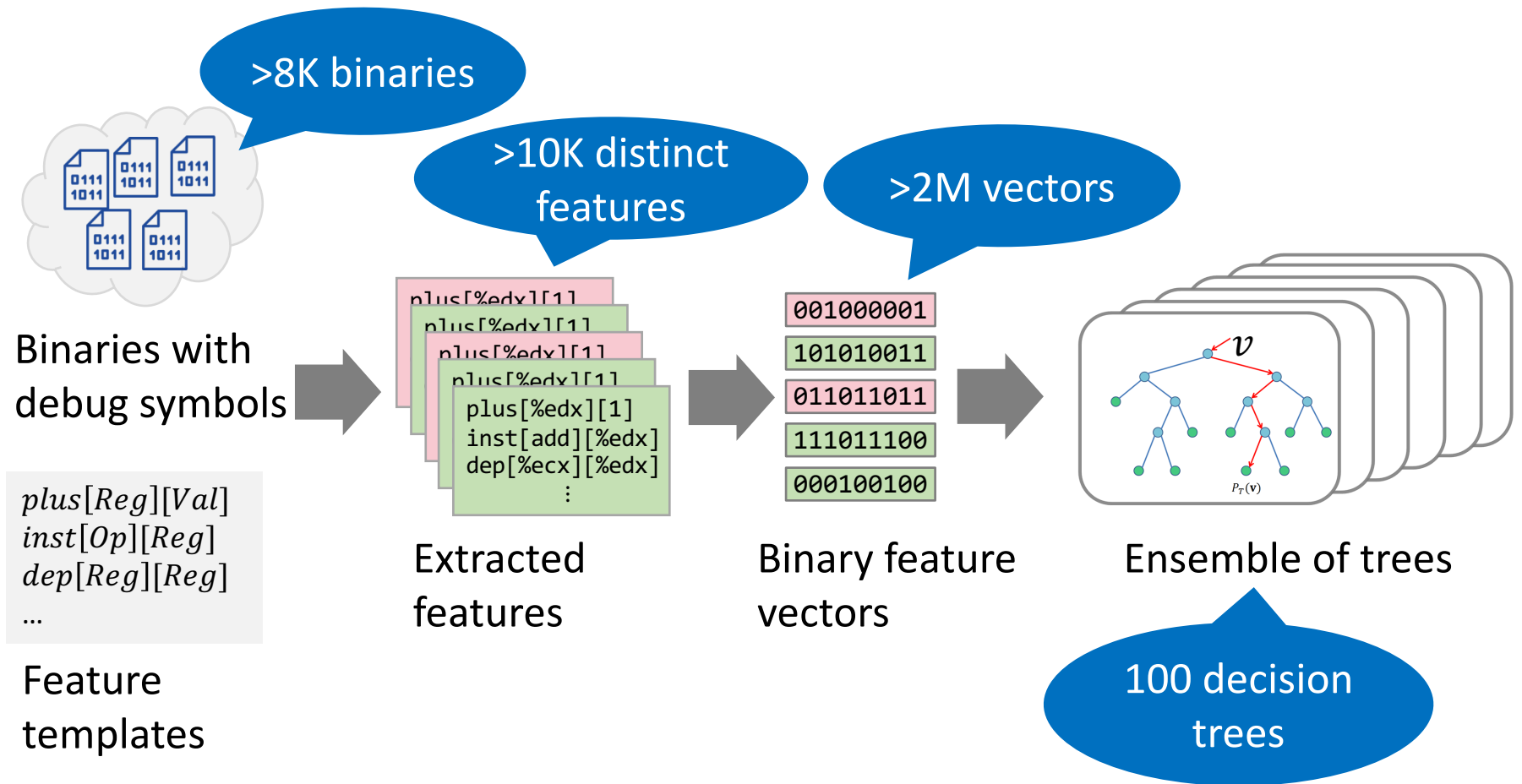or try samples:    lcrack.x86    chgpasswd.x64    ls.ARM

# How does DeBIN work?

# DeBIN: System overview

# Step 1: Recovering variables

# Learning how to recover variables



7

# Variable recovery

Assembly

```
mov    4(%esp), %ecx
mov    $0, %eax
mov    $1, %edx
add    %edx, %eax
add    $1, %edx.2
cmp    %ecx, %edx
jne    8048400
repz   ret
```

%edx.2

Register

```
plus[%edx][1]
inst[add][%edx]
⋮
```

Features

`00100101010001`  Feature vector $v$



Extremely randomized trees

**sem**
(DeBIN will predict name and type)

**tmp**
(stores an intermediate value)

*Extremely randomized trees*, Pierre Geurts, Damien Ernst, and Louis Wehenkel, Machine Learning 2006

6

# Step 2: Predicting names and types

# Probabilistic graphical model

| | EDX.3 | ECX.1 | weight |
|---|---|---|---|
| $f_1$ | i | n | 0.5 |
| $f_2$ | p | s | 0.3 |
| $f_3$ | a | b | 0.1 |

| | EDX.2 | EDX.3 | weight |
|---|---|---|---|
| $f_4$ | i | i | 0.8 |
| $f_5$ | i | j | 0.6 |
| $f_6$ | p | p | 0.3 |

cond-NE-EDX-ECX

EDX.3 ──── ECX.1

dep-EDX-EDX

EDX.2

| 1 | EDX.2 | EDX.3 | weight |
|---|---|---|---|
| 1 | i | i | 0.8 |
| 1 | j | i | 0.6 |
| 1 | p | p | 0.3 |

1

Known elements    1

Unknown elements    ECX.1, …

Binary features    $f_1, f_2, …$

Factors

# Probabilistic graphical model

| | |
|---|---|
| $f_4$ | i |
| $f_5$ | i |
| $f_6$ | p |

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

ght

.8

.6

.3

Know

Unkn

Binar

Factors

# Probabilistic graphical model

| | EDX.3 | ECX.1 | weight |
|---|---|---|---|
| $f_1$ | i | n | 0.5 |
| $f_2$ | p | s | 0.3 |
| $f_3$ | a | b | 0.1 |

| | EDX.2 | EDX.3 | weight |
|---|---|---|---|
| $f_4$ | i | i | 0.8 |
| $f_5$ | i | j | 0.6 |
| $f_6$ | p | p | 0.3 |

dep-EDX

**Next**

How are the features and
their weights learned?

| EDX.3 | weight |
|---|---|
| i | 0.8 |
| i | 0.6 |
| p | 0.3 |

1

Known elements    1

Unknown elements    ECX.1, …

Binary features    $f_1, f_2, …$
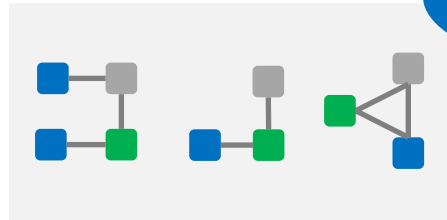
Factors

8

# Learning how to predict names and types



Dependency graphs

> 8,000 binaries

Actual graphs have >1K nodes

Binaries with debug symbols

Static analysis

Train model

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

Feature templates

23 templates

## 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 |

### 3-factor

| | | |
|---|---|---|
| 1 | i | i |
| 1 | j | i |
| 1 | p | p |

### 4-factor

| | | | |
|---|---|---|---|
| 1 | i | i | k |
| 1 | j | i | a |
| 1 | p | p | v |

Binary features and factors

| | 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 $(\vec{u_i}, \vec{k_i})$

9

# End-to-end recovery of debug information

# Recovering 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

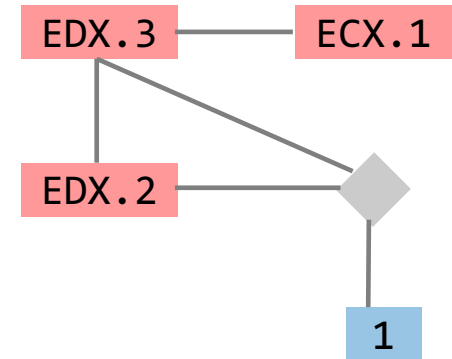Known elements

0    1    mov

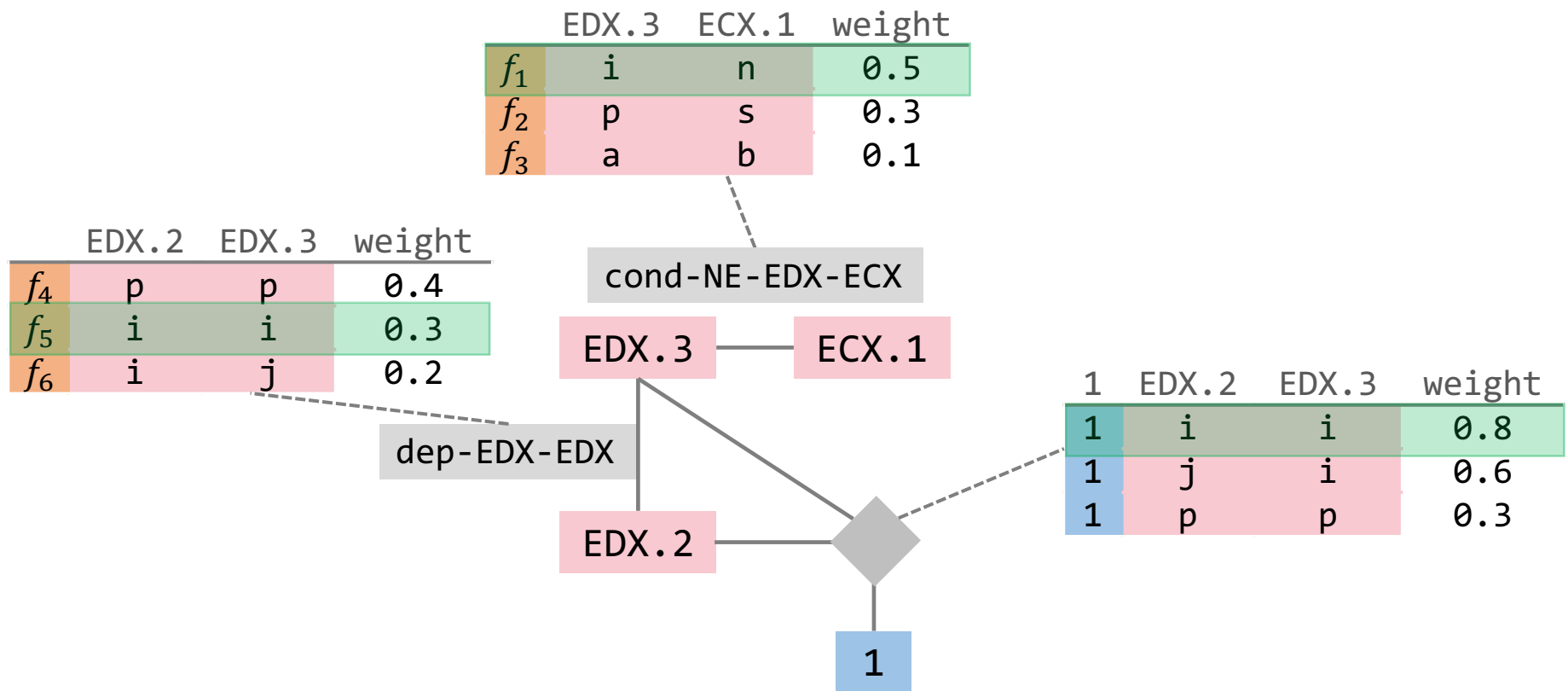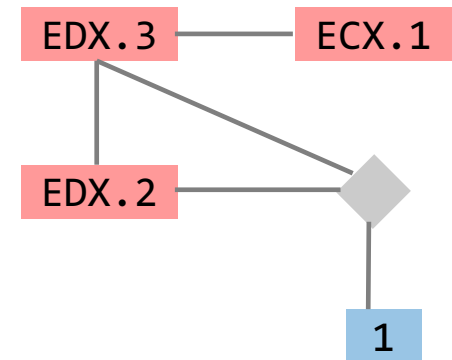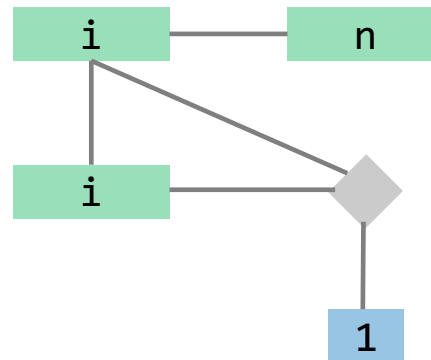Semantic variables

EDX.2    EDX.3    ECX.1

Temporary

EDX.1

Known elements

0    1    mov

EDX.3 —— ECX.1

EDX.2

1

# Recovering debug information

## MAP inference

|     | EDX.3 | ECX.1 | weight |
|-----|-------|-------|--------|
| $f_1$ | i     | n     | 0.5    |
| $f_2$ | p     | s     | 0.3    |
| $f_3$ | a     | b     | 0.1    |

cond-NE-EDX-ECX

|     | EDX.2 | EDX.3 | weight |
|-----|-------|-------|--------|
| $f_4$ | p     | p     | 0.4    |
| $f_5$ | i     | i     | 0.3    |
| $f_6$ | i     | j     | 0.2    |

dep-EDX-EDX

EDX.3 ——— ECX.1

EDX.2

| 1 | EDX.2 | EDX.3 | weight |
|---|-------|-------|--------|
| 1 | i     | i     | 0.8    |
| 1 | j     | i     | 0.6    |
| 1 | p     | p     | 0.3    |

1

# Recovering debug information



Stripped binary

Registers / mem offsets

Known elements

Semantic variables

Temporary

Known elements

Debug information

| Loc | Name | Type |
|-----|------|------|
|     | sum  | int  |
|     | n    | uint |
|     | i    | uint |
|     | res  | int  |

# DeBIN implementation

# DeBIN implementation

Static analysis: BAP

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

## Learning and inference

 http://scikit-learn.org
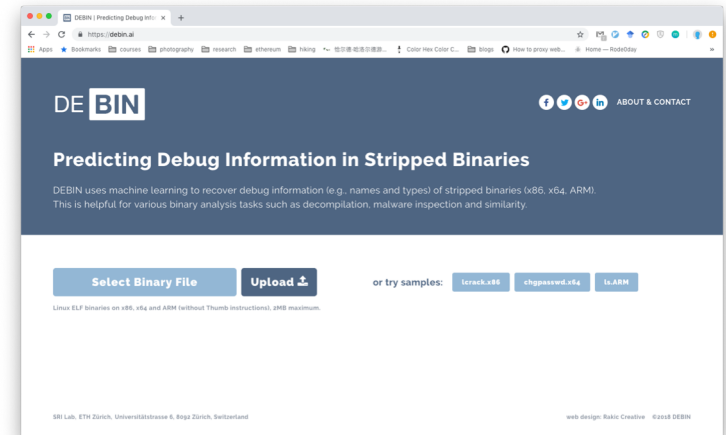
 http://nice2predict.org

 830 Linux packages
x86, x64, ARM
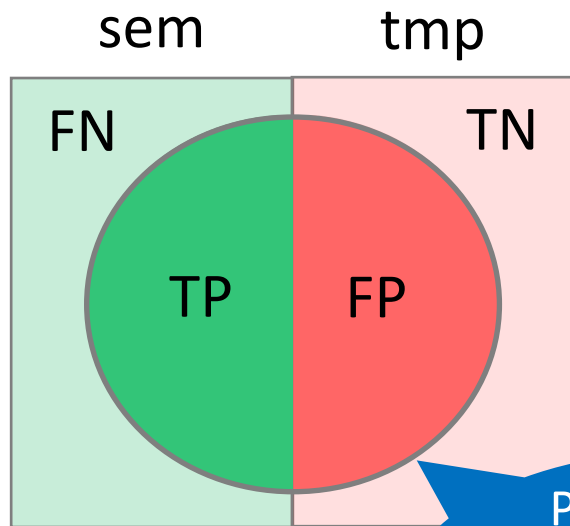


https://debin.ai

# DeBIN evaluation

1. How accurate is DeBIN's variable recovery?
2. How accurate is DeBIN's name and type prediction?
3. Is DeBIN useful for malware inspection?

# Variable recovery accuracy



$$\text{Accuracy} = \frac{|TP|+|TN|}{|sem|+|tmp|}$$

Predicted as semantic registers and memory offsets

Results

| Arch | Accuracy |
|------|----------|
| x86 | 87.1% |
| x64 | 88.9% |
| ARM | 90.6% |

DeBIN recovers variables with nearly 90% accuracy

# Name and type prediction accuracy

P          N



Predicted names and types

Correctly predicted names and types

$$\text{Precision} = \frac{|CP|}{|PN|} = \frac{|\ \ |}{|\ \ |}$$

$$\text{Recall} = \frac{|CP|}{|P|} = \frac{|\ \ |}{|\ \ |}$$

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

Total names and types (P) =

Predicted names and types (PN) =

Correct Predictions (CP) =

# 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 |

## Consistent precision/recall of roughly 65%

# Malware inspection

We inspected 35 x86 malware samples 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
    }
}
```
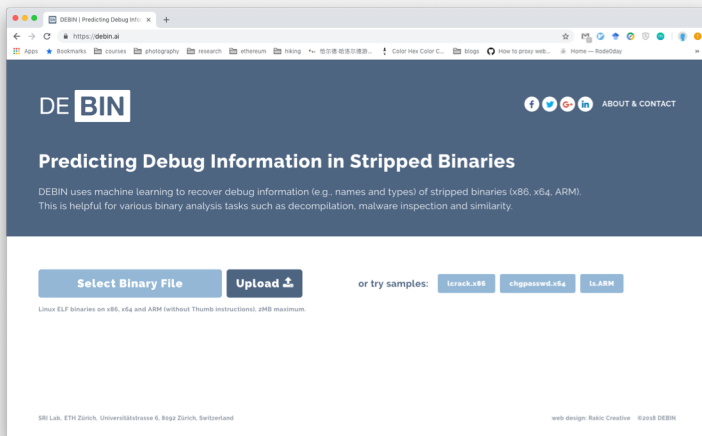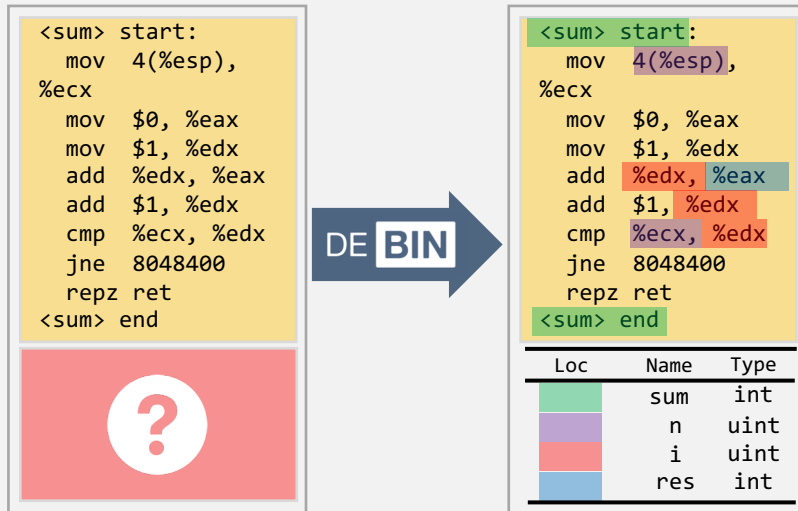
Leakage of sensitive data

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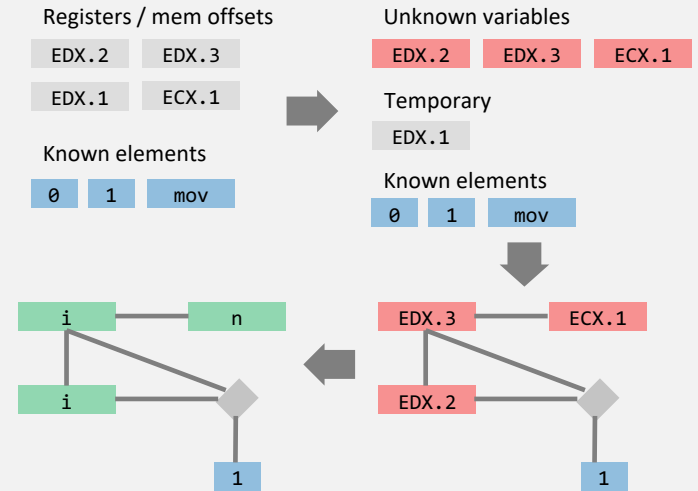
**DE BIN**

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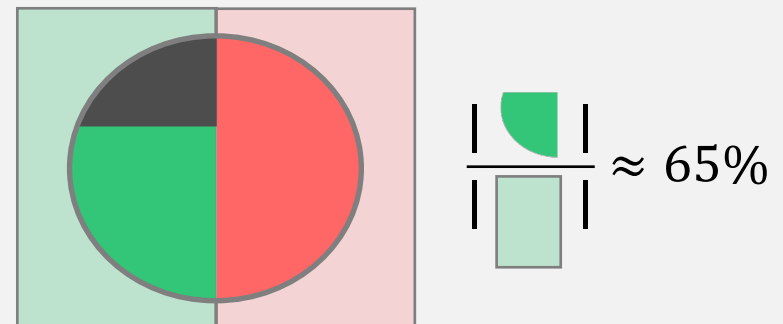
# Summary



```
<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
```

DE BIN

```
<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
```

| Loc | Name | Type |
|-----|------|------|
|     | sum  | int  |
|     | n    | uint |
|     | i    | uint |
|     | res  | int  |

**Predicting Debug Information in Stripped Binaries**

DEBIN uses machine learning to recover debug information (e.g., names and types) of stripped binaries (x86, x64, ARM). This is helpful for various binary analysis tasks such as decompilation, malware inspection and similarity.

Select Binary File   Upload ⬆

or try samples:  lcrack.x86   chgpasswd.x64   ls.ARM

Linux ELF binaries on x86, x64 and ARM (without Thumb instructions), 2MB maximum.

ABOUT & CONTACT

SRI Lab, ETH Zürich, Universitätstrasse 6, 8092 Zürich, Switzerland       web design: Rakic Creative   ©2018 DEBIN

Try online: https://debin.ai

Registers / mem offsets
EDX.2   EDX.3
EDX.1   ECX.1

Known elements
0   1   mov

Unknown variables
EDX.2   EDX.3   ECX.1

Temporary
EDX.1

Known elements
0   1   mov

EDX.3 — ECX.1

EDX.2

i — n

i

Two-stage prediction process

$\approx 65\%$

High precision and accuracy