NeurIPS 2019





Coda: An End-to-End Neural Program Decompiler

Cheng Fu¹, Huili Chen¹, Haolan Liu¹, Xinyun Chen³, Yuandong Tian², Farinaz Koushanfar¹, Jishen Zhao¹,

¹UC San Diego, ²Facebook AI Research, ³UC Berkeley



.

Background: Decompilation

Goal of Decompilation:

Binary executables to High-level programming language

Decompilation for SW defense:

Malware analysis

Vulnerability detection and fixing

Binary comparison/verification

Decompilation for SW attacking:

 Reverse engineering software binary with copyright protection for illegal usage.



Challenges

Prior decompilers (e.g. Hex-rey [1], RetDec [2]...) focus on reverse

RetDec

engineering the *functionality* of binary executables:

Semantics not guaranteed

Source Code

int	a	=	atoi(argv[1]);;
int	b	=	atoi(argv[2]);;
int	с	=	<pre>atoi(argv[3]);;</pre>
a =	b	*	c - 1 ;
if (a	>	1){
	a	=	b + c;
	с	=	a * c - b;

[1] HexRey, 2018, <u>https://www.hex-rays.com/products/decompiler/</u> [2] RetDec, 2019, <u>https://retdec.com/</u>

Decompiled Code

```
int32_t v1 = (int32_t)argv;
atoi((char*)*(int32_t*)(v1 + 4));
int32_t v2 = *(int32_t*)(v1 + 8);
int32_t v3 = *(int32_t*)(v1 + 12);
int32_t result;
if (v3 * v2 >= 3) {
  result = (v3 + v2) * v3 - v2;
} else {
  result = v3 * v2 < 3;
}
return result;
```

Challenges

Many hardware architectures (ISA): x86, MIPS, ARM

Many Programming Languages (PL)

 Extra Human effort to extend to the new version of the hardware architectures or programming languages

- Many formats of binary files
 - o .efl, .bin, .exe ...

Decompilation as Machine translation

Intuitively, decompilation is translation problem and can be solved using auto-encoder for machine translation:



 However, a naïve <u>sequence-to-sequence</u> model is hard to capture the meaning of low-level code and learn the grammar of high-level PL.

Coda Design *Leverage both syntax and dynamic information*



Stage 1: Code Sketch Generation

- What should the encoder captures?
 - Inter and intra instruction dependencies
- Instruction-aware encoder :
 - Coda leverages N-ary Tree Encoder [3] to capture
 - inter and intra dependencies of the low-level code.
 - Opcode and its operands are encoded together.
 - \circ Preserve the order of operands
 - Different encoders are used for encoding different instruction types, namely, memory (mem), branch (br) and arithmetic (art).

# : a i	sourc = b f(a	e C * c > c	program ;){		
,	c =	a * (c - b;		
ل	1	61	24/65-	mem	h0
1	1	ې۱,	24 (\$FP)	mem	h1
	I W	\$2,	20 (\$fp)	art	h2
2	mul	ŞT,	\$1, \$2	Dmem	h3
3	SW	\$ 1,	28(\$fp)	mem	h4
4	1 w	Ş1,	28(Şfp)	mem	h5
5	1 w 1	\$2,	20(\$fp)	art	h6
6	slt	\$1,	\$2, \$1		h7
7	beqz	\$1,	\$BB0_3	br	n /
8	j \$E	32		br	nð
9	\$B2:			br	h9
10	1 w	\$1,	28(\$fp)	mem	h1
11	1 w	\$2,	20(\$fp)	mem	h1
12	mu l	\$1.	\$1, \$2	art	h1
13	1 w	\$2.	24(\$fp)	mem	h1
14	subu	\$1.	\$1, \$2	art	h1
15	i	\$B3		br	h1
16	SW	\$1	20(\$fp)	mem	h1

Code

Sketch

^[3] Tai, Kai Sheng, Richard Socher, and Christopher D. Manning.

[&]quot;Improved semantic representations from tree-structured long short-term memory networks."

Stage 1: Code Sketch Generation

- Tree decoder for Abstract Syntax Tree (AST) generation:
 - AST can be equivalently translated into its corresponding high level Program
 - Advantages:
 - Prevent error propagation/ Preserve node dependency / easy to capture PL grammar
 - Boundaries are more explicit (terminal nodes)
- Parent and input attention feeding mechanism



 $s_k^{20} = \frac{\exp\{h_k^T \cdot h_{20}^{'}\}}{\sum_{j=0}^{16} \exp\{h_j^T \cdot h_{20}^{'}\}}$ Attention Probability $c_{20} = \sum_{k=0}^{16} h_k \cdot s_k^{20}$ Attention expectation $e_{20} = \tanh(W_1 c_t + W_2 h_{20}^{'})$ Attention Vec $t_{20} = \arg\max softmax(We_{20})$ Prediction

Code

Sketch

Stage 2: Iterative Error Correction

- The sketch generated in <u>Stage 1</u> may contain errors:
 - Mispredicted tokens, missing lines, redundant lines

Golden programMispredictedMissing linesRedundant lines $If(a > c) {$ $If(a > b) {$ $If(a > c) {$ $If(a > c) {$ $If(a > c) {$ a = b + c * a;a = b + c * a;a = b + c * a;a = b + c * a;b = a - c;b = a - b; $}$ b = a; $\}$ $\}$ $\}$ b = a;

Dynamic information that can be leveraged:

- $\circ~$ I/O pair to identify the correctness of the functionality
- Recompile the program back into low-level code (ϕ')---> should match with the golden low-level input (ϕ).

Error

Correction

Stage 2: Iterative Error Correction

 $(\phi, \phi') \longrightarrow$ Error Predictor

Error Predictor \longrightarrow (*Err Type*, *Err location*)

- Correct the error using an <u>Error Correction machine</u> (EC machine) guided by <u>the Error Predictor (EP)</u>.
- Optimization techniques:
 - Prevent the false alarm by recompile the updated sketch code and check its <u>Levenshtein edit loss</u> from the golden input.

FD

Ensemble multiple error predictors to cover more potential errors for updates.
 GACG GACG GACG GACG GACG GACG GACG

Error

Correctio

Stage 2: Iterative Error Correction

Algorithm 1 Workflow of iterative EC Machine.

INPUT: N_{EP} Ensembled Error Predictors EP; Source assembly ϕ ; Decompiled Sketch program P'; Compiler Γ ; Maximum iterations S_{max} and steps in each iteration c_{max} ;

OUTPUT: Error corrected program P'_f .

```
1: s_i \leftarrow 0
 2: while s_i < S_{max} do
           Q \leftarrow [], \phi' = \Gamma(P'), \ \Delta' \leftarrow Edit\_loss(\phi, \Gamma(P'))
 3:
           if \Delta' = 0 then break
 4:
           Q \leftarrow EP_i(P') for i = 1, ..., N_{EP} // Attach all the detected error to queue Q
 5:
           Q \leftarrow Prob\_sort(Q, c_{max}) // Rank Q using output probabilities, keep c_{max} results.
 6:
           while \widetilde{Q} is not empty do
 7:
                 err, node \leftarrow \tilde{Q}.pop()
 8:
                 P'_t \leftarrow FSM\_Error\_Correct(P', err, node)
                                                                                   // correct the error in the program
 9:
                 \Delta = \Delta' - Edit \ loss(\phi, \Gamma(P'_t))
10:
                 if \Delta > 0 then
11:
                       P' \leftarrow P'_t
12:
13: Return: P'_f \leftarrow P'
```

Iterative updating workflow

Error

Correction

Experimental Setup

- Compiler configuration : *clang –O0* (disabled optimization)
- Benchmarks:
 - Synthetic programs:
 - **Karel library (Karel)** only function calls (control graph)
 - Math library (Math) function calls with arguments (Data dependency + control graph)
 - Normal expressions (NE) (^,&,*,-,<<,>>,|,%) (Data dependency + control graph)
 - Math library + Normal expressions (Math + NE) replaces the variables in NE with a return value of math function. (Data dependency + control graph)

• Metrics:

- Token Accuracy : The percentage of predicted tokens that match with the groundtruth ones.
- **Program Accuracy: The percentage of programs that yields 100% Token accuracy.**

Results – Stage 1 Performance

• Token accuracy across benchmarks

Benchmarks	Seq2Seq	Seq2Seq+Attn	Seq2AST+Attn	Inst2seq+Attn	Inst2AST+Attn] _	
Karel _S	51.61	97.13	99.81	98.83	99.89	Baseline	
Math _S	23.12	94.85	99.12	96.20	99.72]	
NE _S	18.72	87.36	90.45	88.48	94.66		
$(Math+NE)_S$	14.14	87.86	91.98	89.67	97.90		
Karel _L	33.54	94.42	98.02	98.12	98.56	1	
Math _L	11.32	91.94	96.63	93.16	98.63	1	
NEL	11.02	81.80	85.92	85.97	91.92	1	
$(Math+NE)_L$	6.09	81.56	85.32	86.16	93.20]	

 X_S short programs, X_L long programs

 Coda yields the highest token accuracy across all benchmarks (96.8% on average) compared to all the other methods.

 Coda engenders 10.1% and 80.9% margin over a naive Seq2Seq model with and without attention.

More tolerant to the growth of program length.

Examples

#Golden: #
TurnOn();
TurnOff();
While(leftIsClear){
PutBeeper();
TurnLeft();
if(notFacingNorth){
if(notFacingNorth){
PickBeeper();
PickBe

#Decompiled TurnOn(); TurnOff(); while(leftIsClear){ PutBeeper(); TurnLeft(); if(notFacingNorth){ PickBeeper(); PickBeeper(); PickBeeper(); continue;} } PickBeeper();

TARGET CODE
48 46
d=isgreaterequal(a,d);
g=sin(f);
while(islessgreater(c,f)){
b=sin(b);
e=atan2(a,e);
e=fmax(e,c);
}
b=acos(f);
e=ceil(a);

Math.h

TRANSLATE CODE
d=isgreaterequal(a,d);
g=sin(f);
if(islessgreater(c,f)){
 b=sin(b);
 e=atan2(a.e);
 e=fmax(e,e);
 }
 b=acos(f);
 e=ceil(a);

Karel.h

err

Examples

#aaldaa	# d = 1 = _ 1	#aaldap	
#golden	#decompile	#golden	#decompiled
if(d <g*f){< td=""><td>if(d<g*f){< td=""><td>b=fmax(l,c)*asin(j);</td><td>b=fmax(1 c)*asin(i)</td></g*f){<></td></g*f){<>	if(d <g*f){< td=""><td>b=fmax(l,c)*asin(j);</td><td>b=fmax(1 c)*asin(i)</td></g*f){<>	b=fmax(l,c)*asin(j);	b=fmax(1 c)*asin(i)
a=a*a-d;	a=a*a-d:	f=fmax(k,h);	f = fmax(k, b)
b=a/a+f:	b=a/a+f	if(isless(k.c) islessequal(f.q)){	I = I P I d X (K, II);
o-o/f:		h-tan(f)	lf(lsless(k,c) lslessequal(f,g)){
c-a/i,	e=a-1;		b= <u>tan(f);</u>
a=a-a-g;	a=a-a;	a=isless(b,a)*islessgreater(c,r);	a=atan(b,a)*islessgreater(c,f);
}	}	a=fdim(i,l)/islessequal(l,k);	a=fdim(i,l):
else{	else{	}	1
d=c*f-a:	d=c*f-a:	else{	
d=a/f·	d = a/f	i = exp(i) + pow(i h)	else
		$\sum_{n=1}^{\infty} \left(\frac{1}{n} \right) = \sum_{n=1}^{\infty} \left(\frac{1}{n} \right)$	j=exp(j)+isgreaterequal(j,g);
a=c+r+c;	a=c+r+c;	a=pow(d,a)/stn(g);	a=pow(d,a)/sin(g);
}	f=g;	h=islessequal(i,b)/tan(j);	h=islessequal(i,b)/tan(i):
a=d*f;	}	}	3
a=a+e;	a=d*f;		

Math.h + Normal Expressions

Normal Expressions

a=a+e; a=a+e:

Results – Stage 2 Performance

- Part (i): By ensemble 10 EP, Coda achieves 90.8% error detection rate.
- Part (ii): Coda's EC machine increases the program accuracy from <u>30% to</u> <u>82%</u> on average for Inst2AST-based code sketch generation, respectively.

BenchMarks	(i) Error Detection		(ii) Befor EC		After EC	
Denemviarks	s2s,10	i2a,10	s2s	i2a	s2s	$\imath 2a$
Math _S	91.4	94.2	40.1	64.8	91.2	100.0
NE _S	83.5	88.7	6.6	12.2	53.0	78.6
$(Math+NE)_S$	83.6	90.1	3.5	43.2	63.6	89.2
$Math_L$	87.5	91.3	21.7	51.8	83.9	99.5
NE _L	78.1	84.5	0.2	2.6	33.1	56.4
$(Math+NE)_L$	80.2	85.3	0.1	4.9	38.3	67.2



s2s = sequence-to-sequence with attention

I2a = instruction encoder to AST decoder with attention

Results -- Overall

Coda vs. traditional decompiler (RetDec)

- Lines of code: ~10K vs. ~500K -- 50x reduction
- Toolkit size: ~10MB Neural network size vs. ~5GB toolkit size -- 500x reduction
- Program accuracy: 82% vs. no semantics guarantee

Discussion

Extremely long programs

- LSMT is not good at remembering long sequence.
- Unlike nature language, low-level PL does not have explicit breakup position.

Sensitive to ISA

- Complicated data type / structure / class
- Compiler Optimizations / Obfuscation ...

Summary of Coda

- The first neural-based decompilation framework, which preserves both the semantics and the <u>functionality</u> of the high-level program
- Decomposes the decompilation task into of two key phases -- <u>code sketch</u> <u>generation</u> and <u>iterative errors correction</u>
- Significantly outperforms the Seq2Seq model and traditional decompilers

NeurIPS 2019





Coda: An End-to-End Neural Program Decompiler

Cheng Fu¹, Huili Chen¹, Haolan Liu¹, Xinyun Chen³, Yuandong Tian², Farinaz Koushanfar¹, Jishen Zhao¹,

¹UC San Diego, ²Facebook AI Research, ³UC Berkeley