## Using Logic Programming to Recover C++ Classes and Methods from Compiled Executables

Edward J. Schwartz, Cory F. Cohen, Michael Duggan, Jeffrey Gennari, Jeffrey S. Havrilla, Charles Hines



Software Engineering Institute Carnegie Mellon University Pittsburgh, PA 15213

Carnegie Mellon University Software Engineering Institute

Copyright 2018 Carnegie Mellon University. All Rights Reserved.

This material is based upon work funded and supported by the Department of Defense under Contract No. FA8702-15-D-0002 with Carnegie Mellon University for the operation of the Software Engineering Institute, a federally funded research and development center.

The view, opinions, and/or findings contained in this material are those of the author(s) and should not be construed as an official Government position, policy, or decision, unless designated by other documentation.

NO WARRANTY. THIS CARNEGIE MELLON UNIVERSITY AND SOFTWARE ENGINEERING INSTITUTE MATERIAL IS FURNISHED ON AN "AS-IS" BASIS. CARNEGIE MELLON UNIVERSITY MAKES NO WARRANTIES OF ANY KIND, EITHER EXPRESSED OR IMPLIED, AS TO ANY MATTER INCLUDING, BUT NOT LIMITED TO, WARRANTY OF FITNESS FOR PURPOSE OR MERCHANTABILITY, EXCLUSIVITY, OR RESULTS OBTAINED FROM USE OF THE MATERIAL. CARNEGIE MELLON UNIVERSITY DOES NOT MAKE ANY WARRANTY OF ANY KIND WITH RESPECT TO FREEDOM FROM PATENT, TRADEMARK, OR COPYRIGHT INFRINGEMENT.

[DISTRIBUTION STATEMENT A] This material has been approved for public release and unlimited distribution. Please see Copyright notice for non-US Government use and distribution.

This material may be reproduced in its entirety, without modification, and freely distributed in written or electronic form without requesting formal permission. Permission is required for any other use. Requests for permission should be directed to the Software Engineering Institute at permission@sei.cmu.edu.

DM18-0828



**Carnegie Mellon University** Software Engineering Institute

Using Logic Programming to Recover C++ Classes and Methods from Compiled Executables © 2018 Carnegie Mellon University



**Carnegie Mellon University** Software Engineering Institute Using Logic Programming to Recover C++ Classes and Methods from Compiled Executables © 2018 Carnegie Mellon University



**Carnegie Mellon University** Software Engineering Institute

Using Logic Programming to Recover C++ Classes and Methods from Compiled Executables © 2018 Carnegie Mellon University

#### Goals: How?

#### Static

- Analyze program without executing it
- Can be applied to unknown software
- No need for test cases

# •Recover <u>all</u> classes and <u>all</u> methods

# Do not rely on RTTI metadata

- Only available for polymorphic classes
- Often removed or corrupted in malware

# **Existing Work**

System	Static	All classes	All methods	Works w/o RTTI
SmartDec Fokin 2010	$\checkmark$	×	$\checkmark$	$\checkmark$
SecondWrite Yoo 2014	$\checkmark$	×	×	×
Various CFI systems Pawlowski 2017, etc.	$\checkmark$	×	×	$\checkmark$

# These systems primarily recover information from virtual function tables

Using Logic Programming to Recover C++ Classes and Methods from Compiled Executables © 2018 Carnegie Mellon University

#### **Virtual Function Tables**

Shape's VFTable

0: Shape::draw

1: Shape::getArea

Square's VFTable 0: Square::draw

1: Square::getArea

VFTables enumerate the virtual functions of each class... but only for polymorphic classes.

**Carnegie Mellon University** Software Engineering Institute Using Logic Programming to Recover C++ Classes and Methods from Compiled Executables © 2018 Carnegie Mellon University

# **Existing Work**

System	Static	All classes	All methods	Works w/o RTTI
SmartDec Fokin 2010	$\checkmark$	×	$\checkmark$	$\checkmark$
SecondWrite Yoo 2014	$\checkmark$	×	×	×
Various CFI systems Pawlowski 2017, etc.	$\checkmark$	×	×	$\checkmark$
<b>Lego</b> Srinivasan 2014	×	$\checkmark$	$\checkmark$	$\checkmark$
<b>ObjDigger</b> Jin 2014	$\checkmark$	$\checkmark$	$\checkmark$	$\checkmark$

#### ObjDigger

ObjDigger uses <u>object pointer tracking</u> to recover non-polymorphic classes too



#### If the Obj class does not inherit, it defines foo and bar

**Carnegie Mellon University** Software Engineering Institute Using Logic Programming to Recover C++ Classes and Methods from Compiled Executables © 2018 Carnegie Mellon University

#### **Ambiguous Scenarios**



**Carnegie Mellon University** Software Engineering Institute

Using Logic Programming to Recover C++ Classes and Methods from Compiled Executables © 2018 Carnegie Mellon University

## **ObjDigger's Shortcomings**

•ObjDigger makes mistakes on ambiguous scenarios

- We tuned it so it was "usually" correct
- But mistakes propagated and caused more serious problems
- •ObjDigger was written procedurally in C++
  - if (ambiguous\_scenario) { classes[foo].add\_method(bar); } // guess!
  - There was no mechanism to detect mistakes and repair them

# Improving accuracy became more and more difficult

#### **OOAnalyzer**

OOAnalyzer was designed to overcome ObjDigger's shortcomings

- Reasoning implemented in Prolog
- •Hypothetical reasoning
  - Detect ambiguous scenarios, and make best guess first
  - <u>New</u>: detect when guesses introduce inconsistency and undo them
  - Implemented using Prolog backtracking
- •Reasoning is defined using <u>declarative</u> rules
  - Much easier to understand and maintain

**Carnegie Mellon University** Software Engineering Institute

#### **OOAnalyzer System Design**

**Carnegie Mellon University** Software Engineering Institute Using Logic Programming to Recover C++ Classes and Methods from Compiled Executables © 2018 Carnegie Mellon University



Apolying class (1s 1998a33c to 1998b34

#### Fact Exporter

- Uses conventional binary analysis to produce *initial* facts about the program
  - Initial facts describe low-level program behaviors
- Simple symbolic analysis
  - Intentionally favors scalability over accuracy
  - Does not use constraint solvers
  - Symbolic memory aliases if memory addresses are equal after simplification
  - Path sensitive up to a threshold
- Sufficient because reasoning system can cope with mistakes

#### **Initial Facts**

Initial facts describe low-level program behaviors and form the basis upon which OOAnalyzer's reasoning system operates

Fact Name	Description
ObjPtrAllocation(I, F, P, S)	Instruction I in function F allocates S bytes of memory for the object pointed to by P.
ObjPtrInvoke(I, F, P, M)	Instruction I in function F calls method M on the object pointed to by P.
ObjPtrOffset(P <sub>1</sub> , O, P <sub>2</sub> )	Object pointer $P_2$ points to $P_1 + O$ .
MemberAccess(I, M, O, S)	Instruction I in method M accesses S bytes of memory at offset O from the current object's pointer.
ThisCallMethod(M, P)	Method M receives the object pointed to by P in the ecx register.
NoCallsBefore(M)	No methods are called on any object pointer before method M.
ReturnsSelf(M)	Method M returns the object pointer that was passed as a parameter.
UninitializedReads(M)	Method M reads memory that was not written to by M.
PossibleVFTableEntry(VFT, O, M)	Method M may be at offset O in vftable VFT.
PossibleVFTableEntry(VFT, O, M)	Method M may be at offset O in vftable VFT.

**Carnegie Mellon University** Software Engineering Institute Using Logic Programming to Recover C++ Classes and Methods from Compiled Executables © 2018 Carnegie Mellon University

#### **Entity Facts**

Entity facts are produced during the reasoning process and describe the high-level model of the program being analyzed

Fact Name	Description
Method(M)	Method M is an OO method on a class or struct.
Constructor(M)	Method M is an object constructor.
Destructor(M)	Method M is an object destructor.
$CI_a = CI_b$	The sets of methods Cl <sub>a</sub> and Cl <sub>b</sub> both represent methods from the same class. These sets should be combined into a single class.
$Cl_a \le Cl_b$	The sets of methods Cl <sub>a</sub> and Cl <sub>b</sub> either both represent methods from the same class, or, the methods in Cl <sub>b</sub> are inherited from Cl <sub>a</sub> .
$M \in CI$	Method M is defined directly on class Cl.
ClassCallsMethod(Cl, M)	An instance of class Cl calls method M.
Other estagarias include virtual functions	aleas relationships, and sizes of aleases and tables

Other categories include virtual functions, class relationships, and sizes of classes and tables.

#### Reasoning



#### Forward reasoning

- Unambiguous scenarios
- Interpretation: If *P*<sub>1</sub>, *P*<sub>2</sub>, ... and *P*<sub>n</sub> are satisfied, then *C* is true
- If inconsistency is reached,  $P_1, P_2, \ldots$  or  $P_n$  must not be true

- Hypothetical reasoning
  - Ambiguous scenarios
  - Interpretation: If *P*<sub>1</sub>, *P*<sub>2</sub>, ... and *P*<sub>n</sub> are satisfied, then guess *C* is true
  - If inconsistency is reached, then retract *C* and assume  $\neg C$
  - If inconsistency is still reached,  $P_1, P_2, \ldots$  or  $P_n$  must not be true

#### **Forward Reasoning**

If a method is directly called on a base class, it cannot be directly defined on the derived class.

$$\begin{array}{lll} & \text{Constructor}(\mathsf{M}_d) & \mathsf{M}_d \in \mathsf{Cl}_d \\ & \text{Constructor}(\mathsf{M}_b) & \mathsf{M}_b \in \mathsf{Cl}_b \\ & \text{ClassCallsMethod}(\mathsf{Cl}_d, \mathsf{M}) \\ & \text{ClassCallsMethod}(\mathsf{Cl}_b, \mathsf{M}) & \mathsf{M}_d \neq \mathsf{M}_b \\ & \mathsf{M} \in \mathsf{Cl}_m & \mathsf{Cl}_d \neq \mathsf{Cl}_b & \text{DerivedClass}(\mathsf{Cl}_d, \mathsf{Cl}_b, \_) \end{array}$$

$$CI_m \neq CI_d$$

**Carnegie Mellon University** Software Engineering Institute Using Logic Programming to Recover C++ Classes and Methods from Compiled Executables © 2018 Carnegie Mellon University

#### Hypothetical Reasoning

If a method is called on a derived class but not a base class, (first) assume the method is defined on the derived class.

# $\begin{aligned} \mathsf{ClassCallsMethod}(\mathsf{Cl}_d,\mathsf{M}) & \neg\mathsf{ClassCallsMethod}(\mathsf{Cl}_b,\mathsf{M}) \\ \mathsf{M}\in\mathsf{Cl} & \mathsf{DerivedClass}(\mathsf{Cl}_d,\mathsf{Cl}_b,\_) \end{aligned}$

$$CI_d = CI$$

**Carnegie Mellon University** Software Engineering Institute Using Logic Programming to Recover C++ Classes and Methods from Compiled Executables © 2018 Carnegie Mellon University

#### Hypothetical Reasoning



Using Logic Programming to Recover C++ Classes and Methods from Compiled Executables © 2018 Carnegie Mellon University

#### Hypothetical Reasoning



**Carnegie Mellon University** Software Engineering Institute Using Logic Programming to Recover C++ Classes and Methods from Compiled Executables © 2018 Carnegie Mellon University

#### **Research Questions & Evaluation**

**Carnegie Mellon University** Software Engineering Institute Using Logic Programming to Recover C++ Classes and Methods from Compiled Executables © 2018 Carnegie Mellon University

#### **Research Questions**

- How well does OOAnalyzer identify C++ classes and their constituent methods?
  - OOAnalyzer classes may not naturally correspond to ground truth classes
  - New metric: edit distance

#### Edit Distance Example



#### **Carnegie Mellon University** Software Engineering Institute

Using Logic Programming to Recover C++ Classes and Methods from Compiled Executables © 2018 Carnegie Mellon University

#### ObjDigger vs OOAnalyzer Edit Distances on ObjDigger Programs



Using Logic Programming to Recover C++ Classes and Methods from Compiled Executables © 2018 Carnegie Mellon University

## ObjDigger vs OOAnalyzer Edit Distances on Cleanware

Program	# Class	# Method	ObjDigger Edits	ObjDigger Edits (%)	OOAnalyzer Edits	OOAnalyzer Edits (%)
Firefox.exe	141	638	507	79.5%	212	33.2%
Log4cpp Debug	139	893	829	92.8%	239	26.8%
Log4cpp Release	76	378	272	72.0%	75	19.8%
muParser Debug	180	1437	1361	94.7%	483	33.6%
muParser Release	94	598	369	61.7%	183	30.6%
MySQL cfg_editor.dll	190	1266	00	00	391	30.9%
MySQL mysql.exe	202	1395	$\infty$	$\infty$	439	31.5%
TinyXML Debug	35	415	268	64.6%	69	16.6%
TinyXML Release	33	283	174	61.5%	55	19.4%

**Carnegie Mellon University** Software Engineering Institute Using Logic Programming to Recover C++ Classes and Methods from Compiled Executables © 2018 Carnegie Mellon University

## ObjDigger vs OOAnalyzer Edit Distances on Cleanware

Program	# Class	# Method	ObjDigger Edits	ObjDigger Edits (%)	OOAnalyzer Edits	OOAnalyzer Edits (%)
Firefox.exe	141	638	507	79.5%	212	33.2%
Log4cpp Debug	139	893	829	92.8%	239	26.8%
Log4cpp Release	76	378	272	72.0%	75	19.8%
muParser Debug	180	1437	1361	94.7%	483	33.6%
muParser Release	94	598	369	61.7%	183	30.6%
MySQL cfg_editor.dll	190	1266	8	8	391	30.9%
MySQL mysql.exe	202	1395	$\infty$	$\infty$	439	31.5%
TinyXML Debug	35	415	268	64.6%	69	16.6%
TinyXML Release	33	283	174	61.5%	55	19.4%

OOAnalyzer recovers between 67 and 84% of methods on medium to large programs

## ObjDigger vs OOAnalyzer Edit Distances on Malware

Program	# Class	# Method	ObjDigger Edits	ObjDigger Edits (%)	OOAnalyzer Edits	OOAnalyzer Edits (%)
Malware Ofaaa3d3	21	135	121	89.6%	21	15.6%
Malware 29be5a33	19	130	91	70.0%	15	11.5%
Malware 6098cb7c	55	339	131	38.6%	29	8.6%
Malware 628053dc	207	1920	1245	64.8%	378	19.7%
Malware 67b9be3c	400	2072	1299	62.7%	670	32.3%
Malware cfa69fff	39	184	125	67.9%	37	20.1%
Malware d597bee8	19	133	68	51.1%	17	12.8%
Malware deb6a7a1	283	2712	1900	70.1%	639	23.6%
Malware f101c05e	169	1601	987	61.6%	329	20.5%

**Carnegie Mellon University** Software Engineering Institute

Using Logic Programming to Recover C++ Classes and Methods from Compiled Executables © 2018 Carnegie Mellon University

#### **Research Questions**

- How well does OOAnalyzer identify C++ classes and their constituent methods?
  - OOAnalyzer classes may not naturally correspond to ground truth classes
  - New metric: edit distance
- How well does OOAnalyzer classify methods as constructors, destructors, and virtual methods?

#### OOAnalyzer Method Classification on ObjDigger Programs

Pr	ogram		Constr	uctor		
		Recall	Prec.	F		
CI	mg	44/51	44/53	0.85		
Li	ght-pop3-smtp	41/52	41/44	0.85		
0	otionparser	10/11	10/10	0.95		
Pi	coHttpD	117/142	117/126	0.87		
X3	c	6/7	6/6	0.92		
	6 of the 7 were i	constru dentifie	ctors d		Of the 6 constructors identified, 6 were correct	$2 * \frac{recall * precision}{recall + precision}$

Using Logic Programming to Recover C++ Classes and Methods from Compiled Executables © 2018 Carnegie Mellon University

#### OOAnalyzer Method Classification on Cleanware

Program		Constr	uctor		Destr	uctor		VFT	ables	V	/irtual Met	hods
	Recall	Prec.	F	Recall	Prec.	F	Recall	Prec.	F	Recall	Prec.	F
Firefox.exe	40/51	40/54	0.76	1/39	1/1	0.05	18/33	18/18	0.71	85/101	85/98	0.85
Log4cpp Debug	192/209	192/197	0.95	40/118	40/40	0.51	18/18	18/18	1.00	84/101	84/86	0.92
Log4cpp Release	135/165	135/170	0.81	24/73	24/36	0.44	18/21	18/18	0.92	84/101	84/86	0.90
muParser Debug	293/325	293/314	0.92	28/156	28/30	0.30	12/12	12/13	0.96	35/47	35/43	0.78
muParser Release	197/252	197/269	0.76	15/91	15/21	0.27	12/14	12/13	0.89	35/47	35/37	0.83
MySQL cfg_editor.dll	260/290	260/311	0.87	107/281	107/111	0.55	69/69	69/69	1.00	321/427	321/325	0.85
MySQL mysql.exe	282/314	282/341	0.86	115/300	115/121	0.55	75/75	75/75	1.00	341/453	341/345	0.85
TinyXML Debug	53/60	53/57	0.91	0/39	0/3	0.00	24/24	24/24	1.00	101/119	101/102	0.91
TinyXML Release	49/60	49/53	0.87	27/39	27/36	0.72	24/24	24/24	1.00	101/119	101/103	0.91

#### Conclusion: OOAnalyzer

- Can <u>scalably</u> and <u>statically</u> recover abstractions about <u>all</u> classes and <u>all</u> methods from C++ executables
- Combines a light-weight symbolic analysis with a Prolog-based reasoning system that enables <u>hypothetical reasoning</u> about ambiguous scenarios
- Recovers <u>70% of classes/methods</u> or more on most cleanware and malware
- Classifies constructors, virtual methods, and virtual function tables with high accuracy

Open source tool in the Pharos binary analysis framework: <u>https://github.com/cmu-sei/pharos</u>

#### **Questions?**

#### Presenter

Edward J. Schwartz Research Scientist Carnegie Mellon University SEI CERT/CC

#### Contact

Email: <a href="mailto:eschwartz@cert.org">eschwartz@cert.org</a>

#### Open source tool in the Pharos binary analysis framework: <u>https://github.com/cmu-sei/pharos</u>

**Carnegie Mellon University** Software Engineering Institute Using Logic Programming to Recover C++ Classes and Methods from Compiled Executables © 2018 Carnegie Mellon University

Functions window 🛛 💿 🔞 [ IDA V	liew-A		- O C
unction name ∫ sub_10001000 ∫ sub_100011A0 ∫ sub_100011A0 ∫ objdigger Classes ♥ Cls_1000a238 ♥ Methods Cls_1000a238 ♥ Methods Cls_100015d0 Meth1_100015d0 Meth3_100017d0 Meth3_100017d0 Meth5_100018c0 Meth5_10001900 ▶ Members ▶ Cls_1000a27c ▶ Cls_1000a28c ▶ Cls_1000a28c ▶ Cls_1000a33c ▶ Cls_1000a33c ▶ Cls_1000a33c	<pre>.text:100014F0 ; Cls_1000a238: .text:100014F0 Ctor_100014f0 .text:100014F0 000 .text:100014F1 004 .text:100014F3 004 .text:100014F3 004 .text:100014F3 004 .text:100014F8 008 .text:10001506 004 .text:10001508 004 .text:10001508 004 .text:10001510 004 .text:10001510 004 .text:10001510 004 .text:10001510 004 .text:10001510 004 .text:10001510 004 .text:10001510 004 .text:10001510 004 .text:10001510 004 .text:10001520 004 .text:10001527 004 .text:10001528 000 .text:10001529 ;</pre>	<pre>::Ctor_100014f0 (constructor) proc near ; (     ; s push esi mov esi, ecx push offset LibFileName mov [esi+Cls_1000a238.x call ds:LoadLibraryA xor ecx, ecx mov [esi+Cls_1000a238.x mov [esi+Cls_</pre>	CODE XREF: sub_1 sub_10004A40+231 ; "cabinet.dll" rfptr_0], offset sem_4], eax sem_6], ecx sem_10], ecx sem_10], ecx sem_14], ecx sem_14], ecx sem_12], ecx sem_20], ecx sem_24], ecx
Output window			

finalClass(ClassID, VFTable, MinSize, MaxSize, RealDestructor, MethodList).

- finalClass(0x412190, 0, 0x5, 0x5, 0, [0x412190, 0x4122f0, 0x4123c0, 0x4124c0]).
- finalClass(0x412400, 0, 0x4, 0x4, 0, [0x412400]).
- finalClass(0x417824, 0x417824, 0x54, 0x54, 0, [0x411860, 0x411970, 0x4119f0]). finalClass(0x417860, 0x417860, 0xc, 0xc, 0, [0x4117e0, 0x411a70, 0x411ac0]). finalClass(0x41787c, 0x41787c, 0x4, 0x4, 0, [0x4118f0, 0x411b40, 0x411b90]). finalClass(0x41789c, 0x41789c, 0x4, 0x4, 0, [0x412800, 0x412830]). finalClass(0x41b2f8, 0, 0, 0, 0, [0x41b2f8, 0x41b2fc, 0x41b324, 0x41b330]). finalClass(0x41b304, 0, 0, 0, 0, [0x41b304, 0x41b310, 0x41b348]). finalClass(0x41b308, 0, 0, 0, 0, [0x41b308, 0x41b320, 0x41b33c, 0x41b340]). finalClass(0x41b30c, 0, 0, 0, 0, [0x41b30c, 0x41b31c, 0x41b334, 0x41b338]). finalClass(0x41b328, 0, 0, 0, 0, [0x41b328, 0x41b32c, 0x41b344]).

finalInheritance(DerivedClassID, BaseClassID, ObjectOffset, VFTable, Virtual).

- finalInheritance(0x412190, 0x412400, 0, 0, false).
- finalInheritance(0x417824, 0x417860, 0, 0x417824, false).
- finalInheritance(0x417824, 0x41787c, 0xc, 0x417814, false).

#### finalMember(Class, Offset, Sizes, Certainty).

- finalMember(0x412190, 0x4, [0x1], certain).
- finalMember(0x412400, 0, [0x4], certain).
- finalMember(0x417824, 0x50, [0x4], certain).
- finalMember(0x417860, 0, [0x4], certain).
- finalMember(0x41787c, 0, [0x4], certain).
- finalMember(0x41789c, 0, [0x4], certain).

#### finalMethodProperty(Method, Type, Certainty).

- finalMethodProperty(0x411860, constructor, certain).
- finalMethodProperty(0x4118f0, virtual, certain).
- finalMethodProperty(0x411970, virtual, certain).
- finalMethodProperty(0x4119f0, virtual, certain).
- finalMethodProperty(0x411a70, constructor, certain).
- finalMethodProperty(0x411ac0, virtual, certain).
- finalMethodProperty(0x411b40, constructor, certain).
- finalMethodProperty(0x411b90, virtual, certain).
- finalMethodProperty(0x412190, constructor, certain).
- finalMethodProperty(0x412400, constructor, certain).
- finalMethodProperty(0x412800, constructor, certain).
- finalMethodProperty(0x412830, virtual, certain).