

JARVIS: New Challenge Problems for Software Reverse Engineering

BSidesDC

October 26th, 2019

evm

@evm_sec

A Future Vision for Software Reverse Engineering

**Imagine if Software RE
could be faster, more
effective and accessible**



- You have a tool that produces meaningful, descriptive labels in your disassembly of a completely new binary
- You automatically get a well-labeled software architecture diagram with a description of each module when analyzing a new binary
- People with little or no training can be given an accurate description of a new binary produced by an automated system

Talk Outline

- What problems in software reverse engineering should we consider “solved” ?
- What are the gaps between our solved problems and our existing process?
- What problems should we work on next?
- How do we get started?

My Background – The Embedded Systems RE World

- One “device” – but many boards/processors/firmwares
- Bare metal and RTOS environments
 - Large fully-linked binaries
 - No distinction between OS/libraries/application code
 - No clearly labeled system calls for clues
 - Vulnerability analysis: Takes much longer to get to the entry point
- Dynamic analysis is rarely possible
 - Unreliable JTAG/debugger access
 - Debug fuses
 - Debug BGA balls not pulled out on the motherboard
 - Integrated cores with no physical connection (debugged in development on dev boards)

**What problems in software
reverse engineering should
we consider “solved” ?**

Solved Problems?

- Decompilation
- Function-to-function Matching
- *Combined Static & Dynamic Analysis Approaches**

**since dynamic analysis is rarely possible in my world I am not qualified to evaluate the state of the toolsets (which mostly exist in other environments)*

Decompilation

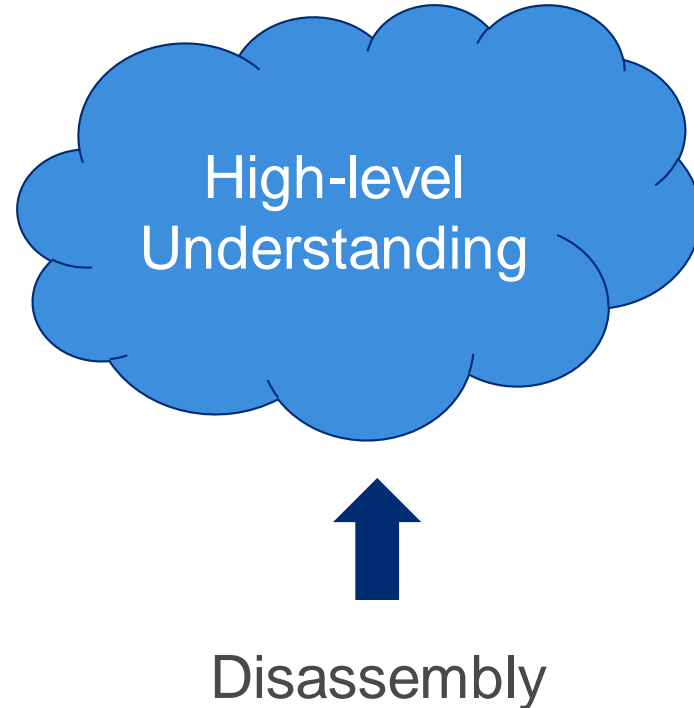
- IDA Pro,^[1] GHIDRA,^[2] RetDec,^[3] Rellic,^[4] angr-management,^[5] JEB^[6]
- Follow process first described by C. Cifuentes in 1994^[7]
 - Syntax Analysis
 - Semantic Analysis
 - Intermediate Code Generation
 - SLEIGH/p-code (GHIDRA) ^[8]
 - Microcode (IDA) ^[9]
 - Interlude (JEB) ^[10]
 - LLVM IR (RetDec and Rellic) ^{[3][4]}
 - VEX (angr-management) ^[5]
 - Control Flow Graph Generation
 - Data Flow Analysis
 - Control Flow Analysis
 - Code Generation
- Decompilers always produce blank code!
- Incremental improvements in decompilation will not lead to fundamental changes in the speed & effectiveness of RE

Function Matching

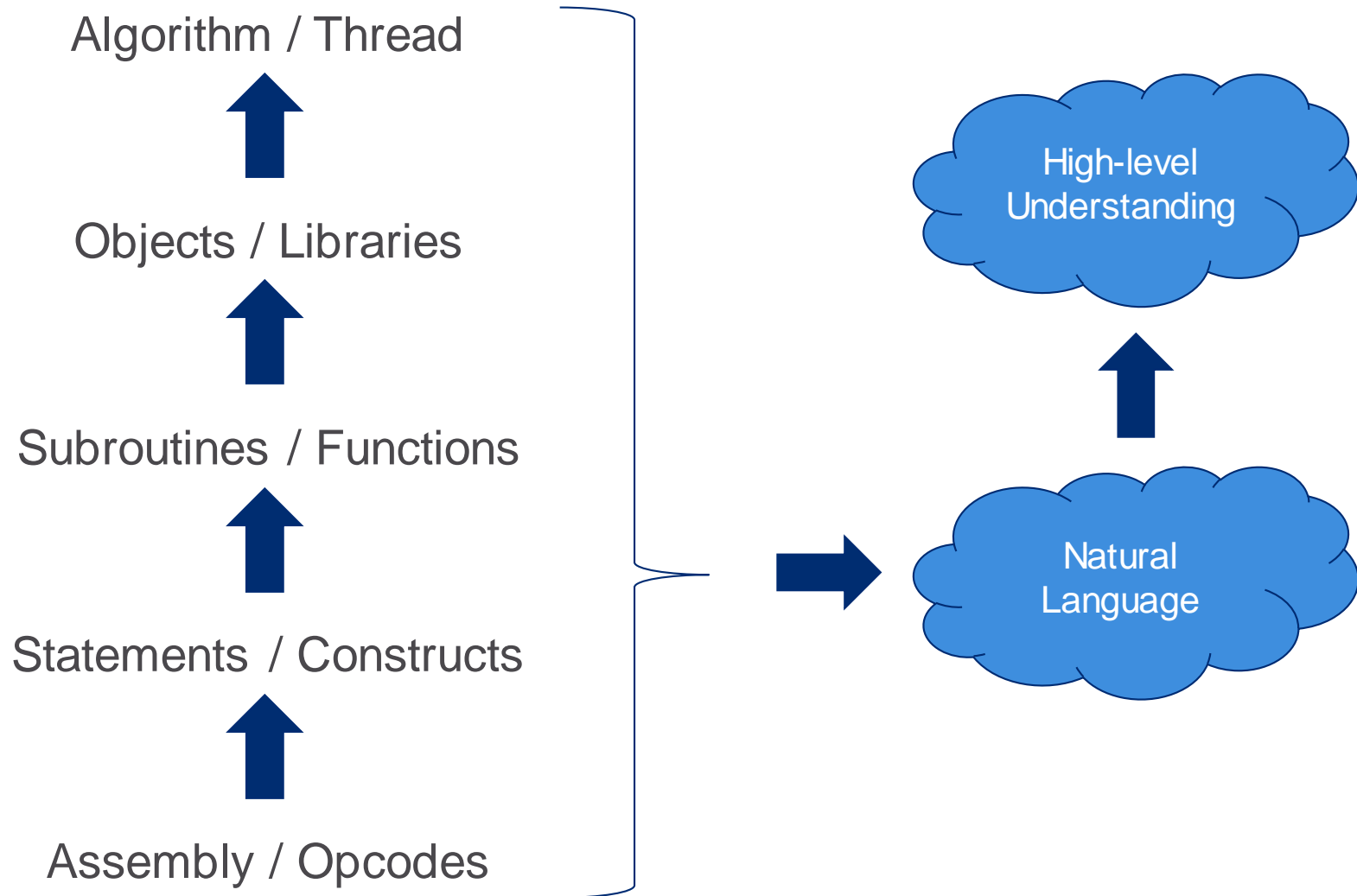
- BinDiff,^[11] Diaphora,^[12] GHIDRA Version Tracker^[13]
- Algorithm
 - Lossy compression function at the basic block level
 - Small primes product ^[14]
 - Graph-based comparison function for function flow graph
 - MD Index ^{[15][16]}
 - KOKA hash ^[17]
 - Graph-based comparison functions for call graph
- Works fairly well in most situations where library code is pre-identified
- Incremental improvements to function matching will likewise not lead to fundamental changes in the speed & effectiveness of RE

**What are the gaps between
our solved problems and
our existing process?**

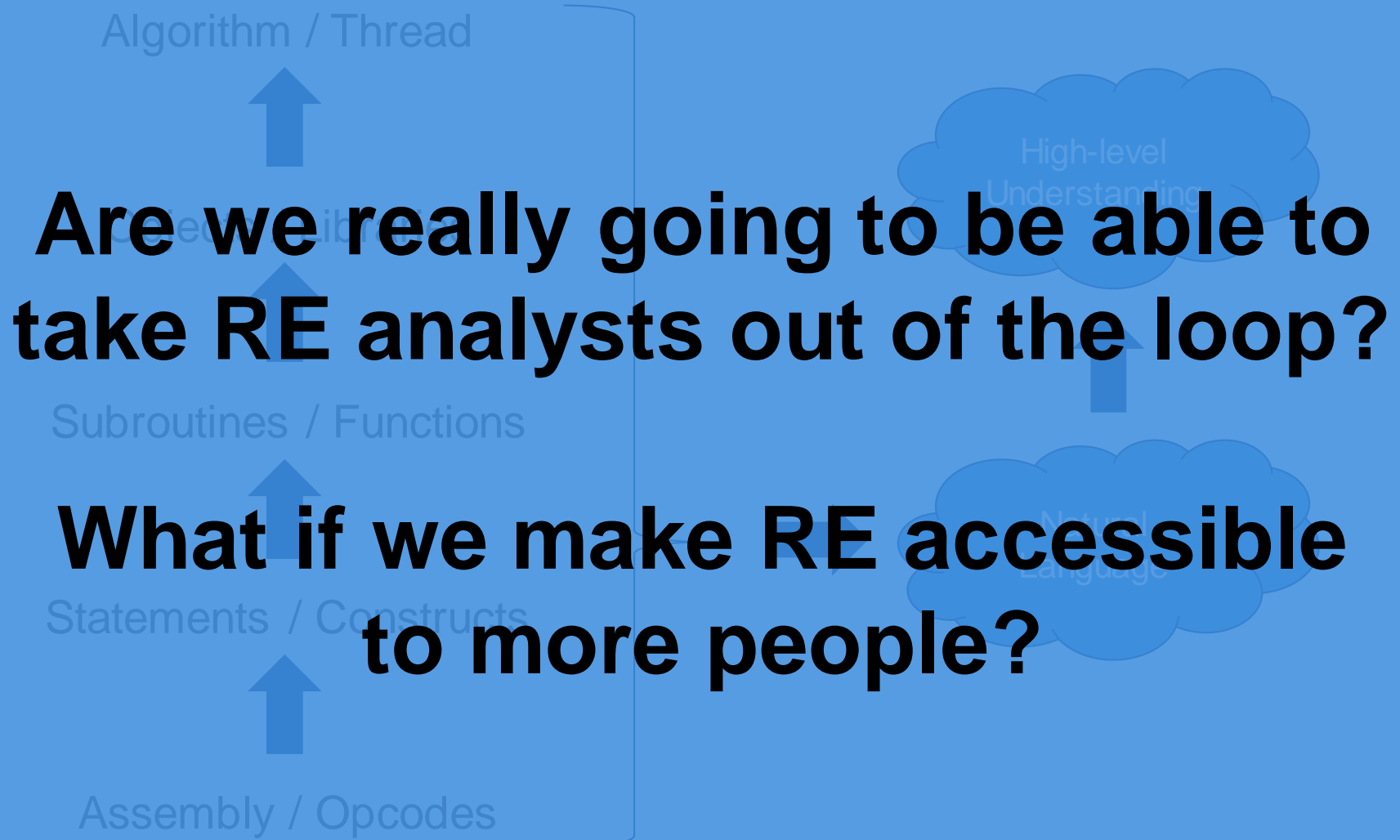
What is Software Reverse Engineering?



Software Reverse Engineering Process



Software Reverse Engineering Process



**What kind of problems
should we work on next?**

What's Next? – Challenge Problems

- **Variable Name Prediction** – Given a blank decompiled function, output meaningful names for the variables in that function
 - Promising research from CMU, ETH Zurich using statistical machine translation (SMT) approach,^[18] neural network approach,^[19] probabilistic model approach^[20]
 - Inspired by work recovering identifiers in obfuscated JavaScript^[21]
- **Statement Commenting** – Given a blank decompiled code statement or fragment, output comments in natural language describing it
 - What is this basic block doing?
 - Existing work in code snippet tagging/labeling^[22]
 - Feed in context information from function
 - Build on variable name prediction
- **Function Summarization** – Given a blank decompiled function, output comments in natural language summarizing the function
 - Language summarization is a difficult problem due to lack of datasets
 - Existing research in source code summarization^{[23][24] [25]}
 - Needs to be built on variable name prediction, possibly statement commenting?
 - Rely on man page / documentation descriptions?

What's Next? – Challenge Problems

- **Library/Object Organization (“CodeCut”)** – Given a fully linked binary with no symbols, locate the boundaries between the original object files – or the boundaries between related sections of code.
 - Preliminary solution using a function affinity metric and a weighted cut graph approach^[26]
 - Source code available for IDA Pro^[27]
- **Library Summarization** – Given an identified set of related decompiled functions, output a name or labels describing it
 - Extension of function-level summarization work
 - Extension of code snippet tagging work
- **Foundational NLP Work**
 - asm2vec^[28]
 - code2vec^[29]

How do we get started?

Assembled Labeled Library for Static Analysis Research (ALLSTAR)

- NLP research will require a large (publicly available) database of:
Open Source Code → Compile Artifacts → Binaries (with Symbols)
- There is a TON of open source code
- There is a TON of firmware out there
- But....there are ZERO debugging symbols
(ok almost zero)
- And *ideally* it would be cross-architecture...(so we avoid overtraining on x86)

ALLSTAR Build Process

- Uses Debian “jessie” distribution
- Try to remove packages that are documentation, debug, data-only, Python, Java, etc. – leaves over 32,000 packages
- Build for all Debian architectures (x86, amd64, ARM, PPC, MIPS, s390x) using Dockcross (Docker containers with cross compilers)
- Fairly simple technical details:
 - CFLAGS: -g -fdump-final-insns -fdump-tree-gimple
 - CXXFLAGS: -g -fdump-final-insns -fdump-tree-gimple -fdump-class-hierarchy
 - DEB_BUILD_OPTIONS: nostrip debug
 - Build with ‘debuild’ setting architecture
- Takes 5 to 6 weeks (?) to build, 35-55 packages per hour

ALLSTAR Stats*

- 1.1 TB total storage
- Binaries:



- i386, MIPS and PPC numbers reflect Debian packages that properly conform to cross-arch build spec
- ARM and s390x lower due to Dockcross' cross compiler (uses a sysroot)
- ~20K cross-platform binaries is still a nice dataset

*projected, build currently in process

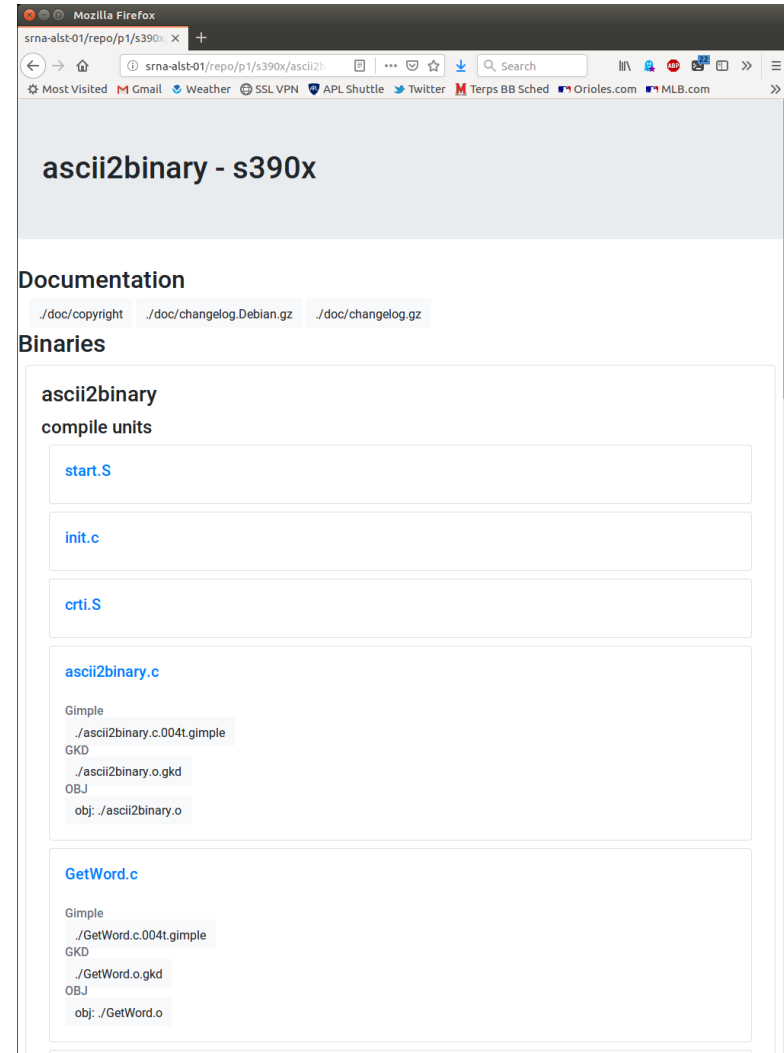
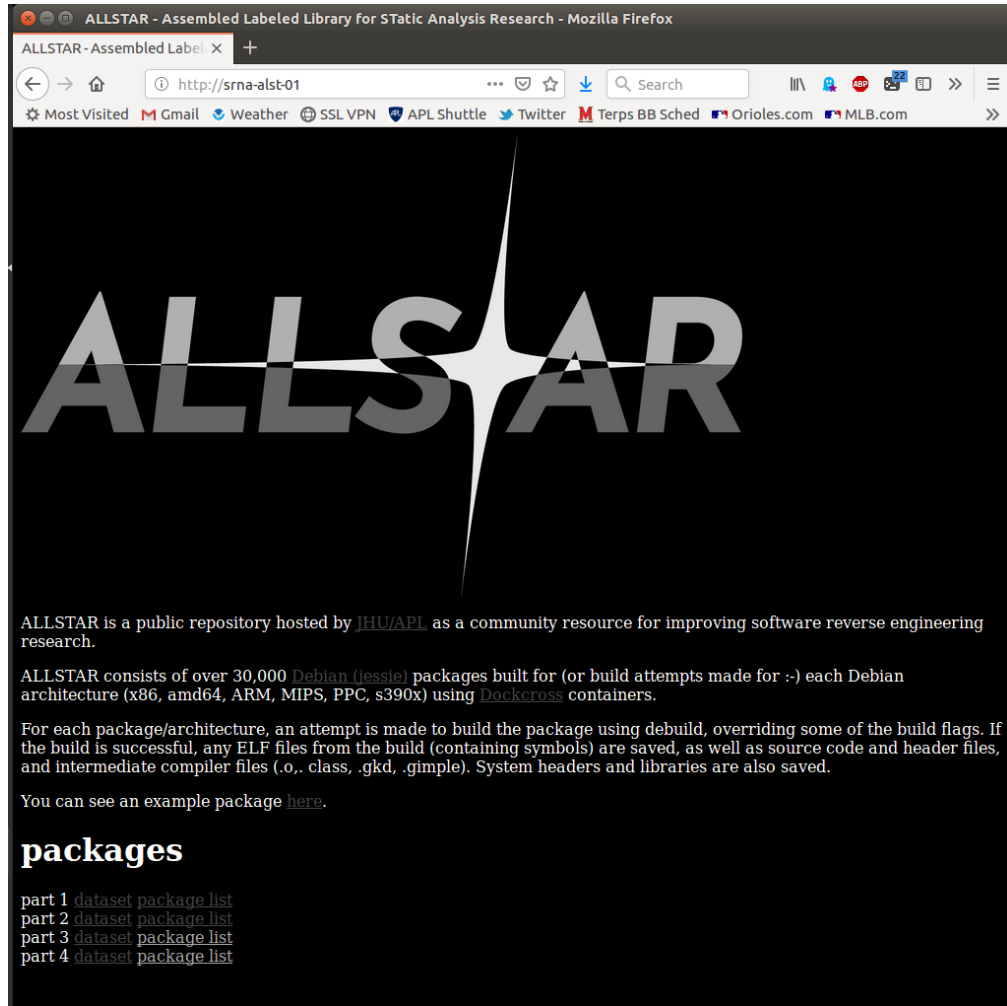
Debian vs. Github

| | |
|--|--|
| ~32K packages (mostly C/C++) | ~400K packages in C ^[18] |
| Structured build process (debbuild) | Build by trying “configure” and “make” |
| Packages build cross-architecture (theoretically, if they follow spec) | Less likely to build for non-amd64 |
| Less duplication | More duplication (up to 70%! ^[30]) |
| GPL or similar open licensing | Licensing unclear / non-obvious |
| More serious/polished code? | A bit more random code? |

What's in ALLSTAR? A Single Data Record

- Binary (fully linked ELF, either runnable or .so)
 - Source file
 - Header files
 - GIMPLE file (like an abstract syntax tree)
 - .class file (if C++)
 - gkd file (final internal representation in RTL)
 - Object file
 - Man page
 - Symbols (in binary) – can be used to produce a .map file
 - Any system library dependencies
- Package Documentation
- HTML index for browsing
- JSON index for code parsing

ALLSTAR – Browser Interface



Building on ALLSTAR

- Open source / publish by end of 2019
- 2 internal research projects for FY20 planned
- Hope to see a lot of research building on ALLSTAR in the future



JOHNS HOPKINS
APPLIED PHYSICS LABORATORY

Special Thanks To:

- halvarflake

Additional Thanks To:

- Igor Skochinsky
- Joxean Koret
- Joan Calvet

Questions?

evm

@evm_sec

evm.ftw@gmail.com

References (1/4)

- 1) Hex-Rays. "IDA: About." <https://www.hex-rays.com/products/ida>.
- 2) National Security Agency. "Ghidra." <https://ghidra-sre.org/>.
- 3) Avast Software. "RetDec::Home." <https://retdec.com/>.
- 4) Trail of Bits. "trail-of-bits/relic," GitHub. <https://github.com/trailofbits/relic>.
- 5) angr Developers. "angr/angr-management," GitHub. <https://github.com/angr/angr-management>.
- 6) PNF Software. "JEB Decompiler by PNF Software." <https://www.pnfsoftware.com>.
- 7) C. Cifuentes. "Reverse Compilation Techniques." Phd diss. Queensland University of Technology, 1994.
- 8) National Security Agency. "Ghidra Pcode." *Ghidra source code*. <https://github.com/NationalSecurityAgency/ghidra/tree/master/Ghidra/Framework/SoftwareModeling/src/main/java/ghidra/pcode>.
- 9) I. Guilfanov. "Decompiler internals: microcode." *RECON 2018*. Brussels, Belgium. <http://www.hex-rays.com/products/ida/support/ppt/recon2018.ppt>.
- 10) J. Calvet. "The (Long) Journey to a Multi-Architecture Disassembler." *RECON 2018*. Montreal, QC. <https://cfp.recon.cx/reconmtl2019/talk/MCEWBB/>

References (2/4)

- 11) zynamics. “BinDiff.” <https://www.zynamics.com/bindiff>.
- 12) J. Koret. “Diaphora.” <https://diaphora.re>.
- 13) National Security Agency. “Version Tracking” *Ghidra Intermediate Training Course*.
https://ghidra.re/courses/GhidraClass/Intermediate/VersionTracking_withNotes.html#VersionTracking.html.
- 14) T. Dullien, and R. Rolles. “Graph-based comparison of executable objects.” *SSTIC 5, no. 1*. 2005.
- 15) T. Dullien, et al. “Automated Attacker Correlation for Malicious Code.” *Information Systems and Technology Panel (IST) Symposium*. Talinn, Estonia. 2010.
- 16) H. Flake and S. Porst. “ShaREing Is Caring.” *CANSEC West 2010*. Vancouver, Canada. 2010. <https://www.slideshare.net/cblichmann/shareing-is-caring-3856188>.
- 17) C. Karamitas, and A. Kehagias. “Efficient features for function matching between binary executables.” *In 2018 IEEE 25th International Conference on Software Analysis, Evolution and Reengineering (SANER)*, pp. 335-345. IEEE, 2018.

References (3/4)

- 18) A. Jaffe, et al. "Meaningful variable names for decompiled code: A machine translation approach." In *Proceedings of the 26th Conference on Program Comprehension* (pp.20-30). ACM.
- 19) J. Lacomis, et al. "DIRE: A Neural Approach to Decompiled Identifier Naming." In *The 34th IEEE/ACM International Conference on Automated Software Engineering*. 2019.
- 20) J. He, et al. "Debin: Predicting debug information in stripped binaries." In *Proceedings of the 2018 ACM SIGSAC Conference on Computer and Communications Security*, pp. 1667-1680. ACM, 2018.
- 21) B. Vasilescu, et al. "Recovering clear, natural identifiers from obfuscated JavaScript names," in *Joint Meeting of the European Software Engineering Conference and the Symposium on the Foundations of Software Engineering*. 2017, pp.683-693.
- 22) B. Gelman et al. "A language-agnostic model for semantic source code labeling." In *Proceedings of the 1st International Workshop on Machine Learning and Software Engineering Symbiosis*, pp.36-44. ACM, 2018.
- 23) J. Moore et al. "A Convolutional Neural Network for Language-Agnostic Source Code Summarization." *arXiv preprint arXiv:1904.00805*, 2019.
- 24) S. Iyer et al. "Summarizing source code using a neural attention model." In *Proceedings of the 54th Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, 2016.

References (4/4)

- 25) X. Hu et. Al. “Deep code comment generation.” *In Proceedings of the 26th Conference on Program Comprehension*, pp.200-210. ACM, 2018.
- 26) evm. “A Code Pirate’s Cutlass: Recovering Software Architecture from Embedded Binaries.” *Shmoocon 2019*. Washington, DC. 2019.
- 27) evm. “CodeCut: Locating Object File Boundaries in IDA Pro.”
<https://github.com/JHUAPL/CodeCut>.
- 28) S. D’Antoine. “asm2vec: Binary Learning for Vulnerability Discovery.” *Jailbreak Security Summit 2018*. Laurel, MD. 2018.
- 29) U. Alon et al. “code2vec: Learning distributed representations of code.” *Proceedings of the ACM on Programming Languages* 3, 2019.
- 30) P. Martins, et al. “50K-C: A dataset of compilable and compiled, Java projects.” In *2018 IEEE/ACM 15th International Conference on Mining Software Repositories (MSR)*, pp 1-5. IEEE, 2018.