



allvm - Binary Decompile

Sandeep Dasgupta
University of Illinois Urbana Champaign
March 30, 2016



Goal & Motivation

Possible Approaches

Our Approach

Ongoing Work

Questions ?



Research Goal

- Research Goal
 - Obtain “richer” LLVM IR than native machine code.
 - Enable advanced compiler techniques (e.g. pointer analysis, information flow tracking, automatic vectorization)



Why “richer” LLVM IR

- Source code analysis not possible
 - IP-protected software
 - Malicious executable
 - Legacy executable
- Source code analysis not sufficient
 - What-you-see-is-not-what-you-execute
- Platform aware optimizations



Goal & Motivation

Possible Approaches

Our Approach

Ongoing Work

Questions ?



The 3 Possible Approaches

- Decompile Machine Code \rightarrow LLVM IR
- “Annotated” Machine Code \rightarrow LLVM IR
- Ship LLVM IR



Decompile Machine Code \rightarrow LLVM IR

Benefits	Challenges
<ul style="list-style-type: none">• Easy to adopt• No compiler support needed	<ul style="list-style-type: none">• Reconstructing code and control flow• Variable recovery• Function & ABI rules recovery

- Tools Available: QEMU, BAP, Dagger, Mcsema, Fracture



“Annotated” Machine Code \rightarrow LLVM IR

Benefits	Challenges
<ul style="list-style-type: none">• Effective reconstruction• Minimal compiler support needed	<ul style="list-style-type: none">• Annotations must be<ul style="list-style-type: none">• Minimal• Compiler & IR independent• Adoption

- Tools Available: None



Ship LLVM IR

Benefits	Challenges
<ul style="list-style-type: none">• <i>No loss</i> of information	<ul style="list-style-type: none">• Adoption in Non LLVM based compilers• Code size bloat

- Tools Available: Portable Native Client, Renderscript, iOS, watchOS and tvOS apps.



Goal & Motivation

Possible Approaches

Our Approach

Ongoing Work

Questions ?



Our Approach

Long term goal

Minimal compiler-independent annotations to reconstruct high-quality IR

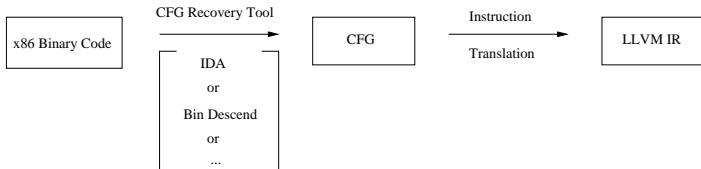
Short term goals

- ❶ Experiment with Machine Code \rightarrow LLVM IR, to **understand** the challenges better
 - To select from existing decompilation frameworks.
 - Experiment with different variable and type recovery strategies
- ❷ Design suitable annotations for what cannot be inferred without them



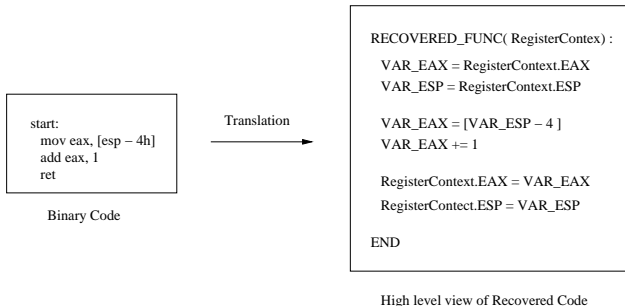
Selecting mcsema

- Actively supported and open sourced
- Well documented
- Functional LLVM IR
- Separation of modules: CFG recovery and CFG \rightarrow LLVM IR





Instruction Translation: Memory Model





mcsema: Demo



Support & Limitations

- What Works
 - Integer Instructions
 - FPU and SSE registers
 - Callbacks, External Call, Jump tables
- In Progress
 - FPU and SSE Instructions: Not fully supported
 - Exceptions
 - Better Optimizations



Goal & Motivation

Possible Approaches

Our Approach

Ongoing Work

Questions ?



Variable & Function Parameter Recovery

- Enables
 - Fundamental analysis (Dependence, Pointer analysis)
 - Optimizations (register promotion)
- State of the art
 - Divine
 - State of the art variable recovery
 - Second Write
 - Heuristics for function parameter detection
 - Scalable variable and type recovery
 - TIE
 - Type recovery



Goal & Motivation

Possible Approaches

Our Approach

Ongoing Work

Questions ?



Questions ?



Goal & Motivation

Possible Approaches

Our Approach

Ongoing Work

Questions ?



What-you-see-is-not-what-you-execute

The following compiler (Microsoft C++ .NET) induced vulnerability was discovered during the Windows security push in 2002

```
memset(password, '\0', len);  
free(password);
```

—————>

```
free(password);
```