

May 6-7, 2021

ARSENAL



### **Ghidra-EVM: Reversing Smart Contracts with Ghidra**

Antonio de la Piedra



### Motivation



- Importance of validating the security of deployed contracts
- Push the limits of Ghidra for analyzing EVM byte code



### Agenda



#### Introduction

Ethereum and the EVM Ghidra extension capabilities

#### **Ghidra-EVM Architecture**

Related work
Design of ghidra-EVM
Modules and dependencies
Challenges and limitations

#### **Demo**

Recovering the CFG Analyzing creation bytecode Solving a challenge

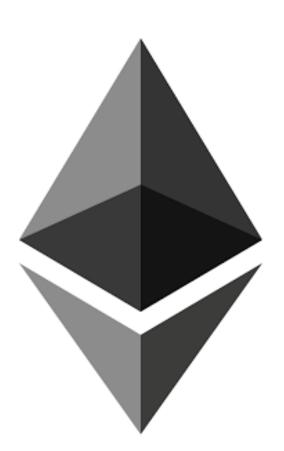


### Introduction



#### Ethereum

- The Ethereum blockchain is the second largest blockchain after bitcoin.
- Support of smart contracts, executable code located at an address of the blockchain that is triggered through a transaction.
- Deployed contracts cannot be changed and updates are usually complex.
- In the past, a considerable amount of funds was stolen due to
- programming errors, such as the DAO attack, or the King of the Ether Throne attack.





#### **Smart contracts**

- It is possible to deploy smart contracts at an address in the Ethereum network.
- It provides different functions that are triggered through a transaction in the blockchain
- The contracts are generally written in languages such as solidity although in the last few years security-oriented constrained languages have been proposed such as LLL and Vyper.





#### The Ethereum Virtual Machine (EVM)

- Compiled smart contracts into bytecode are executed in the Ethereum Virtual Machine (EVM).
- Stack machine, that means that typically each instruction pops the operands from the stack and pushed back the result after the instruction execution.
- The machine works with 256-bit words and utilized two memory units: "memory" and storage.





#### Ghidra

• In 2019, the National Security Agency released the Ghidra reverse engineering tool, equipped with a decompiler and plugin support.





#### **Ghidra Extension Capabilities**

- **Processor module**: where the different instructions of a specific processor are described in the SLEIGH language.
- Loader: detects the object code that is loaded into Ghidra
- and relates it to a certain processor and compiler configuration.
  Plugins: provides access to the Ghidra API for modifying memory blocks, adding labels, modifying the object code, locating instructions and performing analysis on the disassembled code: Java, Python 2 and Python 3 (via ghidra-bridge).



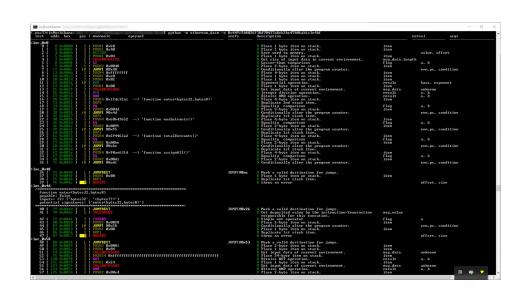


# Architecture of Ghidra-EVM



#### Related work

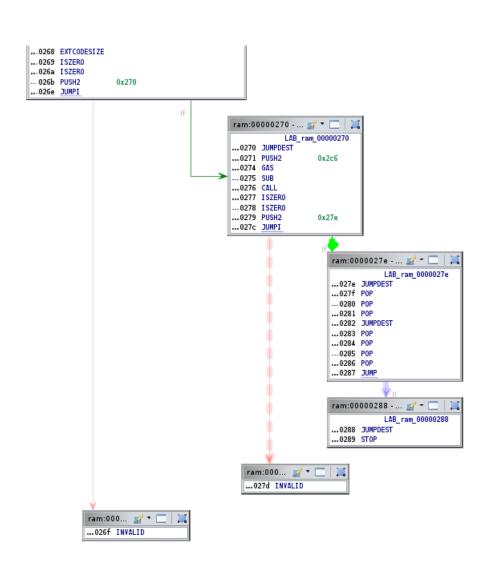
- **Recovering the CFG**: ethereum-graph-debugger, Consensys surya (solidity) and evm cfg builder: retrieves the CFG of a compiled contract in EVM bytecode.
- Disassemble plugins for EVM byte code: IDA Pro and Binary Ninja.
- Binary analysis frameworks: Rattle, ethereum-dasm, pyevmasm (library).





#### Related work

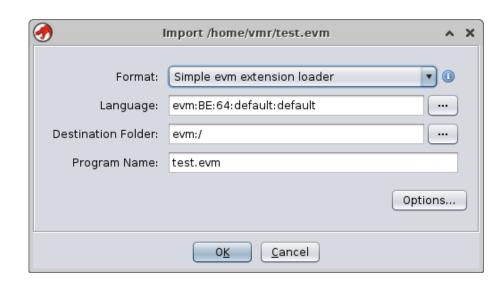
- On Nov 6, 2019 Crytic opened a bounty for building a Ghidra plugin that, using evm cfg builder, could retrieve a show the CFG of a smart contract in Ghidra.
- However, this bounty was cancelled months after without a proposed solution (https://github.com/crytic/evm\_cfg\_builder/issues/23).
- This work, inspired by the cancelled bounty, proposes an architecture to reverse engineering smart contracts in Ghidra.





#### Design of Ghidra-EVM: Loader

- Ghidra-EVM expects an hexadecimal string representing bytecode or a binary file with bytecode:
  - Solc's output
  - Downloaded from the blockchain
- The loader expects a file with the extension .evm or .evm\_h.





#### Design of Ghidra-EVM: Processor Module

- Described in 5 different files:
  - **Cspec**: Compiler specification.
  - Ldefs: Processor identification and language definitions.
  - Opinion: links the loader to the compiler definition in Idefs.
  - Psec: Defines the size of the memory blocks.
  - **Slaspec**: Register definitions and instructions:
  - define space register type=register\_space size=\$(SIZE) wordsize=1;
     define register offset=0x00 size=\$(SIZE) [sp pc];



#### Design of Ghidra-EVM: Memory model

- In the EVM, a smart contract relies on three main memory structures:
  - **Storage memory**: of 256 bits per slot, utilized by data structures in the smart contract defined as storage.
  - Volatile memory: utilized for storing computations, and data structures
  - defined as memory and the freepointer data structure.
  - Stack: utilized by the majority of instructions since the arguments
  - are typically pushed and poped to/from the stack.



#### Design of Ghidra-EVM: Memory model

• The storage and volatile memories are defined in Ghidra-EVM to aid in the decompilation process.

• The third memory block, evm\_jump\_table that is used to generate, the Control Flow Graph (CFG) and to resolve the jump destinations of the JUMP and JUMPI instructions.



#### Design of Ghidra-EVM: Instruction definition

- Around 140 opcodes in the Ethereum yellow paper.
- Defined in SLEIGH in the slaspec file e.g.

```
:AND is op=0x16 {
    __value1 :$(SIZE) = 0;
    __value2 :$(SIZE) = 0;
    __value3 :$(SIZE) = 0;

    pop(_value1);
    pop(_value2);

    __value3 = __value1 & __value2;

    push(_value3);
}
```



#### Design of Ghidra-EVM: Instruction definition

• Instructions that obtain data from the blockchain are defined as pcodeop e.g.

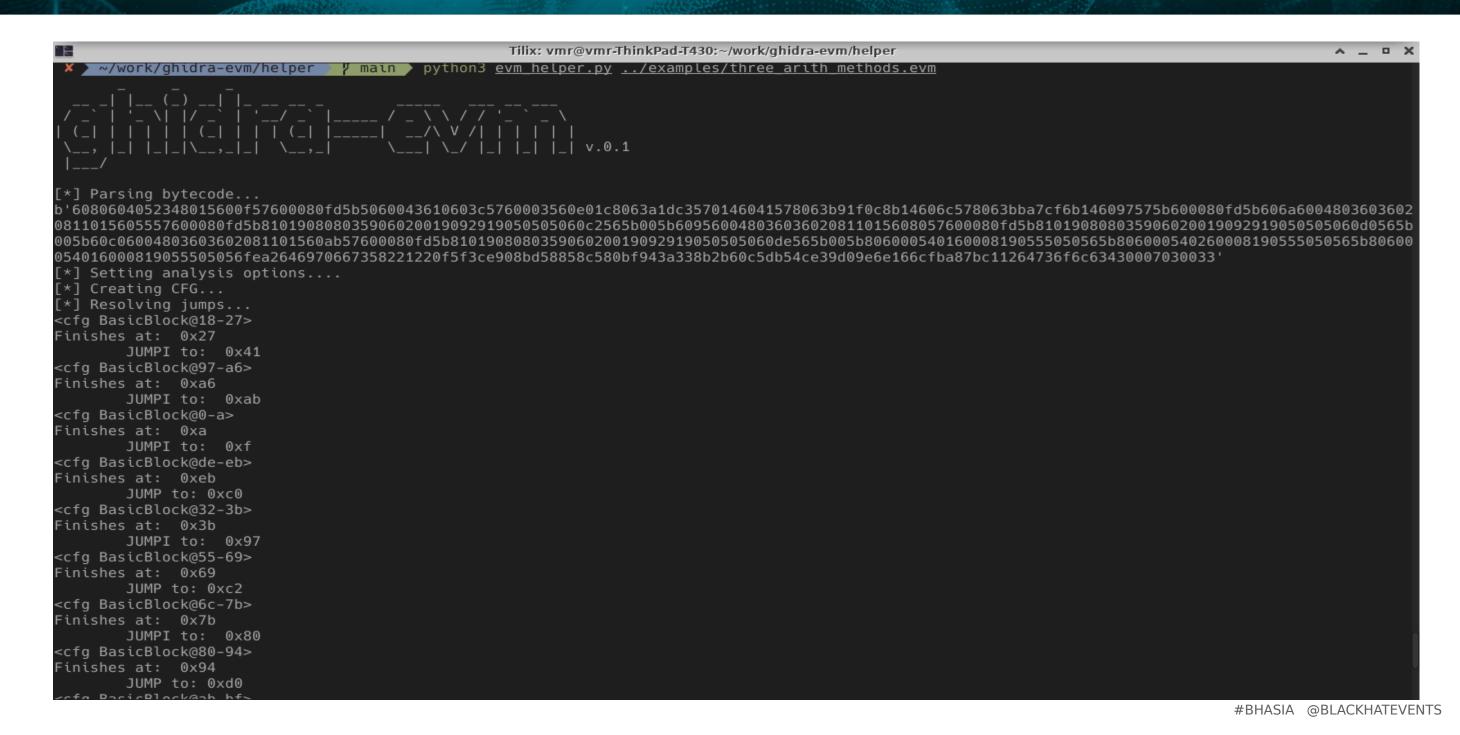
```
define pcodeop get_block_gas_limit;
:GASLIMIT is op=0x45 {
    _value1 :$(SIZE) = 0;
    _value1 = get_block_gas_limit();
    push(_value1);
```



#### Design of Ghidra-EVM: CFG and branches

- evm\_cfg\_builder is used to recover the CFG
  - Python 3 code: ghidra\_bridge is used (RPC proxy for Python objects)
  - A jump table is created during analysis: contains the destination address of every jump in the disassembled code
- The SLEIGH definitions of JUMP and JUMPI recover the JUMPDEST
- location via the jump table.
- Through evm\_cfg\_builder contract methods are defined as functions in Ghidra, method names and properties are extracted and showed in the disassembler.

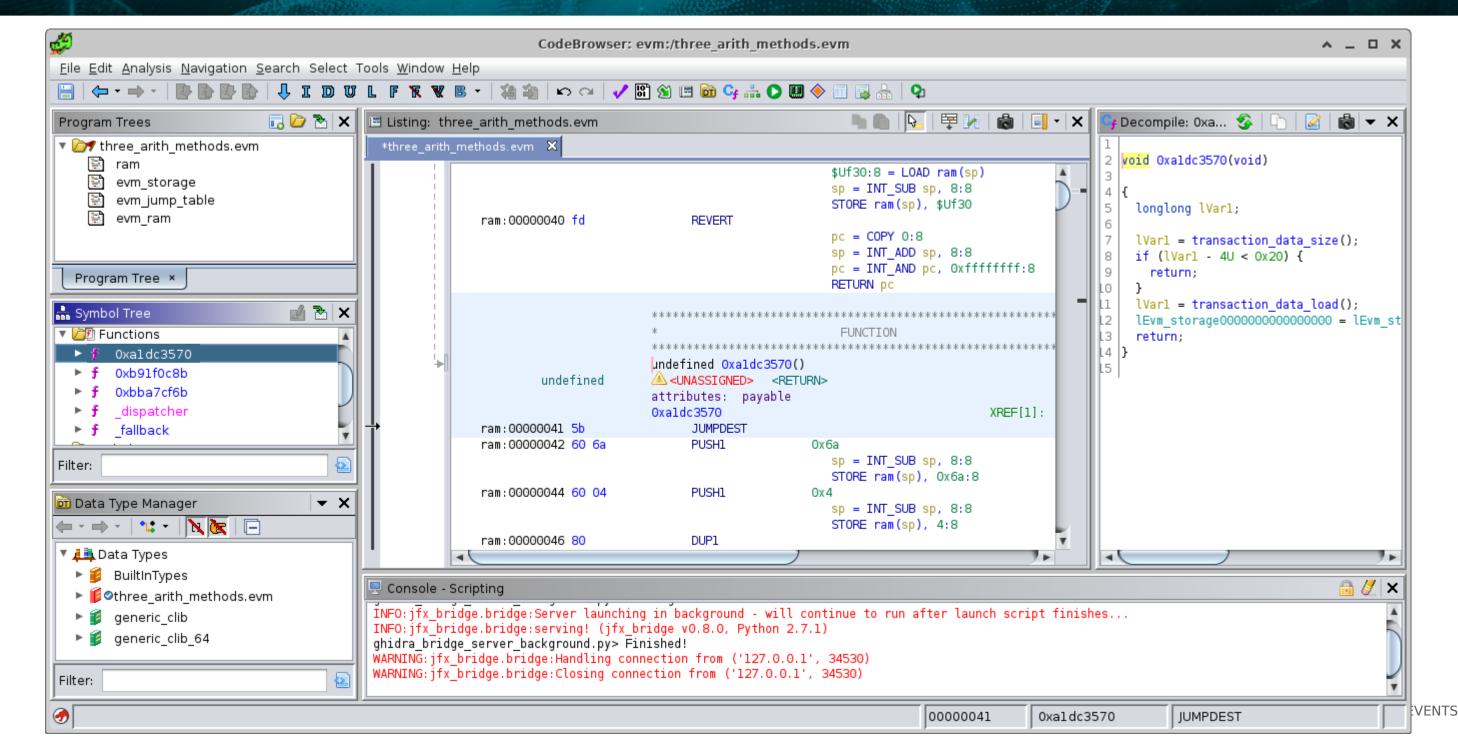




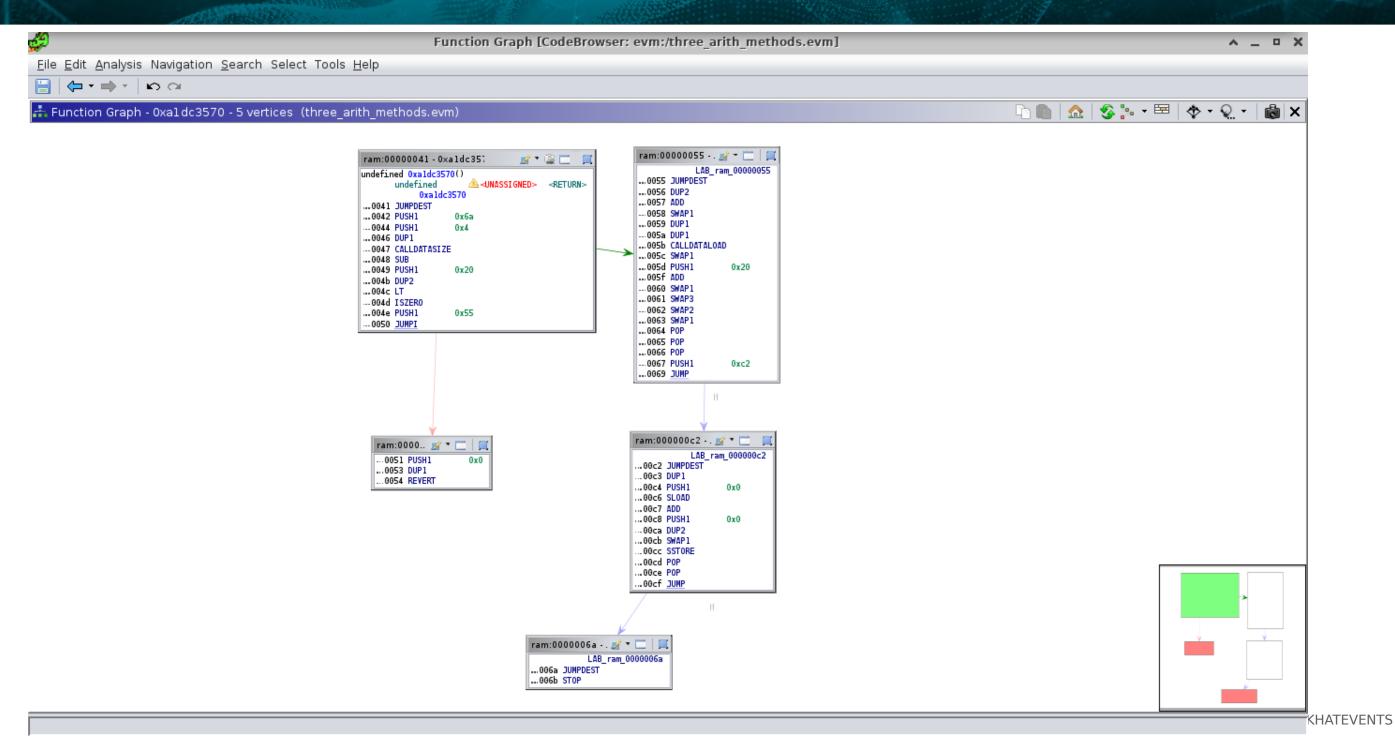


```
Tilix: vmr@vmr-ThinkPad-T430:~/work/ghidra-evm/helper
        JUMP to: 0xc0
<cfg BasicBlock@32-3b>
Finishes at: 0x3b
       JUMPI to: 0x97
<cfg BasicBlock@55-69>
Finishes at: 0x69
        JUMP to: 0xc2
<cfg BasicBlock@6c-7b>
Finishes at: 0x7b
        JUMPI to: 0x80
<cfg BasicBlock@80-94>
Finishes at: 0x94
        JUMP to: 0xd0
<cfg BasicBlock@ab-bf>
Finishes at: 0xbf
        JUMP to: 0xde
<cfg BasicBlock@c2-cf>
Finishes at: 0xcf
       JUMP to: 0x6a
<cfg BasicBlock@28-31>
Finishes at: 0x31
       JUMPI to: 0x6c
<cfg BasicBlock@d0-dd>
Finishes at: 0xdd
       JUMP to: 0x95
<cfg BasicBlock@41-50>
Finishes at: 0x50
       JUMPI to: 0x55
<cfg BasicBlock@f-17>
Finishes at: 0x17
        JUMPI to: 0x3c
[*] Exploring functions...
       Found function _dispatcher
       Found function fallback
       Found function 0xa1dc3570
       Found function 0xb91f0c8b
       Found function 0xbba7cf6b
[*] Analyzing....
*] Disassemble all....
~/work/ghidra-evm/helper / main
```











#### Design of Ghidra-EVM: This version

- Initially designed for Ghidra 9.1.2, tested in Linux
- Disassemble and CFG recovery (via ghidra\_bridge and cfg evm builder)
  Identification of contract methods and labeling as functions in Ghidra
- Identification of method properties (via cfg evm builder)
- Identification of typical method names (via cfg evm builder)



#### Design of Ghidra-EVM: This version

- Identification of problematic instructions and labelingExperimental decompilation



#### Design of Ghidra-EVM: Limitations and challenges

- Ghidra seems to be designed to work with up 64-bit architectures in mind
- It's difficult to support a 256-bit architecture in many ways:
  - Memory definitions with that word size fail to compile. Affects storage and volatile memory.
  - Same for instructions with large operands in SLEIGH: affects decompilation output.
  - This limits the analysis to a restricted size of a EVM contract.
  - Contributions and help appreciated:)



## Demo 1: Simple utilization



# Demo 2: Analyzing creation code



# Demo 3: Extending Ghidra-EVM