

SYMBOLIC PROGRAM SLICING ON SMART CONTRACTS

Zi Xuan Chen, Fang Yu

National Chengchi University
Department of Computer Science, Management Information Systems
104703010@nccu.edu.tw, yuf@nccu.edu.tw

ABSTRACT

We propose a method to do program slicing on stack-based programming language. With slicing, we can analyse properties more efficiently. We take the EVM bytecode[1] as the example language, which is used to write smart contract[2] on Ethereum[1].

Keywords— BlockChain, Ethereum, Slicing, Verification

1. INTRODUCTION

There are many interest properties about Ethereum. Ethereum can be viewed as a transaction-based state machine. We can transit the transaction state by sending transaction or execute smart contract. On Ethereum, all transaction executed on Ethereum Virtual Machine (EVM), which is a simple stack-based architecture, limits stack item size to 1024. The word size of the machine (and thus size of stackitems) is 256-bit. Executions will be failed if stackoverflow occurred while executing the smart contract. Besides stack, *gas* is another interesting property on Ethereum. To limit the cost of the transaction execution, EVM takes the handling fee named *gas* from transaction sender.

To analyze the properties more precisely, we slice the smart contract by constructing dependency graph (DG). With the dependencies, we can slice a smaller program from some interested point. Then we can the the sliced program for other analysis purpose.

2. RELATED WORK

TBC.

3. METHOD

To compute the dependencies between instructions, we need construct the control flow graph (CFG) first. Unlike register-based machine's instruction, which's operand are called as register explicitly, for the stack-based machine, the operand that instructions depended are stored on the stack implicitly. Thus, a CFG for stack simulation is needed.

Algorithm 1: buildCFGandStackDependency

Input: *Opcode*

Output: *CFG, StackDG*

```
1 function buildCFGandStackDependency(opcode)
2   dg, cfg = DG(opcode), CFG(opcode);
3   cfg.buildBasicBlocks();
4   cfg.buildSimpleEdges();
5   cfg.buildFunctions(cfg.basicBlocks.first);
6   for func ∈ cfg.functions do
7     | valueSetAnalysis(cfg, dg, func);
8   end
9   return cfg, dg;
10 end function
    /* state constructor */
11 function State()
12   this.visit = dict(dflt=0);
13   this.stacksIn = dict(dflt=None);
14   this.stacksOut = dict(dflt=None);
15   this.discoveredTargets = dict(dflt=0);
16   this.lastDiscoveredTargets = dict(dflt=0);
17 end function
    /* value set analysis */
18 function valueSetAnalysis(cfg, dg, func)
19   stat = State();
20   toExplore = { func.entry };
21   do
22     outBlocks = { toExplore.pop() };
23     do
24       | outBlocks = outBlocks ∪
25         |   transFuncBlock(cfg, dg, func,
26           |   outBlocks.pop(), stat);
27     while outBlocks;
28     for src, dsts ∈ stat.lastDiscoveredTargets do
29       | cfg.addEdges(src, dsts);
30       | toExplore = toExplore ∪ dsts;
31     end
32     stat.visit = dict(dflt=0);
33     stat.lastDiscoveredTargets = dict(dflt=0);
34   while toExplore;
35 end function
```

The first step to construct CFG is splitting basic blocks. We split basic blocks by **JUMPDEST** and **end instructions**. **JUMPDEST** is normally considered as the beginning of blocks because other blocks can target **JUMPDEST** to connect the edges. The **end instructions** include **STOP**, **SELF-DESTRUCT**, **RETURN**, **REVERT**, **INVALID**, **SUICIDE**, **JUMP**, **JUMPI**.

The second step is building the edges between basic blocks. Some edges can be computed by simply succeeding the pc of the **JUMPI** and other instructions \notin **end instructions**, followed by **JUMPDEST**. Because the property of the stack-based machine, all the jump destinations are pushed to stack implicitly. For this problem, we use value set analysis (VSA) to find all the possible destination values.

Most of smart contracts are written in Solidity, which is an object-oriented (or contract-oriented), high-level language for implementing smart contracts. The contract in Solidity is like an object. Users can call the public functions in the contract, which we can treat as member functions in an object. From lower-level — EVM bytecode, the Solidity compiler will compile a dispatcher to dispatch public functions in the contract. The dispatcher can recognize the function hashes in transactions that users sent via Application Binary Interface (ABI). We compute the function boundaries and apply the VSA on each function to construct complete contract flow graph and instruction dependencies.

In the VSA, we traverse the basic blocks in CFG and continue finding new target address of next block by simulate the execution of instructions with an abstract stack. Abstract stack abstracts all the possible stack state. The n -th item in abstract stack represent a set of all possible value of n -th item in all possible stack state. So it is an over-approximation to compute the jump destinations. For each block, we record the states of the abstract stack before and after executing the instructions in the block. With the states, we can check the convergence of the analysis. If we revisit a block, and its post-execution state is same as last time, we consider it achieve the convergence. We assume no new value would be found. Note that only the **PUSH**, **SWAP**, **DUP**, **AND** are implemented here, for other instructions, we only do the push, pop on the stack based on the operation times it defined. It's because other instructions implementation will not affect the target address computation.

The instructions dependencies are also constructed while doing VSA. To build the dependencies, we need keeping track of the data flow of operands. All the operands are pushed into and popped from the stack. Instead of marking the operands with the instructions which pushed it, we push the instructions to the stack directly, and edges are added when the instructions are popped. Note that we don't use abstract stack here. Instead, we use a set of stack to keep all the states of possible stacks to maintain the accuracy of each operand list. If we use abstract stack here, more combination of operand list will be generated. It will lead more ambiguous result for address evaluation while building memory dependency.

Algorithm 2: ValueSetAnalysisUtility

```

/* trasfer function blocks */
1 function transFuncBlock(cfg, dg, func, block, stat)
2   if (func.id == DISPATCHER_ID
3     and block.reachable)
4     or stat.visit[block] > visitLimit then
5     return;
6   stat.visit[block] += 1;
7   /* save pre-stack to check convergence */
8   prevStack, _ = stat.stacksOut[block]
9   oprdStack, instStack = abtStack(), listStack();
10  inBlocks = [b ∈ block.inBlocks
11    | stat.stacksOut[b] ≠ None]
12  for father ∈ inBlocks do
13    ostk, istk = stat.stacksOut[father];
14    oprdStack = oprdStack.merge(ostk);
15    instStack = oprdStack.merge(istk);
16  end
17  /* explore the block */
18  exploreBlock(dg, block,
19    oprdStack, instStack, stat);
20  /* add branch according the result */
21  if block.end ∈ {JUMP, JUMPI} then
22    oprdStack, _ = stat.stacksIn[end];
23    for dst ∈ oprdStack.top().vals() do
24      if isJumpDest(dst) then
25        addBranch(src, dst, stat);
26    end
27    oprdStack, _ = stat.stacksOut[end];
28  if prevStack ≠ oprdStack then
29    /* not converged */
30    return block.outBlocksByFunc(func.id);
31  return ∅;
32 end function
33 /* explore basic block */
34 function exploreBlock(dg, block, oprdStack,
35  instStack, stat)
36  for inst ∈ block.instructions do
37    stat.stacksIn[inst] = (oprdStack, instStack);
38    stat.stacksOut[inst] = transferFuncInst(
39      dg, inst, oprdStack, instStack, stat);
40  end
41 end function
42 /* add branch to value set analysis */
43 function addBranch(src, dst, stat)
44  if dst ∉ stat.discoveredTargets[src] then
45    if src ∉ stat.lastDiscoveredTargets then
46      stat.lastDiscoveredTargets[src] = ∅;
47      stat.lastDiscoveredTargets[src].add(dst);
48      stat.discoveredTargets[src].add(dst);
49  end

```

Algorithm 3: ValueSetAnalysisUtility

```

/* transfer instruction */
1 function transferFuncInst(dg, inst,
2   oprdStack, instStack, stat)
3   oprdStack = oprdStack.copy();
4   instStack = instStack.copy();
5   if inst ∈ PUSHn[n=1..32] then
6     oprdStack.push(inst.operand);
7     instStack.push(inst);
8   else if inst ∈ SWAPn[n=1..16] then
9     oprdStack.swap(n);
10    instStack.swap(n);
11   else if inst ∈ DUPn[n=1..16] then
12     oprdStack.dup(n);
13     instStack.dup(n);
14   else if inst = AND then
15     v1, v2 = oprdStack.pop(), oprdStack.pop();
16     oprdStack.push(absAnd(v1, v2));
17     v1s, v2s = instStack.pop(), instStack.pop();
18     for v1, v2 ∈ zip(v1s, v2s) do
19       dg.addEdges(inst, [v1, v2]);
20     end
21     instStack.push(inst);
22   else
23     repeat inst.popTimes times
24       | oprdStack.pop();
25     end
26     for args ∈ [instStack.pop()
27       | n ∈ range(inst.popTimes)]T do
28       | dg.addEdges(inst, args);
29     end
30     repeat inst.pushTimes times
31       | oprdStack.push(None);
32       | instStack.push(inst);
33     end
34   return oprdStack, instStack;
35 end function

```

After building the instruction dependency with stacks by value set analysis. we start to build the dependency between instruction by address of memory and storage.

Memory and storage are other two main data read/write mechanisms on Ethereum except stack, where the memory is volatile, the storage is non-volatile. By the way, according the figure below, we can find that the EVM does not follow the standard von Neu-mann architecture. Rather than storing program code in generally-accessible memory or storage, it is stored separately in a virtual ROM interactable only through a specialised instruction.[1]

There is an indexer recording current memory used, it indicates the maximum index of current memory used. The total fee for memory-usage payable is proportional to smallest

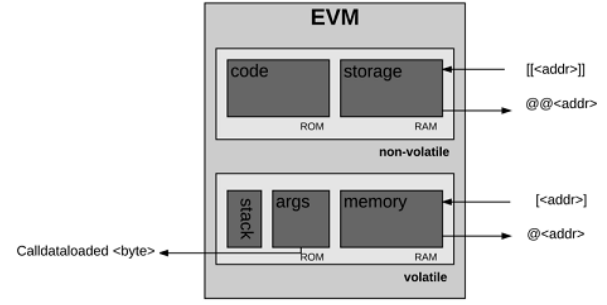


Fig. 1. EVM architecture

multiple of 32 bytes that are required such that all memory indices (whether for read or write) are included in the range[1].

Storage fees have a slightly nuanced behaviour—to incentivise minimisation of the use of storage (which corresponds directly to a larger state database on all nodes), the execution fee for an operation that clears an entry in the storage is not only waived, a qualified refund is given; in fact, this refund is effectively paid up-front since the initial usage of a storage location costs substantially more than normal usage[1].

The main idea to construct address dependency is traversing the CFG from write instructions with a address, and build the dependencies with the instructions which read the address, until meet the next instructions which rewrite the address. With the DG we constructed, we can evaluate some address values and the stored values from known constants by the dependencies. There are some instructions about the read/write operation on memory and storage. For the storage, the write instructions is **SSTORE**, the load instruction is **SLOAD**. For the memory, the write instruction is **MSTORE**, the read instruction include **MLOAD**, **SHA3**, **CREATE**, **CALL**, **RETURN**.

For write instructions, there are three parts we concerned: the address it stored, the range of memory (offset) it covered and the value it wrote. For each part, if we could eval it as constants from other constants (normally come from the terminal of DG — **PUSH** instruction), we save the values as a concrete value set for the part. if not, all the instruction that it needed to do evaluation, would be saved as a dependant instruction set, once all the instruction in the set are evaluated, the part could be evaluated again to get the concrete values. Same as the write instructions, the read instructions also have these part. But for the values it load, are depended on the write instructions' value which have the same address as it, that's just the dependency we want to construct.

To build the address dependency, we new an environment which contains the three parts for the read/write instructions. we eval the storage and memory separately, but in some situations, the **SSTORE** will dependent on some **MLOAD** instructions, so we put both part into same evaluation loop.

Algorithm 4: eval

Input: *instruction, visit***Output:** *concrete values, dependent instructions*

```
1 function eval(inst, visit)
2   if inst ∈ visit then
3     return ∅, {inst};
4   visit.add(inst);
5   concrete, dependent = ∅, ∅;
6   if inst.name.startswith('PUSH') then
7     return {int(op.operand)}, ∅;
8   cons, deps = {map(eval(_, visit), argList)
9                 | argList ∈ inst.argLists}T;
10  for argList ∈ cons do
11    val = None;
12    if None ∈ argList then
13      continue;
14    else if inst.name = 'ADD' then
15      val = let x, y = argList in x + y;
16    else if inst.name = 'SUB' then
17      val = let x, y = argList in x - y;
18    else if inst.name = 'MUL' then
19      val = let x, y = argList in x * y;
20    else if inst.name = 'DIV' then
21      val = let x, y = argList in x / y;
22    else if inst.name = 'EXP' then
23      val = let x, y = argList in xy;
24    else if inst.name = 'ISZERO' then
25      val = let [x] = argList in  $\begin{cases} 0, & \text{if } x = 0 \\ 1, & \text{otherwise} \end{cases}$ 
26    else if inst.name = 'NOT' then
27      val = let [x] = argList
28          in (1 << 256) - 1 - x;
29    else if inst.name = 'AND' then
30      val = let x, y = argList in x & y;
31    else if inst.name = 'OR' then
32      val = let x, y = argList in x | y;
33    else if inst.name = 'EQ' then
34      val = let x, y = argList in x = y;
35    else if inst.name ∈
36      {'MLOAD', 'SLOAD', 'SHA3'} then
37      concrete.update(env.conVals[inst]);
38    else
39      /* SHA3 not impl yet */
40      throw Exception("not handle the inst yet");
41    if val ≠ None then
42      concrete.add(val);
43      dependent.update(concat(deps));
44  end
45  visit.remove(inst);
46  return concrete, dependent;
47 end function
```

In the environment, there is a set name *evald*, which contains the instructions which's all possible values we have already evald and stored in the concrete set, and their dependent instruction set are empty. For every evaluation round of storage and memory we check the *evald* set to determine if it achieved the convergence. If there are no new instruction added to *evald*, the evaluation will terminate.

Algorithm 5: AnalysisEnvironment

```
/* Environment constructor */
1 function Environment(stackDg)
2   rInsts = stackDg.rInsts;
3   wInsts = stackDg.wInsts;
4   vals = { i, eval(i.vals, ∅) | i ∈ wInsts };
5   addrs = { i, eval(i.addrs, ∅) | i ∈ rInsts ∪ wInsts };
6   offsets = { i, eval(i.offsets, ∅) | i ∈ rInsts ∪ wInsts };
7   /* eval instructions' parameters (val) */
8   this.conVals = { i: con | (i, (con, _)) ∈ vals };
9   this.valDepInsts = { i: dep | (i, (_, dep)) ∈ vals };
10  /* eval instructions' parameters (addr) */
11  this.conAddrs = { i: con | (i, (con, _)) ∈ addrs };
12  this.addrDepInsts = { i: dep | (i, (_, dep)) ∈ addrs };
13  /* eval instructions' parameters (offset) */
14  this.conOffsets = { i: con | (i, (con, _)) ∈ offsets };
15  this.offsetDepInsts = { i: dep | (i, (_, dep)) ∈ offsets };
16  /* write insts that can't be re-evald */
17  this.evald = { i ∈ addrDepInsts | addrDepInsts[i] = ∅ }
18              ∪ { i ∈ offsetDepInsts | offsetDepInsts[i] = ∅ }
19              ∪ { i ∈ valDepInsts | valDepInsts[i] = ∅ }
20 end function
21 /* Return All dependant program counters */
22 function depInsts(env, inst)
23   return env.addrDepInsts[inst]
24         ∪ env.offsetDepInsts[inst]
25         ∪ env.valDepInsts[inst]
26 end function
27 /* check insts if have same addr parameters */
28 function addrOverlap(env, instA, instB)
29   rangeA = product(env.conAddrs[instA],
30                   env.conOffsets[instA]);
31   rangeB = product(env.conAddrs[instB],
32                   env.conOffsets[instB]);
33   for (Aa, Ao), (Ba, Bo) ∈ product(rangeA, rangeB) do
34     if {Aa..Aa + Ao} ∩ {Ba..Ba + Bo} ≠ ∅ then
35       return True;
36   end
37   return False;
38 end function
```

To build the address dependency, we need to evaluate the address of the instruction first. The eval function return both the concrete value set and the dependent instruction set. The argument list is a list of instruction that the current instruc-

Algorithm 6: buildAddressDependency

Input: *CFG, StackDG***Output:** *AddressDG*

```
1 import AnalysisEnvironment as env
2 function buildAddressDependency(cfg,
3   stackDg, opcode)
4   /* declare and alias variables */
5   addrDg = DG(opcode);
6   visit, alter =  $\emptyset, \emptyset$ ;
7   swrites = { inst  $\in$  stackDg | inst = SSTORE };
8   sreads = { inst  $\in$  stackDg | inst = SLOAD };
9   mwrites = { insts  $\in$  stackDg | inst = MSTORE };
10  mreads = { inst  $\in$  stackDg | inst  $\in$  { MLOAD,
11    SHA3, CREATE, CALL, RETURN } };
12  /* new a environment */
13  env = Environment(stackDg);
14  /* evaluate until it converge */
15  while True do
16    ealed = env.ealed.copy();
17    buildDependency(addrDg,
18      swrites, sreads, visit);
19    buildDependency(addrDg,
20      mwrites, mreads, visit);
21    if exist write inst can be re-ealed then
22      for inst  $\in$  re-ealed do
23        update Environment variables
24        in env with
25        eval(instruction parameters)
26      end
27    if env.ealed  $\setminus$  ealed =  $\emptyset$  then
28      break;
29  end
30  return addrDg;
31 end function
32 /* build dependency with write instructions */
33 function buildDependency(addrDg, writes, reads, visit)
34  concrete = { inst  $\in$  (writes  $\setminus$  visit)
35    | depInsts(env, inst) =  $\emptyset$  };
36  while concrete  $\neq \emptyset$  do
37    for inst  $\in$  concrete do
38      block = CFG.blockOf(inst);
39      dfsCFG(addrDg, inst, block,
40        writes, reads,  $\emptyset$ );
41    end
42    visit.update(concrete);
43    for inst  $\in$  (writes  $\setminus$  visit) do
44      if (depInsts(env, inst)  $\setminus$  env.ealed) =  $\emptyset$ 
45      then
46        update Environment variables
47        in env with
48        eval(instruction parameters)
49      end
50    end
51    concrete = { inst  $\in$  (writes  $\setminus$  visit)
52      | depInsts(env, inst) =  $\emptyset$  };
53  end
54 end function
```

tion dependent, which is generated from the step of constructing stack dependency graph. After having some concrete addresses of write instructions, we start the CFG traversing with deep first search (DFS) from the addresses.

Algorithm 7: dfsCFG

```
1 import AnalysisEnvironment as env
2 /* do CFG dfs for building dependency */
3 function dfsCFG(addrDg, wInst, block, writes,
4   reads, visit)
5   if block  $\in$  visit then
6     return;
7   else
8     visit.add(block);
9     rwInsts = block.insts  $\cap$  (writes  $\cup$  reads);
10    if wInst  $\in$  block then
11      rwInsts = rwInsts  $\setminus$ 
12        { inst  $\in$  block.insts | inst.pc < wInst.pc };
13    for inst  $\in$  rwInsts do
14      /* The or part check the probability
15        to re-write the same address */
16      if inst.name = wInst.name
17      and (addrOverlap(env, wInst, inst)
18        or env.addrDepInsts[inst]  $\neq \emptyset$ ) then
19        visit.remove(block);
20        return;
21      if inst  $\in$  reads then
22        deps = env.addrDepInsts[inst]  $\cup$ 
23          env.offsetDepInsts[inst];
24        if deps  $\neq \emptyset$  and
25        deps  $\setminus$  env.ealed =  $\emptyset$  then
26          update Environment variables in
27          env with eval(instruction
28            parameters,  $\emptyset$ )
29        if addrOverlap(env, wInst, inst) then
30          env.ealed.add(inst);
31          addrDg.addEdge(wInst, inst);
32          env.conVals[inst].update(
33            env.conVals[wInst]);
34    end
35    for nextBlock  $\in$  block.outBlock do
36      dfsCFG(wInst, nextBlock,
37        writes, reads, visit);
38    end
39    visit.remove(block);
40  end function
```

In the DFS, we must check if we revisit the block. if did not visit, continue the DFS, or return. The other detail we need to notice is that, we need to remove the beginning instruction's previous instruction when we are exploring first block. For all the read/write instruction in the blocks, we check the address if overlap with the address of original write

address. If the overlapping occurred on write instruction, we need check if there exists the probability to rewrite the beginning address. If the probability exists, the DFS will terminate. If the overlapping occurred on read instruction, the concrete values of beginning instruction will be added to it. And the address dependency between the two instruction will be constructed.

4. EXPERIMENT

5. CONCLUSION

6. REFERENCES

- [1] Gavin Wood, "Ethereum: A secure decentralised generalised transaction ledger eip-150 revision (759dccc - 2017-08-07)," 2017, Accessed: 2018-01-03.
- [2] Nick Szabo, "Advances in distributed security," 2003, Accessed: 2016-04-31.