

## Descriptions for files and folders found in this repo

### Data

- binutils\_compiled folder: compiled binaries for binutils
  - SW: binutils v2.26.1
  - ARCH: ARM, MIPS, x86
  - OP: O0-O3
- final\_data folder: all json files containing info on ACFGs. The files in this folder are used for input into Gemini. Info contains namely the function name ("f\_name"), successor(s) for every node with the index in the overarching list corresponding to the index of a node in the ACFG ("succs"), number of nodes in ACFG ("n\_num"), and the attributes for each node in the ACFG ("features")
  - SW: OpenSSL v1.0.1u, v1.0.1f
  - ARCH: ARM, MIPS, x86
  - OP: O0-O3
  - GCC v5.4
- openssl\_compiled folder: all compiled binaries corresponding to the specifications in final\_data. This is the beginning point where the binaries in this folder are input into the ACFG extractor and the output are supposed to be the files in the final\_data folder
- openssl\_decomp folder: decompiled from binaries in openssl\_compiled folder. For input for LLM approaches
- partial\_test\_fn\_100.txt: names of the functions randomly sampled to form the reduced test dataset
- decomp\_test.json: to be used together with partial\_test\_fn\_100.txt to get all files needed to be used to obtain code explanations using gpt-3.5-turbo (see get\_summaries.py)
- summaries.csv: output from get\_summaries.py
- ssl\_data.csv: all files to be used to get retrieval results for both Gemini and LLM
- for\_llm.csv: all files to be used to get AUC for both Gemini and LLM (random pairwise, 50-50 positive and negative samples)

### Code (LLM approaches)

- get\_summaries.py: obtain code explanations using gpt-3.5-turbo with pseudo C functions as input. All code explanations stored in summaries.csv
- get\_similarities\_sbert.py: get sentence embeddings from code explanations using SBERT
- get\_sim\_codebertscore.py: obtain F1 score from pseudo C functions using CodeBERTScore

### Code (Gemini)

- Refer to <https://github.com/xiaojunxu/dnn-binary-code-similarity> for official code to implement Gemini
- Refer to <https://github.com/qian-feng/Gencoding> for unofficial code to implement ACFG extractor

- See section below for a detailed explanation on how the ACFG extraction works and how to modify the code to run it correctly.

## ACFG Extraction Code (Gencoding/qian-feng, as per Genius)

- Note that the procedure for ACFG extraction is the same for Genius and Gemini
- Requirements
  - IDA Pro 7.1
  - Python 2.7.18
  - Networkx 2.2
- Only scripts in the raw-feature-extractor folder are relevant for this task
- How to run the extraction code
  - C:\path\to\ida64.exe  
-SC:\path\to\Gencoding-master\raw-feature-extractor\preprocessing\_ida.py -c C:\path\to\elf
    - This is the command to generate the ida file. When this is run, IDA Pro will open up and starts to analyse the input file
    - -S switch is to execute a script file when the database is open
    - -c switch is to disassemble a new file and to delete the old database, if any
    - <https://hex-rays.com/products/ida/support/idadoc/417.shtml> -> documentation for IDA Pro command line switches
  - raw\_graphs.py
    - Line 3 sys.path.insert(0, 'C:/Python27/Lib/site-packages')
      - This is for IDA Pro to detect the networkx package installed for Python 2.7.x
      - Change to your Python 2.7.x path
  - preprocessing\_ida.py
    - Line 21 Set 'path' to the path where you want to output the ida files to
- What was the ACFG extraction process like?
  - Reference to this Chinese blog: <https://www.cnblogs.com/lqerio/p/15572511.html#1.3%E8%BE%93%E5%87%BAACFG>
  - preprocessing\_ida.py is the main programme to extract the ACFG
  - From preprocessing\_ida.py, line 22 cfigs = get\_func\_cfigs\_c(FirstSeg()) is the main function which does (1) retrieval of CFGs, and (2) extraction of attributes which will make up the ACFG

```

13  if __name__ == '__main__':
14
15      args = parse_command()
16      path = args.path
17      analysis_flags = idc.GetShortPrm(idc.INF_START_AF)
18      analysis_flags &= ~idc.AF_IMMOFF
19      # turn off "automatically make offset" heuristic
20      idc.SetShortPrm(idc.INF_START_AF, analysis_flags)
21      idaapi.autoWait()
22      cfigs = get_func_cfigs_c(FirstSeg())
23      binary_name = idc.GetInputFile() + '.ida'
24      fullpath = os.path.join(path, binary_name)
25      pickle.dump(cfigs, open(fullpath, 'w'))
26      print binary_name
27      idc.Exit(0)

```

get\_func\_cfigs\_c found is the first function called in preprocessing\_ida.py

```

110 def get_func_cfigs_c(ea):
111     binary_name = idc.GetInputFile()
112     raw_cfigs = raw_graphs(binary_name)
113     externs_eas, ea_externs = processpltSegs()
114     i = 0
115     for funcea in Functions(SegStart(ea)):
116         funcname = get_unified_funcname(funcea)
117         func = get_func(funcea)
118         print i
119         i += 1
120         icfg = cfig.getCfg(func, externs_eas, ea_externs)
121         func_f = get_discoverRe_feature(func, icfg[0])
122         raw_g = raw_graph(funcname, icfg, func_f)
123         raw_cfigs.append(raw_g)
124
125     return raw_cfigs

```

get\_func\_cfigs\_c as seen in func.py

- As seen in the screenshot above, the function first gets the input binary as specified in the command line, then instantiates an raw\_graphs object with the binary\_name.
- The function processpltSegs() first gets the list of all segments (line 224) and for each segment (line 225), it gets the linear address of the start of the segment (226), then uses this address to get the name of the segment (line 227). If the segment name contains '.plt', 'extern' or '.MIPS.stubs', the start and end addresses of that segment is used (lines 228-230) to get the function name and the current address of that function. The 2 dictionaries, funcdata and datafunc, are used to store the function name and address information, where for funcdata the key is the function name while the value is the function address whereas it is the opposite for datafunc (as seen in the reversed dictionary names), where the key is the function address and the value is the function name. Hence, processpltSegs() returns a dictionary of function names and their corresponding addresses, and a dictionary of function addresses and their corresponding names (the dictionaries are the reverse of each other).

```

221 def processpltSegs():
222     funcdata = {}
223     datafunc = {}
224     for n in xrange(idaapi.get_segm_qty()):
225         seg = idaapi.getnseg(n)
226         ea = seg.startEA
227         segname = SegName(ea)
228         if segname in ['.plt', 'extern', '.MIPS.stubs']:
229             start = seg.startEA
230             end = seg.endEA
231             cur = start
232             while cur < end:
233                 name = get_unified_funcname(cur)
234                 funcdata[name] = hex(cur)
235                 datafunc[cur] = name
236                 cur = NextHead(cur)
237     return funcdata, datafunc

```

processpltSegs() function in func.py

- After getting these 2 dictionaries from processpltSegs(), every function in the input binary is iterated through from the start of the first segment (line 115). For each function, the function name 'funcname' is retrieved from get\_unified\_funcname (line 116) and the pointer to the function 'func' is retrieved from get\_func (line 117).
- Next, the CFG is obtained from the getCfg function which can be found in cfg\_constructor.py (line 120). The pointer to the function 'func', and the 2 dictionaries obtained from processpltSegs() are passed in as arguments to getCfg. The CFG is a networkx DiGraph (line 18). For every block in control blocks, unvisited nodes are visited and are added to the CFG. At the end of getCfg, the attributingRe function is called which adds 9 attributes to every node in the newly generated CFG, namely numIns, numCalls, numLIs, numAs, numNc, consts, strings, externs, and numTIs. To note that this isn't the final list of attributes to be used for the ACFG.
- After the CFG is obtained, the get\_discoverRe\_feature function from discovRe.py is called (line 121), returning a list of 11 attributes including FunctionCalls, LogicInstr, Transfer, Locals, BB, Edges, Incoming, Instrs, between, strings, and consts. Note that this is where the 'betweenness' attribute mentioned in the paper comes from. This list of features stored in the variable 'func\_f' in the main get\_func\_cfgs\_c function is then passed as an argument into the raw\_graph constructor along with the function name and the CFG of that function to generate the raw ACFG. (cont. below after the screenshots)

```

15 def getCfg(func, externs_eas, ea_externs):
16     func_start = func.startEA
17     func_end = func.endEA
18     cfg = nx.DiGraph()
19     control_blocks, main_blocks = obtain_block_sequence(func)
20     i = 0
21     visited = {}
22     start_node = None
23     for bl in control_blocks:
24         start = control_blocks[bl][0]
25         end = control_blocks[bl][1]
26         src_node = (start, end)
27         if src_node not in visited:
28             src_id = len(cfg)
29             visited[src_node] = src_id
30             cfg.add_node(src_id)
31             cfg.node[src_id]['label'] = src_node
32         else:
33             src_id = visited[src_node]
34
35         #if end in seq_blocks and GetMnem(PrevHead(end)) != 'jmp':
36         if start == func_start:
37             cfg.node[src_id]['c'] = "start"
38             start_node = src_id
39         if end == func_end:
40             cfg.node[src_id]['c'] = "end"
41         #print control_ea, 1
42         refs = CodeRefsTo(start, 0)
43         for ref in refs:
44             if ref in control_blocks:
45                 dst_node = control_blocks[ref]
46                 if dst_node not in visited:
47                     visited[dst_node] = len(cfg)
48                     dst_id = visited[dst_node]
49                     cfg.add_edge(dst_id, src_id)
50                     cfg.node[dst_id]['label'] = dst_node
51                 #print control_ea, 1
52                 refs = CodeRefsTo(start, 1)
53                 for ref in refs:
54                     if ref in control_blocks:
55                         dst_node = control_blocks[ref]
56                         if dst_node not in visited:
57                             visited[dst_node] = len(cfg)
58                             dst_id = visited[dst_node]
59                             cfg.add_edge(dst_id, src_id)
60                             cfg.node[dst_id]['label'] = dst_node
61                 #print "attributing"
62                 attributingRe(cfg, externs_eas, ea_externs)
63                 # removing deadnodes
64                 #old_cfg = copy.deepcopy(cfg)
65                 #transform(cfg)
66                 return cfg, 0

```

getCfg(func, externs\_eas, ea\_externs) from cfg\_constructor.py

```

139 def attributingRe(cfg, externs_eas, ea_externs):
140     for node_id in cfg:
141         bl = cfg.node[node_id]['label']
142         numIns = calInsts(bl)
143         cfg.node[node_id]['numIns'] = numIns
144         numCalls = calCalls(bl)
145         cfg.node[node_id]['numCalls'] = numCalls
146         numLIs = calLogicInstructions(bl)
147         cfg.node[node_id]['numLIs'] = numLIs
148         numAs = calArithmeticIns(bl)
149         cfg.node[node_id]['numAs'] = numAs
150         strings, consts = getBBconsts(bl)
151         cfg.node[node_id]['numNc'] = len(strings) + len(consts)
152         cfg.node[node_id]['consts'] = consts
153         cfg.node[node_id]['strings'] = strings
154         externs = retrieveExterns(bl, ea_externs)
155         cfg.node[node_id]['externs'] = externs
156         numTIs = calTransferIns(bl)
157         cfg.node[node_id]['numTIs'] = numTIs

```

attributingRe function from cfg\_constructor.py

```

41 def get_discoverRe_feature(func, icfg):
42     start = func.startEA
43     end = func.endEA
44     features = []
45     FunctionCalls = getFuncCalls(func)
46     #1
47     features.append(FunctionCalls)
48     LogicInstr = getLogicInsts(func)
49     #2
50     features.append(LogicInstr)
51     Transfer = getTransferInsts(func)
52     #3
53     features.append(Transfer)
54     Locals = getLocalVariables(func)
55     #4
56     features.append(Locals)
57     BB = getBasicBlocks(func)
58     #5
59     features.append(BB)
60     Edges = len(icfg.edges())
61     #6
62     features.append(Edges)
63     Incoming = getIncommingCalls(func)
64     #7
65     features.append(Incoming)
66     #8
67     Instrs = getIntrs(func)
68     features.append(Instrs)
69     between = retrieveGP(icfg)
70     #9
71     features.append(between)
72
73     strings, consts = getfunc_consts(func)
74     features.append(strings)
75     features.append(consts)
76     return features

```

get\_discoverRe\_feature function from discovRe.py

- During the creation of the ACFG, the raw\_graph object is instantiated with funcname as 'funcname' (name of function), the CFG g as 'old\_g' and the list of features obtained from get\_discoverRe\_feature as 'fun\_features'. **Hence to note that 'fun\_features' isn't the list of attributes we're looking for in the ACFG.** The attributing() function is called, which then creates the ACFG 'g' and the associated features by calling retrieveVec. A list of 8 attributes is returned from retrieveVec, namely consts, strings, offs, numAs, numCalls, numIns, numLIs, and numTIs. **To reiterate, the ACFG is 'g' and not 'old\_g', while the attributes of the ACFG are found in 'g'** (self.g.node[node\_id]['v'] gives the list of attributes) and were generated from the retrieveVec function, all of which can be found in raw\_graphs.py.

```

12 class raw_graph:
13     def __init__(self, funcname, g, func_f):
14         self.funcname = funcname
15         self.old_g = g[0]
16         self.g = nx.DiGraph()
17         self.entry = g[1]
18         self.fun_features = func_f
19         self.attributing()

```

```

24     def attributing(self):
25         self.obtainOffsprings(self.old_g)
26         for node in self.old_g:
27             fvector = self.retrieveVec(node, self.old_g)
28             self.g.add_node(node)
29             self.g.node[node]['v'] = fvector
30
31         for edge in self.old_g.edges():
32             node1 = edge[0]
33             node2 = edge[1]
34             self.g.add_edge(node1, node2)
35
36     def obtainOffsprings(self, g):
37         nodes = g.nodes()
38         for node in nodes:
39             offsprings = {}
40             self.getOffsprings(g, node, offsprings)
41             g.node[node]['offs'] = len(offsprings)
42         return g
43
44     def getOffsprings(self, g, node, offsprings):
45         node_offs = 0
46         sucs = g.successors(node)
47         for suc in sucs:
48             if suc not in offsprings:
49                 offsprings[suc] = 1
50                 self.getOffsprings(g, suc, offsprings)

```

```

52     def retrieveVec(self, id_, g):
53         feature_vec = []
54         #numC0
55         numc = g.node[id_]['consts']
56         feature_vec.append(numc)
57         #nums1
58         nums = g.node[id_]['strings']
59         feature_vec.append(nums)
60         #offsprings2
61         offs = g.node[id_]['offs']
62         feature_vec.append(offs)
63         #numAs3
64         numAs = g.node[id_]['numAs']
65         feature_vec.append(numAs)
66         # of calls4
67         calls = g.node[id_]['numCalls']
68         feature_vec.append(calls)
69         # of insts5
70         insts = g.node[id_]['numIns']
71         feature_vec.append(insts)
72         # of LIIs6
73         insts = g.node[id_]['numLIIs']
74         feature_vec.append(insts)
75         # of TIIs7
76         insts = g.node[id_]['numTIIs']
77         feature_vec.append(insts)
78         return feature_vec

```



- Eg of the features list: [[8L], ['L'], 15, 0, 0, 7, 0, 0] corresponding to [**consts**, **strings**, offs, numAs, numCalls, numIns, numLIs, numTIs]
  - #1 **consts (not numeric)**: list of long integers → get the length of the list as the numeric value
  - #2 **strings (not numeric)**: list of strings → get the length of the list as the numeric value
  - #3 offs: # offspring
  - #4 numAs: # arithmetic instructions
  - #5 numCalls: # calls
  - #6 numIns: # instructions
  - #7 numLIs: # logic instructions
  - #8 numTIs: # transfer instructions
  - \*\*\* missing betweenness?? → found in 'fun\_features' instead of 'v' \*\*\*

Type	Attribute name
Block-level attributes	String Constants
	Numeric Constants
	No. of Transfer Instructions
	No. of Calls
	No. of Instructions
Inter-block attributes	No. of Arithmetic Instructions
	No. of offspring
	Betweenness

**Table 1: Basic-block attributes**