

DynaPyt:

A Dynamic Analysis Framework for Python

Aryaz Eghbali, Michael Pradel
Software Lab, University of Stuttgart

ESEC/FSE 2022



European Research Council
Established by the European Commission



<https://software-lab.org/>

Motivation



cuda runtime error(2): out of memory

Security:

- Dynamic taint analysis

Performance:

- Memory leak
- Slow operations

Testing:

- Branch coverage

The screenshot shows the PyTorch documentation website. The page title is "FREQUENTLY ASKED QUESTIONS". The main heading is "My model reports 'cuda runtime error(2): out of memory'". The text explains that this error occurs when running out of memory on a GPU and lists common causes and fixes. A code block shows a training loop where `total_loss` is accumulated. Below the code, it notes that `total_loss` accumulates history and suggests using `total_loss += float(loss)` instead.

```
total_loss = 0
for i in range(10000):
    optimizer.zero_grad()
    output = model(input)
    loss = criterion(output)
    loss.backward()
    optimizer.step()
    total_loss += loss
```

Options

Ad-hoc

- Implement an instrumenter with LibCST
- Nalin has 200+ lines of code for assignment tracking

```
55 def leave_SimpleStatementLine(self, node, updated_node):
56     body = []
57     for node_line in node.body:
58         if isinstance(node_line, cst.Assign) or isinstance(node_line, cst.AugAssign):
59             line_number = str(self.get_metadata(node=node, key=cst.metadata.PositionProvider).start.line)
60             all_targets = []
61             # TODO: We currently only consider assignments/augassigns of type a = b and ignore a.b.c = m or b[c] = m.
62             if isinstance(node_line, cst.Assign):
63                 # Go through all targets. Only an Assignment can have multiple targets. Eg. a=b=23 or a,b,c = 2,3,4
64                 for assign_target in node_line.targets:
65                     target_single_var = matchers.AssignTarget(target=matchers.Name())
66                     if matchers.matches(assign_target, target_single_var):
67                         all_targets.append(assign_target.target.value)
68                     # Track values of type a,b,c = 1,2,3
69                     target_tuple = matchers.AssignTarget(target=matchers.Tuple())
70                     if matchers.matches(assign_target, target_tuple):
71                         for elem in assign_target.children:
72                             for v in elem.children:
73                                 if matchers.matches(v, matchers.Element(value=matchers.Name())):
74                                     all_targets.append(v.value.value)
75
76             if not len(all_targets):
77                 return updated_node
78         elif isinstance(node_line, cst.AugAssign):
79             target_single_var = matchers.AugAssign(target=matchers.Name())
80             if matchers.matches(node_line, target_single_var):
81                 all_targets.append(node_line.target.value)
82             else:
83                 return updated_node
84     # One call for each target. Eg a = b = 23. A call each for 'a' & 'b'
85     call_expr_nodes = []
86     for target_var in all_targets:
87         # dc = cst.Dict[
88             #     cst.DictElement(cst.Name('variable_name'), cst.Name(target_var)),
89             #     cst.DictElement(cst.Name('line_number'), cst.Integer(line_number)),
90             #     cst.DictElement(cst.Name('value'), node_line.value)
91         # ]
```

Options

sys.settrace

- 70+ lines of code to read the stack properly
- Need low level bytecode and stack operations

```
29 PTR_SIZE = sizeof(PY_OBJECT)
30 F_VALUESTACK_OFFSET = sizeof(Frame) - 2 * PTR_SIZE
31 F_STACKTOP_OFFSET = sizeof(Frame) - PTR_SIZE
32
33
34 class OpStack(Sequence[Any]):
35     __slots__ = ("_frame", "_len")
36
37     def __init__(self, frame):
38         self._frame = Frame.from_address(id(frame))
39         stack_start_addr = c_ssize_t.from_address(id(frame) + F_VALUESTACK_OFFSET).value
40         stack_top_addr = c_ssize_t.from_address(id(frame) + F_STACKTOP_OFFSET).value
41         self._len = (stack_top_addr - stack_start_addr) // PTR_SIZE
42
43     def __repr__(self) -> str:
44         if not self:
45             return "<OpStack> (empty)"
46         return "<OpStack ({})>{}- {}{}\n".format(
47             len(self),
48             "\n- ".join(repr(o) for o in reversed(self)),
49         )
50
51     def __len__(self):
52         return self._len
53
54     def _preproc_slice(self, idx: Optional[int], default: int) -> int:
55         if idx is None:
56             return default
57         if idx < -self._len or idx >= self._len:
58             raise IndexError(idx)
59         if idx < 0:
60             return idx + self._len
61         return idx
62
63     def __getitem__(self, item: Union[int, slice]) -> Any:
64         if isinstance(item, int):
65             if item < -self._len or item >= self._len:
66                 raise IndexError(item)
67             if item < 0:
68                 return self._frame.f_stacktop[item]
69             return self._frame.f_valuystack[item]
70
71         if isinstance(item, slice):
72             item = slice(
73                 self._preproc_slice(item.start, 0),
74                 self._preproc_slice(item.stop, self._len),
75                 item.step
76             )
77             return self._frame.f_valuystack[item]
78
79         raise TypeError(item)
```

```
1 import sys
2 import opcode
3 from get_stack import OpStack
4
5 branches = dict()
6
7 def show_trace(frame, event, arg):
8     global branches
9     frame.f_trace_opcodes = True
10    code = frame.f_code
11    offset = frame.f_lasti
12
13    if opcode.opname[code.co_code[offset]] in ['POP_JUMP_IF_FALSE', 'POP_JUMP_IF_TRUE']:
14        st = OpStack(frame)
15        branches[(frame.f_lineno, st[0])] = branches.get((frame.f_lineno, st[0]), 0) + 1
16    return show_trace
17
18 sys.settrace(show_trace)
```

Options

DynaPyt (this work)

- Just 2 hooks implemented
- At the exact abstraction level of the analysis

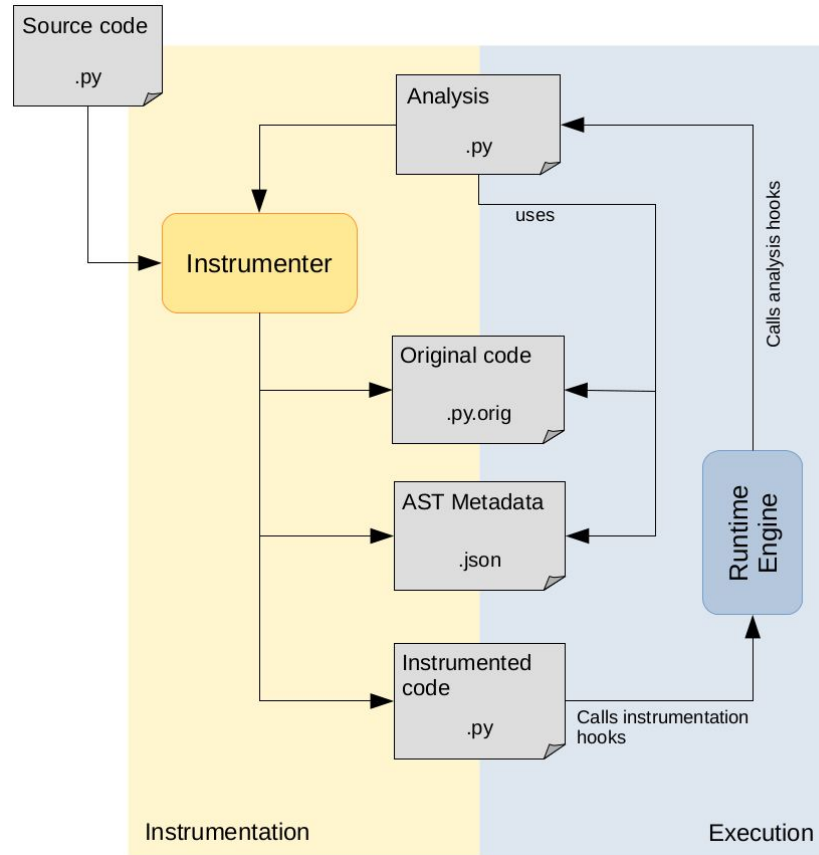
```
1 from typing import Optional
2 from .BaseAnalysis import BaseAnalysis
3
4 class BranchCoverage(BaseAnalysis):
5     def __init__(self):
6         self.branches = dict()
7
8     def enter_control_flow(self, dyn_ast: str, iid: int, cond_value: bool) -> Optional[bool]:
9         self.branches[(iid, bool(cond_value))] = self.branches.get((iid, bool(cond_value)), 0) + 1
10
11     def end_execution(self):
12         for k, v in self.branches.items():
13             print(f'Branch {k[0]} taken with condition {k[1]}, {v} time{" if v == 1 else "s"}')
```

Choose wisely!

	Engineering Effort	Abstraction Level	Extra Runtime Overhead
Ad-hoc	High	Matching the analysis	Low
sys.settrace	Medium	Different from the analysis	High
DynaPyt	Low	Matching the analysis	Low

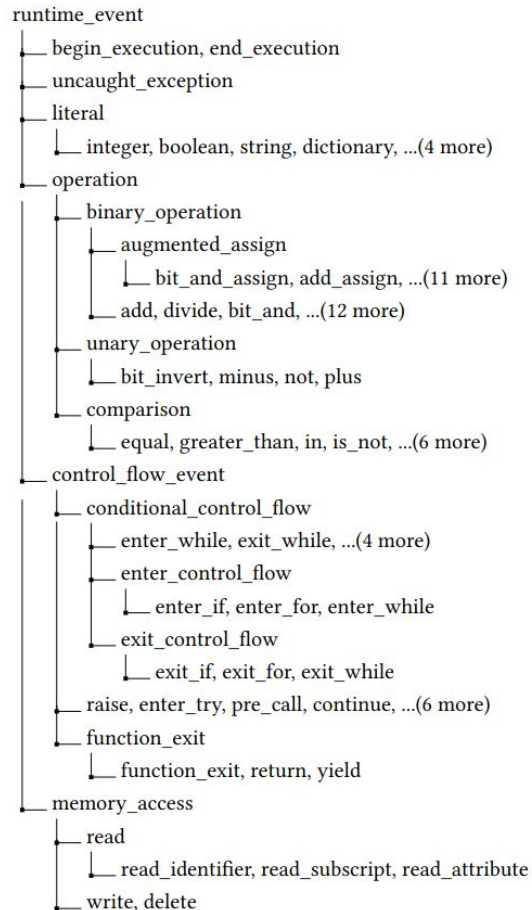
DynaPyt

- Instrumentation
- Runtime engine
- Analysis code



Features

- 97 available hooks
 - Hierarchy: various levels of abstraction
 - Any combination
- Pay-per-use
 - Only the used hooks get instrumented
→ overhead only for used hooks
- Modify execution
 - Runtime values → in e.g. concolic testing

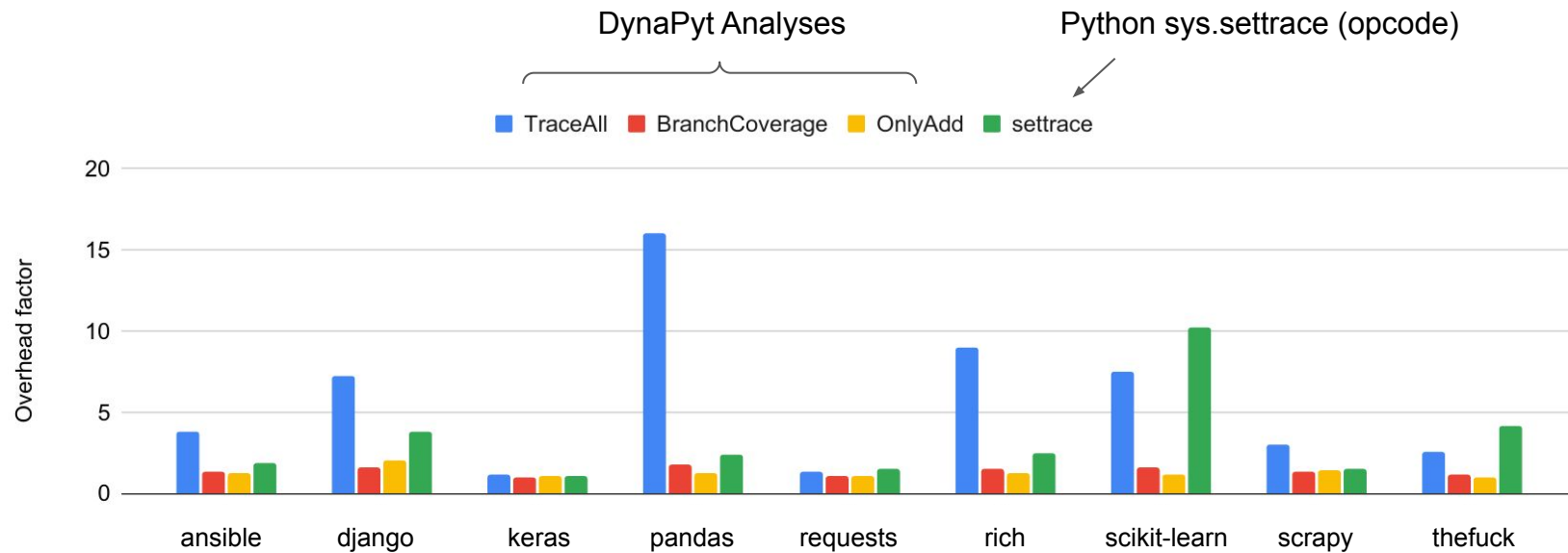


Instrument Time

TraceAll analysis:
most expensive
instrumentation

	Repository	Instrument time (mm:ss)	# of files	Lines of code
1	ansible/ansible	6:59	2,188	176,173
2	django/django	14:07	3,603	318,602
3	keras-team/keras	5:41	678	155,407
4	pandas-dev/pandas	12:32	2,727	358,195
5	psf/requests	0:16	54	6,370
6	Textualize/rich	0:57	178	24,362
7	scikit-learn/scikit-learn	6:52	1,419	180,185
8	scrapy/scrapy	1:49	505	37,181
9	nvbn/thewfuck	1:21	620	12,070

Runtime Overhead



Usage

Install:

```
pip install dynapyt
```

Instrument code:

```
python -m dynapyt.run_instrumentation --directory  
<dir> --analysis MLMemoryAnalysis
```

Run analysis:

```
python -m dynapyt.run_analysis --entry <main.py>  
--analysis MLMemoryAnalysis
```

```
1 # Inspired by https://pytorch.org/docs/stable/notes/faq.html#my-model-reports-cuda-runtime-error-2-out-of-memory  
2 from collections import defaultdict  
3 from .BaseAnalysis import BaseAnalysis  
4  
5 class MLMemoryAnalysis(BaseAnalysis):  
6     def __init__(self) -> None:  
7         super().__init__()  
8         self.in_ctrl_flow = []  
9         self.threshold = 3  
10        self.memory_leak = defaultdict(lambda: 0)  
11        self.last_opr = None  
12  
13        def enter_control_flow(self, dyn_ast, iid, condition):  
14            self.last_opr = None  
15            if (len(self.in_ctrl_flow) > 0) and (self.in_ctrl_flow[-1] != iid):  
16                self.in_ctrl_flow.append(iid)  
17  
18        def exit_control_flow(self, dyn_ast, iid):  
19            self.last_opr = None  
20            self.in_ctrl_flow.pop()  
21  
22        def binary_operation(self, dyn_ast, iid, opr, left, right, res):  
23            if (len(self.in_ctrl_flow) > 0) and right.requires_grad:  
24                self.last_opr = iid  
25            else:  
26                self.last_opr = None  
27  
28        def write(self, dyn_ast, iid, left, right):  
29            if (len(self.in_ctrl_flow) > 0) and right.requires_grad and (self.last_opr is not None):  
30                cur = (iid, self.in_ctrl_flow[-1])  
31                self.memory_leak[cur] += 1  
32                if self.memory_leak[cur] > 3:  
33                    print('Memory issue detected')  
34                    exit(1)  
35            self.last_opr = None
```

Analysis Simplicity

Name	Description	Analysis hooks	LoC
BranchCoverage	Measures how often each branch gets covered	1	6
CallGraph	Computes a dynamic call graph	1	19
KeyInList	Warns about performance anti-pattern of linearly search through a list	2	10
MLMemory	Warns about memory leak issues in deep learning code	4	29
SimpleTaint	Taint analysis useful to, e.g., detect SQL injections	7	53
AllEvents	Implements the runtime_event analysis hook to trace all events	1	4

Future Work

For DynaPyt:

- Write to attributes, as a multi-write, to a tuple
- `async/await`

With DynaPyt:

- Early detection of ML issues
- Creating datasets of Python executions

Q&A

DynaPyt: Dynamic Analysis Framework for Python

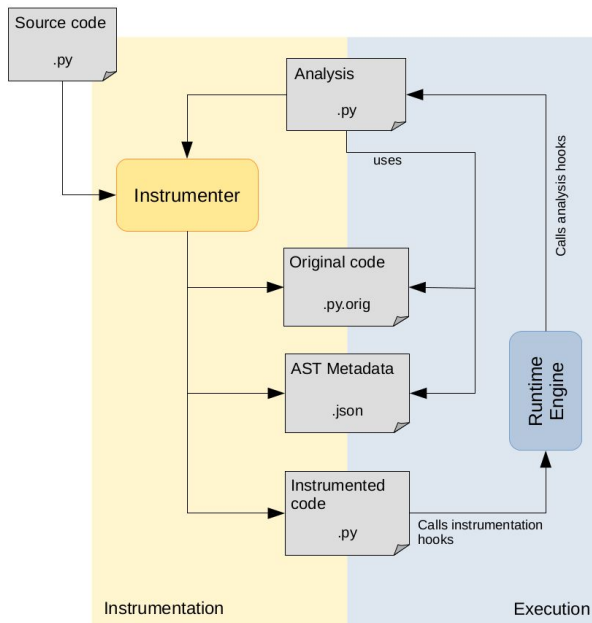
- ★ Ease of analysis implementation
- ★ Low overhead runtime

Install:

```
pip install dynapyt
```

Code & documentation:

```
https://github.com/sola-st/DynaPyt
```



```
runtime_event
├── begin_execution, end_execution
├── uncaught_exception
├── literal
│   ├── integer, boolean, string, dictionary, ...(4 more)
├── operation
│   ├── binary_operation
│   │   ├── augmented_assign
│   │   │   ├── bit_and_assign, add_assign, ...(11 more)
│   │   │   ├── add, divide, bit_and, ...(12 more)
│   ├── unary_operation
│   │   ├── bit_invert, minus, not, plus
├── comparison
│   ├── equal, greater_than, in, is_not, ...(6 more)
├── control_flow_event
│   ├── conditional_control_flow
│   │   ├── enter_while, exit_while, ...(4 more)
│   │   ├── enter_control_flow
│   │   │   ├── enter_if, enter_for, enter_while
│   │   ├── exit_control_flow
│   │   │   ├── exit_if, exit_for, exit_while
│   ├── raise, enter_try, pre_call, continue, ...(6 more)
├── function_exit
│   ├── function_exit, return, yield
├── memory_access
│   ├── read
│   │   ├── read_identifier, read_subscript, read_attribute
│   ├── write, delete
```

Execution Faithfulness

- Preserve original execution
 - All above 97.8% passing tests
 - Part of the difference is due to execution stack accesses