

A Multilanguage Static Analysis of Python/C Programs with Mopsa

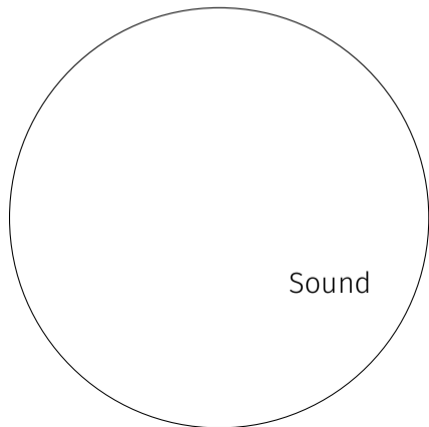
Raphaël Monat, Abdelraouf Ouadjaout, Antoine Miné

Software Languages Lab, VUB
7 July 2022

rmonat.fr

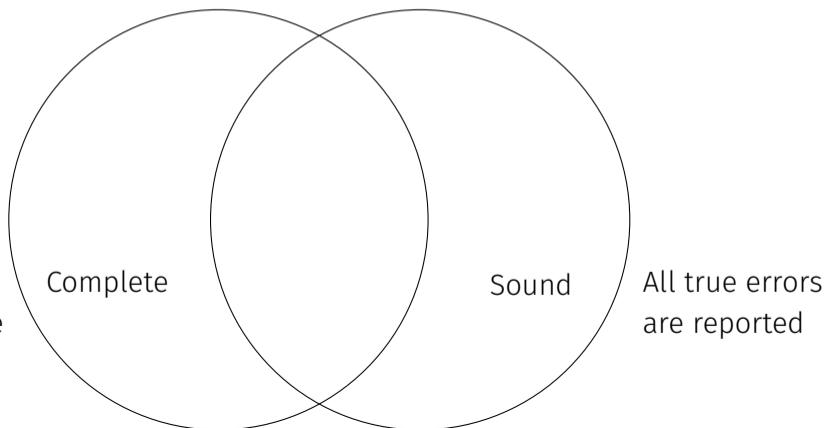


Introduction

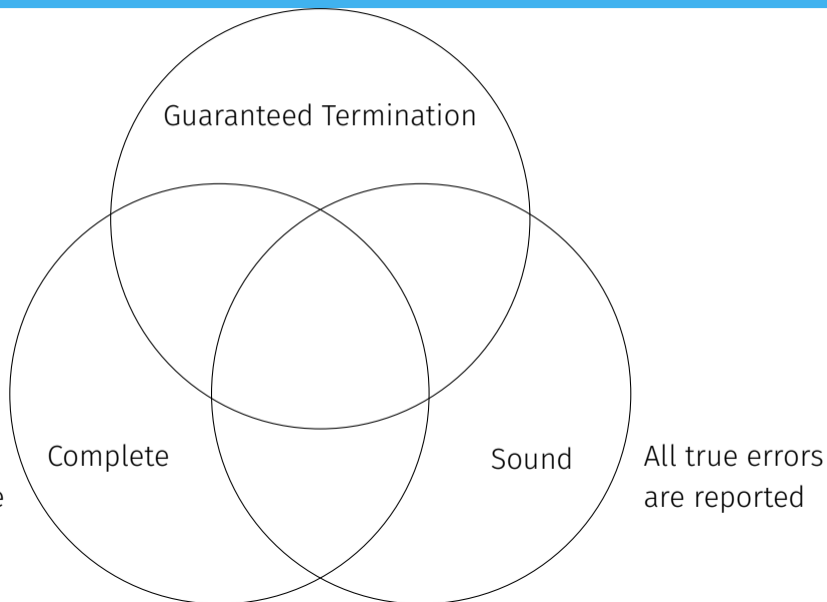


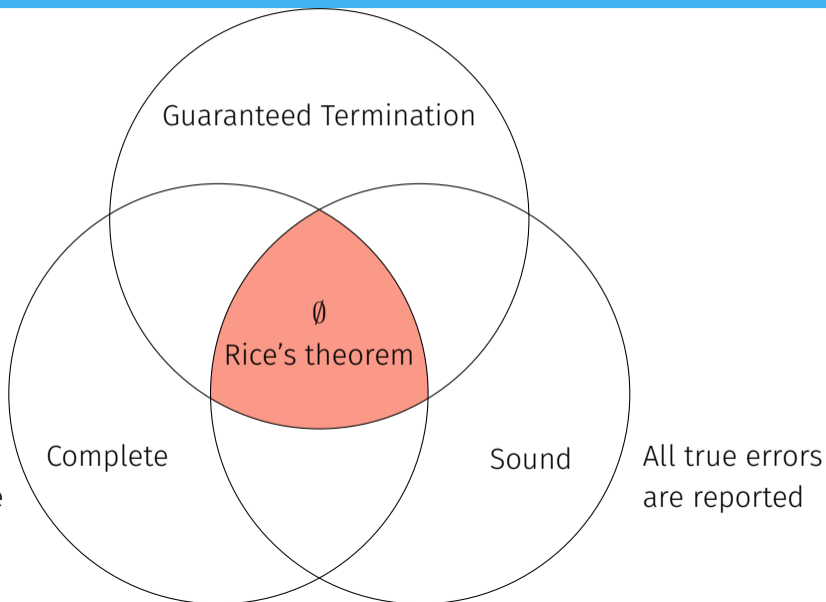
Sound

All true errors
are reported

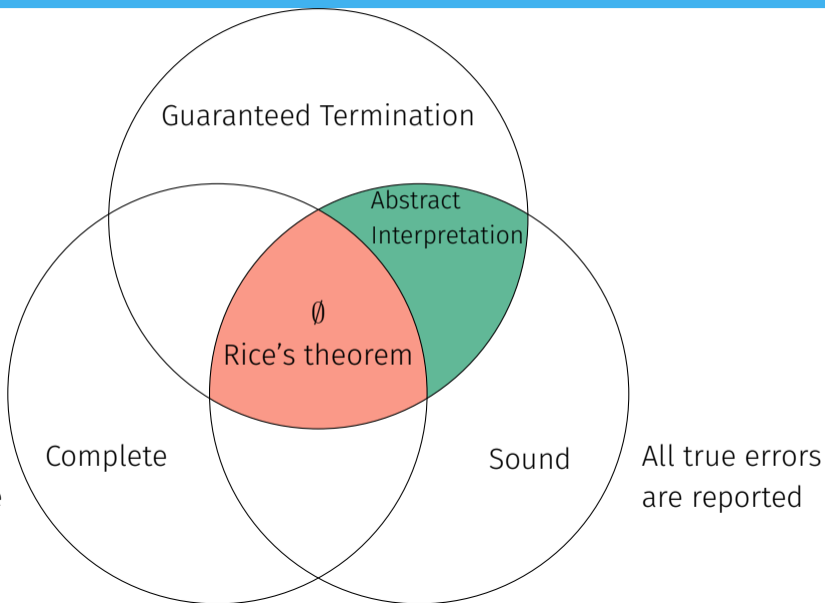


Program verification





Program verification



Conservative static program analysis

average.py

```
1 def average(l):
2     m = 0
3     for i in range(len(l)):
4         m = m + l[i]
5         m = m // (i + 1)
6     return s
7
8 r1 = average([1, 2, 3])
9 r2 = average(['a', 'b', 'c'])
```

TypeError: unsupported operand type(s) for '+': 'int' and 'str'

argslen.c

```
1 #include <string.h>
2
3 int main(int argc, char *argv[]) {
4     int i = 0;
5     for (char **p = argv; *p; p++) {
6         strlen(*p); // valid string
7         i++; // no overflow
8     }
9     return 0;
10 }
```

No alarm

Specifications of the analyzer

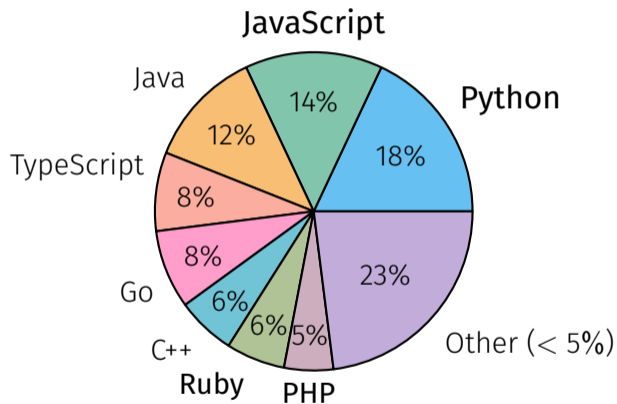
Inference of program properties such as the absence of run-time errors.

Semantic based on a formal modelization of the language.

Automatic no expert knowledge required.

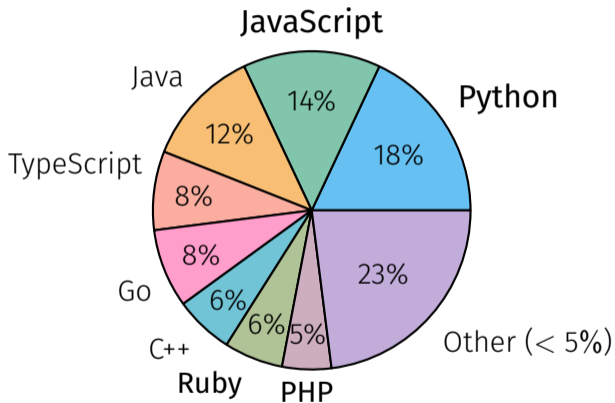
Sound covers all possible executions.

Dynamic programming languages



Most popular languages on GitHub

Dynamic programming languages



Most popular languages on GitHub

Features

- ▶ Object orientation
- ▶ Dynamic typing
- ▶ Dynamic object structure
- ▶ Introspection operators
- ▶ `eval`

One in five of the top 200 Python libraries contains C code

One in five of the top 200 Python libraries contains C code

- ▶ To bring better performance (numpy)

One in five of the top 200 Python libraries contains C code

- ▶ To bring better performance (numpy)
- ▶ To provide library bindings (pygit2)

One in five of the top 200 Python libraries contains C code

- ▶ To bring better performance (numpy)
- ▶ To provide library bindings (pygit2)

Pitfalls

One in five of the top 200 Python libraries contains C code

- ▶ To bring better performance (numpy)
- ▶ To provide library bindings (pygit2)

Pitfalls

- ▶ Different values (\mathbb{Z} vs. `Int32`)

One in five of the top 200 Python libraries contains C code

- ▶ To bring better performance (numpy)
- ▶ To provide library bindings (pygit2)

Pitfalls

- ▶ Different values (\mathbb{Z} vs. `Int32`)
- ▶ Shared memory state

Outline

- 1 Introduction
- 2 A Taste of Python
- 3 Mopsa
- 4 Towards a Multilanguage Semantics
- 5 Implementation & Experimental Evaluation
- 6 Conclusion

A Taste of Python

No standard

- ▶ CPython is the reference

⇒ manual inspection of the source code and handcrafted tests

No standard

- ▶ CPython is the reference
 - ⇒ manual inspection of the source code and handcrafted tests

Operator redefinition

- ▶ Calls, additions, attribute accesses
- ▶ Operators eventually call overloaded `__methods__`

Protected attributes

```
1 class Protected:
2     def __init__(self, priv):
3         self._priv = priv
4     def __getattr__(self, attr):
5         if attr[0] == "_": raise AttributeError("...")
6         return object.__getattr__(self, attr)
7
8 a = Protected(42)
9 a._priv # AttributeError raised
```

Python's specificities (II)

Dual type system

- ▶ Nominal (classes, MRO)

Fspath (from standard library)

```
1 class Path:
2     def __fspath__(self): return 42
3
4 def fspath(p):
5     if isinstance(p, (str, bytes)):
6         return p
7     elif hasattr(p, "__fspath__"):
8         r = p.__fspath__()
9         if isinstance(r, (str, bytes)):
10            return r
11        raise TypeError
12
13 fspath("/dev" if random() else Path())
```

Barrett, Cassels, Haahr, Moon, Playford, and Withington. "A Monotonic Superclass Linearization for Dylan". OOPSLA 1996

Python's specificities (II)

Dual type system

- ▶ Nominal (classes, MRO)
- ▶ Structural (attributes)

Fspath (from standard library)

```
1 class Path:
2     def __fspath__(self): return 42
3
4 def fspath(p):
5     if isinstance(p, (str, bytes)):
6         return p
7     elif hasattr(p, "__fspath__"):
8         r = p.__fspath__()
9         if isinstance(r, (str, bytes)):
10            return r
11        raise TypeError
12
13 fspath("/dev" if random() else Path())
```

Barrett, Cassels, Haahr, Moon, Playford, and Withington. "A Monotonic Superclass Linearization for Dylan". OOPSLA 1996

Python's specificities (II)

Dual type system

- ▶ Nominal (classes, MRO)
- ▶ Structural (attributes)

Exceptions

Exceptions rather than specific values

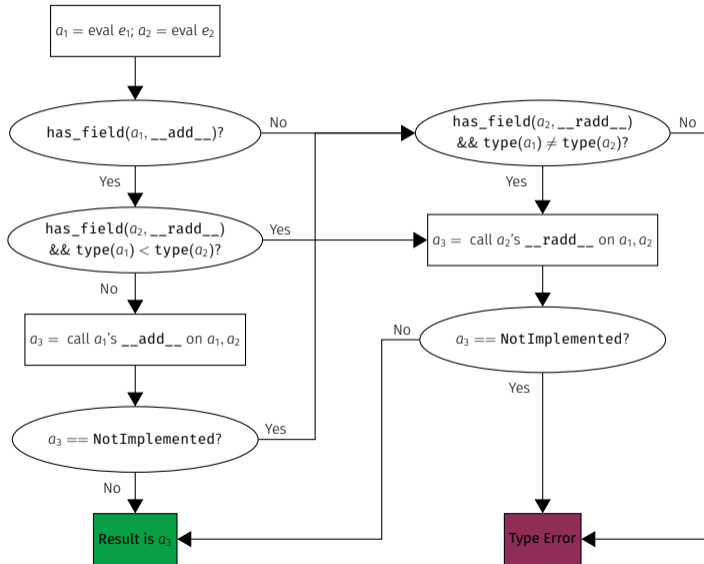
- ▶ `1 + "a" ↪ TypeError`
- ▶ `l[len(l) + 1] ↪ IndexError`

Fspath (from standard library)

```
1 class Path:
2     def __fspath__(self): return 42
3
4 def fspath(p):
5     if isinstance(p, (str, bytes)):
6         return p
7     elif hasattr(p, "__fspath__"):
8         r = p.__fspath__()
9         if isinstance(r, (str, bytes)):
10            return r
11            raise TypeError
12
13 fspath("/dev" if random() else Path())
```

Barrett, Cassels, Haahr, Moon, Playford, and Withington. "A Monotonic Superclass Linearization for Dylan". OOPSLA 1996

Example Semantics – binary operators



Custom infix operators

```
1 class Infix(object):
2     def __init__(self, func): self.func = func
3     def __or__(self, other): return self.func(other)
4     def __ror__(self, other): return Infix(lambda x: self.func(other, x))
5
6 instanceof = Infix(isinstance)
7 b = 5 |instanceof| int
8
9 @Infix
10 def padd(x, y):
11     print(f"{x} + {y} = {x + y}")
12     return x + y
13 c = 2 |padd| 3
```

Credits tomerriliba.com/blog/Infix-Operators/

Mopsa



A program analysis workflow

Averaging numbers

```
1 def average(l):
2     m = 0
3     for i in range(len(l)):
4         m = m + l[i]
5     m = m // (i + 1)
6     return m
7
8 l = [randint(0, 20)
9     for i in range(randint(5, 10))]
10 m = average(l)
```

A program analysis workflow

Averaging numbers

```
1 def average(l):
2     m = 0
3     for i in range(len(l)):
4         m = m + l[i]
5         m = m // (i + 1)
6     return m
7
8 l = [randint(0, 20)
9     for i in range(randint(5, 10))]
10 m = average(l)
```

Proved safe?

- ▶ $m // (i+1)$
- ▶ $l[i]$

Searching for a loop invariant (l. 4)

Environment abstraction

$$m \mapsto @_{\text{int}}^{\#} \quad i \mapsto @_{\text{int}}^{\#}$$

A program analysis workflow

Averaging numbers

```
1 def average(l):
2     m = 0
3     for i in range(len(l)):
4         m = m + l[i]
5     m = m // (i + 1)
6     return m
7
8 l = [randint(0, 20)
9     for i in range(randint(5, 10))]
10 m = average(l)
```

Proved safe?

- ▶ $m // (i+1)$
- ▶ $l[i]$

Searching for a loop invariant (l. 4)

Stateless domains: **list content**,

Environment abstraction

$$m \mapsto @_{\text{int}}^{\#} \quad i \mapsto @_{\text{int}}^{\#} \quad \underline{\text{els}}(l) \mapsto @_{\text{int}}^{\#}$$

A program analysis workflow

Averaging numbers

```
1 def average(l):
2     m = 0
3     for i in range(len(l)):
4         m = m + l[i]
5     m = m // (i + 1)
6     return m
7
8 l = [randint(0, 20)
9     for i in range(randint(5, 10))]
10 m = average(l)
```

Proved safe?

- ▶ $m // (i+1)$
- ▶ $l[i]$

Searching for a loop invariant (l. 4)

Stateless domains: list content,

Environment abstraction

$$m \mapsto @_{\text{int}}^{\#} \quad i \mapsto @_{\text{int}}^{\#} \quad \underline{\text{els}}(l) \mapsto @_{\text{int}}^{\#}$$

Numeric abstraction (intervals)

$$m \in [0, +\infty) \quad \underline{\text{els}}(l) \in [0, 20] \quad i \in [0, +\infty)$$

A program analysis workflow

Averaging numbers

```
1 def average(l):
2     m = 0
3     for i in range(len(l)):
4         m = m + l[i]
5     m = m // (i + 1)
6     return m
7
8 l = [randint(0, 20)
9     for i in range(randint(5, 10))]
10 m = average(l)
```

Proved safe?

- ▶ `m // (i+1)`
- ▶ `l[i]`

Searching for a loop invariant (l. 4)

Stateless domains: list content, **list length**

Environment abstraction

$$m \mapsto @_{\text{int}}^{\#} \quad i \mapsto @_{\text{int}}^{\#} \quad \underline{\text{els}}(l) \mapsto @_{\text{int}}^{\#}$$

Numeric abstraction (intervals)

$$m \in [0, +\infty) \quad \underline{\text{els}}(l) \in [0, 20]$$
$$\underline{\text{len}}(l) \in [5, 10] \quad i \in [0, 10]$$

A program analysis workflow

Averaging numbers

```
1 def average(l):
2     m = 0
3     for i in range(len(l)):
4         m = m + l[i]
5     m = m // (i + 1)
6     return m
7
8 l = [randint(0, 20)
9     for i in range(randint(5, 10))]
10 m = average(l)
```

Proved safe?

- ▶ `m // (i+1)`
- ▶ `l[i]`

Searching for a loop invariant (l. 4)

Stateless domains: list content, list length

Environment abstraction

$$m \mapsto @_{\text{int}}^{\#} \quad i \mapsto @_{\text{int}}^{\#} \quad \underline{\text{els}}(l) \mapsto @_{\text{int}}^{\#}$$

Numeric abstraction (polyhedra)

$$m \in [0, +\infty) \quad \underline{\text{els}}(l) \in [0, 20]$$
$$0 \leq i < \underline{\text{len}}(l) \quad 5 \leq \underline{\text{len}}(l) \leq 10$$

A program analysis workflow

Averaging tasks

```
1 class Task:
2     def __init__(self, weight):
3         if weight < 0: raise ValueError
4         self.weight = weight
5
6     def average(l):
7         m = 0
8         for i in range(len(l)):
9             m = m + l[i].weight
10            m = m // (i + 1)
11        return m
12
13 l = [Task(randint(0, 20))
14      for i in range(randint(5, 10))]
15 m = average(l)
```

Proved safe?

- ▶ $m // (i+1)$
- ▶ $l[i].weight$

Searching for a loop invariant (l. 4)

Stateless domains: list content, list length

Environment abstraction

$$m \mapsto @_{int}^{\#} \quad i \mapsto @_{int}^{\#} \quad \underline{\text{els}(l)} \mapsto @_{Task}^{\#}$$
$$\underline{@_{Task}^{\#} \cdot \text{weight}} \mapsto @_{int}^{\#}$$

Numeric abstraction (polyhedra)

$$m \in [0, +\infty)$$
$$0 \leq i < \underline{\text{len}(l)} \quad 5 \leq \underline{\text{len}(l)} \leq 10$$
$$0 \leq \underline{@_{Task}^{\#} \cdot \text{weight}} \leq 20$$

Attributes abstraction

$$@_{Task}^{\#} \mapsto (\{\text{weight}\}, \emptyset)$$

A program analysis workflow

Averaging tasks

```
1 class Task:
2     def __init__(self, weight):
3         if weight < 0: raise ValueError
4         self.weight = weight
5
6 def average(l):
7     m = 0
8     for i in range(len(l)):
9         m = m + l[i].weight
10        m = m // (i + 1)
11    return m
12
13 l = [Task(randint(0, 10), randint(1, 10))
14      for i in range(randint(5, 10))]
15 m = average(l)
```

Conclusion

- ▶ Different domains depending on the precision
- ▶ Use of auxiliary variables (underlined)

Searching for a loop invariant (l. 4)

Stateless domains: list content, list length

Environment abstraction

$m \mapsto @^{\#} \dots i \dots @^{\#} \dots @^{\#} \text{Task}$

$0 \leq i < \underline{\text{len}}(l) \quad 5 \leq \underline{\text{len}}(l) \leq 10$

$0 \leq \underline{@^{\#}_{\text{Task}} \cdot \text{weight}} \leq 20$

Proved safe?

- ▶ $m // (i+1)$
- ▶ $l[i].\text{weight}$

Attributes abstraction

$@^{\#}_{\text{Task}} \mapsto (\{\text{weight}\}, \emptyset)$



Modular Open Platform for Static Analysis¹

gitlab.com/mopsa/mopsa-analyzer

¹Journault, Miné, Monat, and Ouadjaout. “Combinations of reusable abstract domains for a multilingual static analyzer”. 2019.



Modular Open Platform for Static Analysis¹




gitlab.com/mopsa/mopsa-analyzer




- ▶ One AST to analyze them all
 - 🚩 Multilanguage support
 - 📄 Expressiveness
 - ♻️ Reusability

¹Journault, Miné, Monat, and Ouadjaout. “Combinations of reusable abstract domains for a multilingual static analyzer”. 2019.



Modular Open Platform for Static Analysis¹ gitlab.com/mopsa/mopsa-analyzer

- ▶ One AST to analyze them all
 - ▶  Multilanguage support
 - ▶  Expressiveness
 - ▶  Reusability

- ▶ Unified domain signature
 - ▶  Semantic rewriting
 - ▶  Loose coupling
 - ▶  Observability

¹Journault, Miné, Monat, and Ouadjaout. “Combinations of reusable abstract domains for a multilingual static analyzer”. 2019.



Modular Open Platform for Static Analysis¹

gitlab.com/mopsa/mopsa-analyzer

▶ One AST to analyze them all



Multilanguage support



Expressiveness



Reusability

▶ Unified domain signature



Semantic rewriting



Loose coupling



Observability

▶ DAG of abstract domains



Composition



Cooperation

¹Journault, Miné, Monat, and Ouadjaout. “Combinations of reusable abstract domains for a multilingual static analyzer”. 2019.



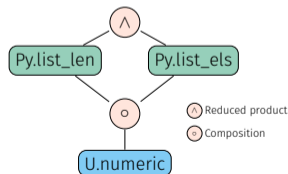
Modular Open Platform for Static Analysis¹

gitlab.com/mopsa/mopsa-analyzer

- ▶ One AST to analyze them all
 - ▶ Multilanguage support
 - ▶ Expressiveness
 - ▶ Reusability

- ▶ Unified domain signature
 - ▶ Semantic rewriting
 - ▶ Loose coupling
 - ▶ Observability

- ▶ DAG of abstract domains
 - ▶ Composition
 - ▶ Cooperation



¹Journault, Miné, Monat, and Ouadjaout. “Combinations of reusable abstract domains for a multilingual static analyzer”. 2019.

Universal.Iterators.Loops

Matches `while(...){...}`

Computes fixpoint using widening

Dynamic, semantic iterators with delegation

```
for(init; cond; incr) body
```

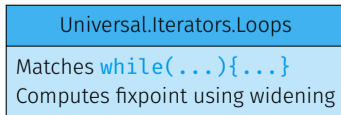
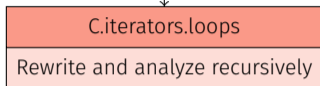
Universal.Iterators.Loops

Matches `while(...){...}`

Computes fixpoint using widening

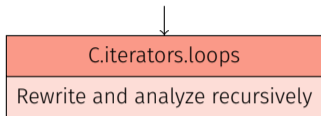
Dynamic, semantic iterators with delegation

`for(init; cond; incr) body`

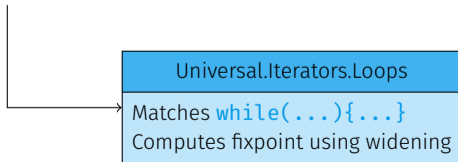


Dynamic, semantic iterators with delegation

```
for(init; cond; incr) body
```



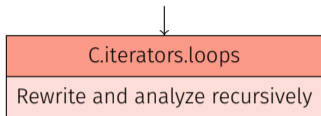
```
init;  
while(cond) {  
  body;  
  incr;  
}  
clean init
```



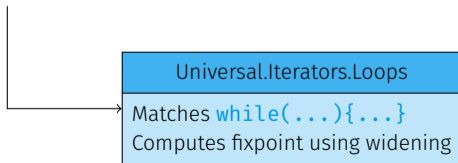
Dynamic, semantic iterators with delegation

```
for(init; cond; incr) body
```

```
for target in iterable: body
```

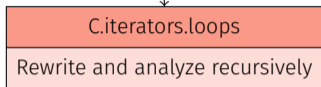


```
init;  
while(cond) {  
    body;  
    incr;  
}  
clean init
```

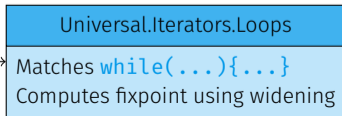


Dynamic, semantic iterators with delegation

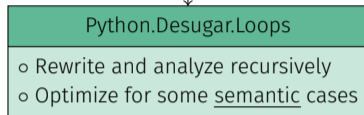
```
for(init; cond; incr) body
```



```
init;  
while(cond) {  
    body;  
    incr;  
}  
clean init
```

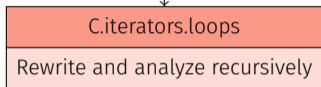


```
for target in iterable: body
```



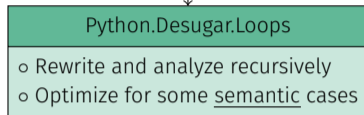
Dynamic, semantic iterators with delegation

`for(init; cond; incr) body`

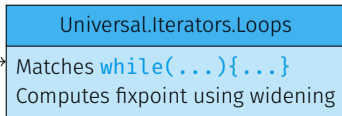


```
init;  
while(cond) {  
    body;  
    incr;  
}  
clean init
```

`for target in iterable: body`



```
it = iter(iterable)  
while(1) {  
    try: target = next(it)  
    except StopIteration: break  
    body  
}  
clean it
```



Towards a Multilanguage Semantics

Multilanguage code – example

counter.c

```
1 typedef struct {
2     PyObject_HEAD;
3     int count;
4 } Counter;
5
6 static PyObject*
7 CounterIncr(Counter *self, PyObject *args)
8 {
9     int i = 1;
10    if(!PyArg_ParseTuple(args, "|i", &i))
11        return NULL;
12
13    self->count += i;
14    Py_RETURN_NONE;
15 }
16
17 static PyObject*
18 CounterGet(Counter *self)
19 {
20     return Py_BuildValue("i", self->count);
21 }
```

count.py

```
22 from counter import Counter
23 from random import randrange
24
25 c = Counter()
26 power = randrange(128)
27 c.incr(2**power-1)
28 c.incr()
29 r = c.get()
```


Multilanguage code – example

counter.c

```
1  typedef struct {
2      PyObject_HEAD;
3      int count;
4  } Counter;
5
6  static PyObject*
7  CounterIncr(Counter *self, PyObject *args)
8  {
9      int i = 1;
10     if(!PyArg_ParseTuple(args, "|i", &i))
11         return NULL;
12
13     self->count += i;
14     Py_RETURN_NONE;
15 }
16
17 static PyObject*
18 CounterGet(Counter *self)
19 {
20     return Py_BuildValue("i", self->count);
21 }
```

count.py

```
22 from counter import Counter
23 from random import randrange
24
25 c = Counter()
26 power = randrange(128)
27 c.incr(2**power-1)
28 c.incr()
29 r = c.get()
```

Multilanguage code – example

counter.c

```
1  typedef struct {
2      PyObject_HEAD;
3      int count;
4  } Counter;
5
6  static PyObject*
7  CounterIncr(Counter *self, PyObject *args)
8  {
9      int i = 1;
10     if(!PyArg_ParseTuple(args, "|i", &i))
11         return NULL;
12
13     self->count += i;
14     Py_RETURN_NONE;
15 }
16
17 static PyObject*
18 CounterGet(Counter *self)
19 {
20     return Py_BuildValue("i", self->count);
21 }
```

count.py

```
22 from counter import Counter
23 from random import randrange
24
25 c = Counter()
26 power = randrange(128)
27 c.incr(2**power-1)
28 c.incr()
29 r = c.get()
```

Multilanguage code – example

counter.c

```
1  typedef struct {
2      PyObject_HEAD;
3      int count;
4  } Counter;
5
6  static PyObject*
7  CounterIncr(Counter *self, PyObject *args)
8  {
9      int i = 1;
10     if(!PyArg_ParseTuple(args, "|i", &i))
11         return NULL;
12
13     self->count += i;
14     Py_RETURN_NONE;
15 }
16
17 static PyObject*
18 CounterGet(Counter *self)
19 {
20     return Py_BuildValue("i", self->count);
21 }
```

count.py

```
22 from counter import Counter
23 from random import randrange
24
25 c = Counter()
26 power = randrange(128)
27 c.incr(2**power-1)
28 c.incr()
29 r = c.get()
```

Multilanguage code – example

counter.c

```
1  typedef struct {
2      PyObject_HEAD;
3      int count;
4  } Counter;
5
6  static PyObject*
7  CounterIncr(Counter *self, PyObject *args)
8  {
9      int i = 1;
10     if(!PyArg_ParseTuple(args, "|i", &i))
11         return NULL;
12
13     self->count += i;
14     Py_RETURN_NONE;
15 }
16
17 static PyObject*
18 CounterGet(Counter *self)
19 {
20     return Py_BuildValue("i", self->count);
21 }
```

count.py

```
22 from counter import Counter
23 from random import randrange
24
25 c = Counter()
26 power = randrange(128)
27 c.incr(2**power-1)
28 c.incr()
29 r = c.get()
```

Multilanguage code – example

counter.c

```
1  typedef struct {
2      PyObject_HEAD;
3      int count;
4  } Counter;
5
6  static PyObject*
7  CounterIncr(Counter *self, PyObject *args)
8  {
9      int i = 1;
10     if(!PyArg_ParseTuple(args, "|i", &i))
11         return NULL;
12
13     self->count += i;
14     Py_RETURN_NONE;
15 }
16
17 static PyObject*
18 CounterGet(Counter *self)
19 {
20     return Py_BuildValue("i", self->count);
21 }
```

count.py

```
22 from counter import Counter
23 from random import randrange
24
25 c = Counter()
26 power = randrange(128)
27 c.incr(2**power-1)
28 c.incr()
29 r = c.get()
```

Multilanguage code – example

counter.c

```
1  typedef struct {
2      PyObject_HEAD;
3      int count;
4  } Counter;
5
6  static PyObject*
7  CounterIncr(Counter *self, PyObject *args)
8  {
9      int i = 1;
10     if(!PyArg_ParseTuple(args, "|i", &i))
11         return NULL;
12
13     self->count += i;
14     Py_RETURN_NONE;
15 }
16
17 static PyObject*
18 CounterGet(Counter *self)
19 {
20     return Py_BuildValue("i", self->count);
21 }
```

count.py

```
22 from counter import Counter
23 from random import randrange
24
25 c = Counter()
26 power = randrange(128)
27 c.incr(2**power-1)
28 c.incr()
29 r = c.get()
```

▶ $\text{power} \leq 30 \Rightarrow r = 2^{\text{power}}$

Multilanguage code – example

counter.c

```
1 typedef struct {
2     PyObject_HEAD;
3     int count;
4 } Counter;
5
6 static PyObject*
7 CounterIncr(Counter *self, PyObject *args)
8 {
9     int i = 1;
10    if(!PyArg_ParseTuple(args, "|i", &i))
11        return NULL;
12
13    self->count += i;
14    Py_RETURN_NONE;
15 }
16
17 static PyObject*
18 CounterGet(Counter *self)
19 {
20     return Py_BuildValue("i", self->count);
21 }
```

count.py

```
22 from counter import Counter
23 from random import randrange
24
25 c = Counter()
26 power = randrange(128)
27 c.incr(2**power-1)
28 c.incr()
29 r = c.get()
```

- ▶ $\text{power} \leq 30 \Rightarrow r = 2^{\text{power}}$
- ▶ $32 \leq \text{power} \leq 64$: OverflowError:
signed integer is greater than maximum
- ▶ $\text{power} \geq 64$: OverflowError:
Python int too large to convert to C long

Multilanguage code – example

counter.c

```
1 typedef struct {
2     PyObject_HEAD;
3     int count;
4 } Counter;
5
6 static PyObject*
7 CounterIncr(Counter *self, PyObject *args)
8 {
9     int i = 1;
10    if(!PyArg_ParseTuple(args, "|i", &i))
11        return NULL;
12
13    self->count += i;
14    Py_RETURN_NONE;
15 }
16
17 static PyObject*
18 CounterGet(Counter *self)
19 {
20     return Py_BuildValue("i", self->count);
21 }
```

count.py

```
22 from counter import Counter
23 from random import randrange
24
25 c = Counter()
26 power = randrange(128)
27 c.incr(2**power-1)
28 c.incr()
29 r = c.get()
```

- ▶ $\text{power} \leq 30 \Rightarrow r = 2^{\text{power}}$
- ▶ $\text{power} = 31 \Rightarrow r = -2^{31}$
- ▶ $32 \leq \text{power} \leq 64$: OverflowError:
signed integer is greater than maximum
- ▶ $\text{power} \geq 64$: OverflowError:
Python int too large to convert to C long

How to analyze multilanguage programs?

Type annotations

```
class Counter:  
    def __init__(self): ...  
    def incr(self, i: int = 1): ...  
    def get(self) -> int: ...
```

How to analyze multilanguage programs?

Type annotations

```
class Counter:  
    def __init__(self): ...  
    def incr(self, i: int = 1): ...  
    def get(self) -> int: ...
```

- ▶ No raised exceptions \implies missed errors

How to analyze multilanguage programs?

Type annotations

```
class Counter:  
    def __init__(self): ...  
    def incr(self, i: int = 1): ...  
    def get(self) -> int: ...
```

- ▶ No raised exceptions \implies missed errors
- ▶ Only types

How to analyze multilanguage programs?

Type annotations

```
class Counter:  
    def __init__(self): ...  
    def incr(self, i: int = 1): ...  
    def get(self) -> int: ...
```

- ▶ No raised exceptions \implies missed errors
- ▶ Only types
- ▶ Typedhed: type annotations for the standard library

How to analyze multilanguage programs?

Type annotations

```
class Counter:  
    def __init__(self): ...  
    def incr(self, i: int = 1): ...  
    def get(self) -> int: ...
```

- ▶ No raised exceptions \implies missed errors
- ▶ Only types
- ▶ Typedhed: type annotations for the standard library, used in previous work: Monat, Ouadjaout, and Miné. “Static Type Analysis by Abstract Interpretation of Python Programs”. ECOOP 2020.

How to analyze multilanguage programs?

Type annotations

Rewrite into Python code

```
class Counter:
    def __init__(self):
        self.count = 0
    def get(self):
        return self.count
    def incr(self, i=1):
        self.count += i
```

How to analyze multilanguage programs?

Type annotations

Rewrite into Python code

```
class Counter:  
    def __init__(self):  
        self.count = 0  
    def get(self):  
        return self.count  
    def incr(self, i=1):  
        self.count += i
```

- ▶ No integer wrap-around in Python

How to analyze multilanguage programs?

Type annotations

Rewrite into Python code

```
class Counter:  
    def __init__(self):  
        self.count = 0  
    def get(self):  
        return self.count  
    def incr(self, i=1):  
        self.count += i
```

- ▶ No integer wrap-around in Python
- ▶ Some effects can't be written in pure Python (e.g., read-only attributes)

How to analyze multilanguage programs?

Type annotations

Rewrite into Python code

Drawbacks of the current approaches

How to analyze multilanguage programs?

Type annotations

Rewrite into Python code

Drawbacks of the current approaches

- ▶ Not the real code

How to analyze multilanguage programs?

Type annotations

Rewrite into Python code

Drawbacks of the current approaches

- ▶ Not the real code
- ▶ Not automatic: manual conversion

How to analyze multilanguage programs?

Type annotations

Rewrite into Python code

Drawbacks of the current approaches

- ▶ Not the real code
- ▶ Not automatic: manual conversion
- ▶ Not sound: some effects are not taken into account

How to analyze multilanguage programs?

Type annotations

Rewrite into Python code

Drawbacks of the current approaches

- ▶ Not the real code
- ▶ Not automatic: manual conversion
- ▶ Not sound: some effects are not taken into account

Our approach

How to analyze multilanguage programs?

Type annotations

Rewrite into Python code

Drawbacks of the current approaches

- ▶ Not the real code
- ▶ Not automatic: manual conversion
- ▶ Not sound: some effects are not taken into account

Our approach

- ▶ Analyze both the C and Python sources

How to analyze multilanguage programs?

Type annotations

Rewrite into Python code

Drawbacks of the current approaches

- ▶ Not the real code
- ▶ Not automatic: manual conversion
- ▶ Not sound: some effects are not taken into account

Our approach

- ▶ Analyze both the C and Python sources
- ▶ Switch from one language to the other just as the program does

How to analyze multilanguage programs?

Type annotations

Rewrite into Python code

Drawbacks of the current approaches

- ▶ Not the real code
- ▶ Not automatic: manual conversion
- ▶ Not sound: some effects are not taken into account

Our approach

- ▶ Analyze both the C and Python sources
- ▶ Switch from one language to the other just as the program does
- ▶ Reuse previous analyses of C and Python

How to analyze multilanguage programs?

Type annotations

Rewrite into Python code

Drawbacks of the current approaches

- ▶ Not the real code
- ▶ Not automatic: manual conversion
- ▶ Not sound: some effects are not taken into account

Our approach

- ▶ Analyze both the C and Python sources
- ▶ Switch from one language to the other just as the program does
- ▶ Reuse previous analyses of C and Python
- ▶ Detect runtime errors in Python, in C, and at the boundary

Difficulty: shared memory

- ▶ Two distinct visions of a shared state
- ▶ Synchronization? We could perform a full state translation, but
 - the cost would be high in the analysis
 - some abstractions can be shared between Python and C

Difficulty: shared memory

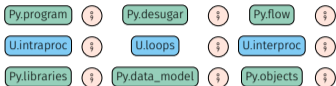
- ▶ Two distinct visions of a shared state
- ▶ Synchronization? We could perform a full state translation, but
 - the cost would be high in the analysis
 - some abstractions can be shared between Python and C

State separation \rightsquigarrow reduced synchronization

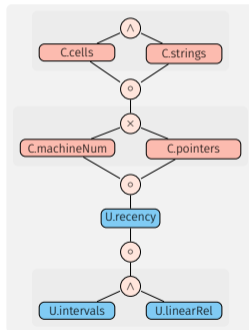
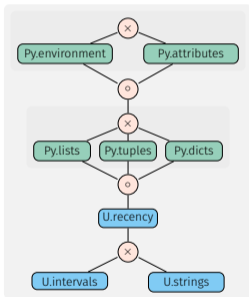
- ▶ Observation: structures are directly dereferenceable by one language only
- ▶ Switch to other language otherwise (`c.incr()` \rightsquigarrow `self->count += 1`)
Additional hypothesis: C accesses to Python objects through the API
- ▶ Synchronization: only when objects change language for the first time

Implementation & Experimental Evaluation

From distinct Python and C analyses...

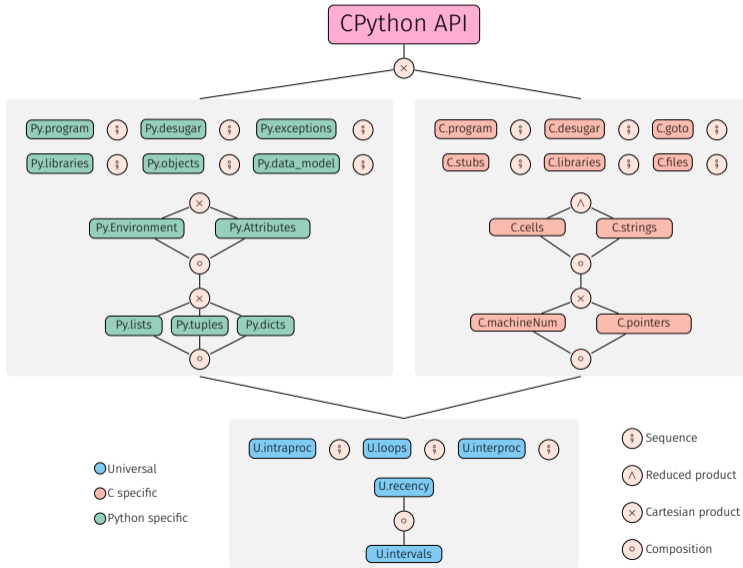


- Universal
- C specific
- Python specific

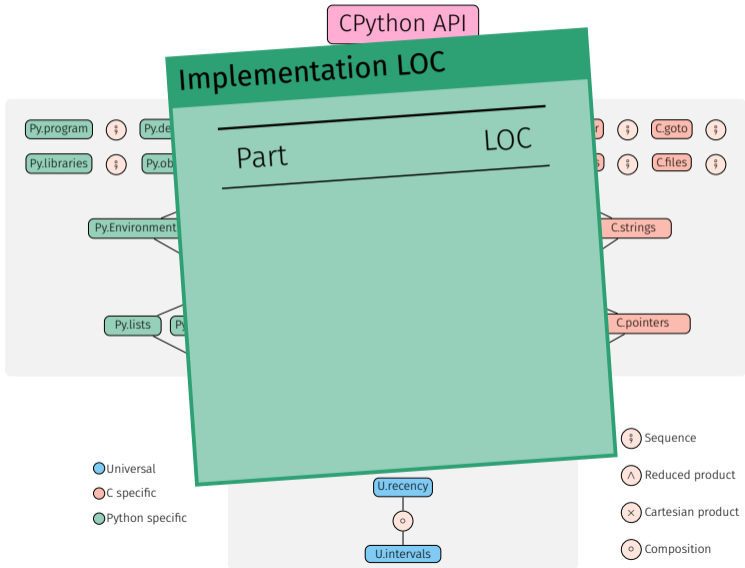


- Sequence
- ^ Reduced product
- x Cartesian product
- o Composition

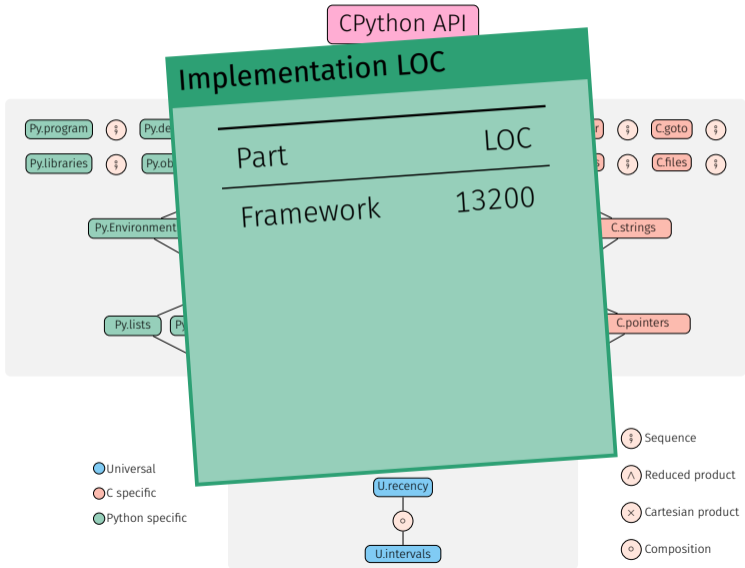
... to a multilanguage analysis!



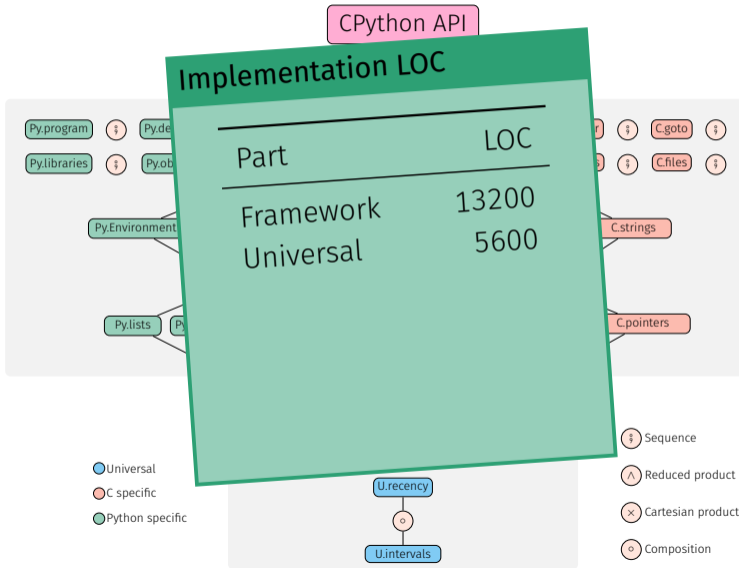
... to a multilanguage analysis!



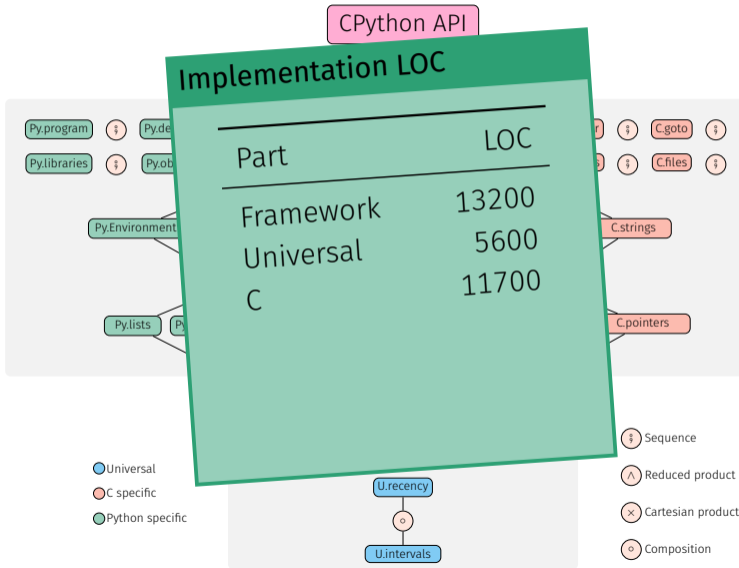
... to a multilanguage analysis!



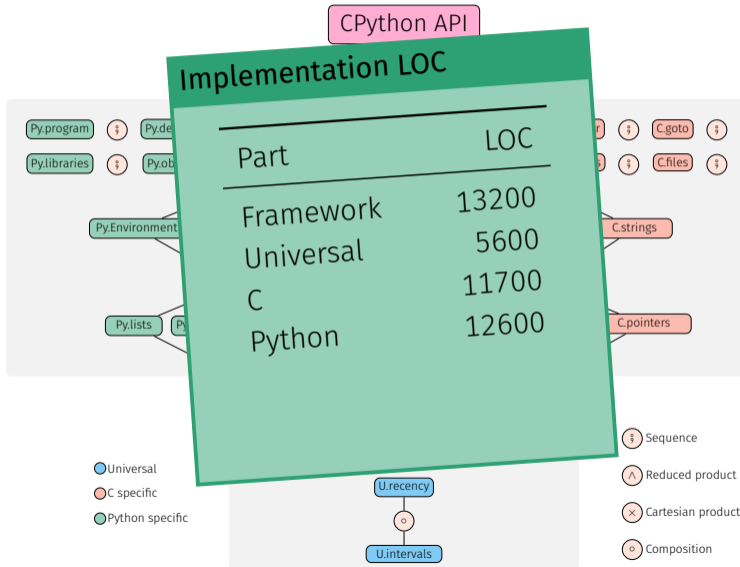
... to a multilanguage analysis!



... to a multilanguage analysis!



... to a multilanguage analysis!

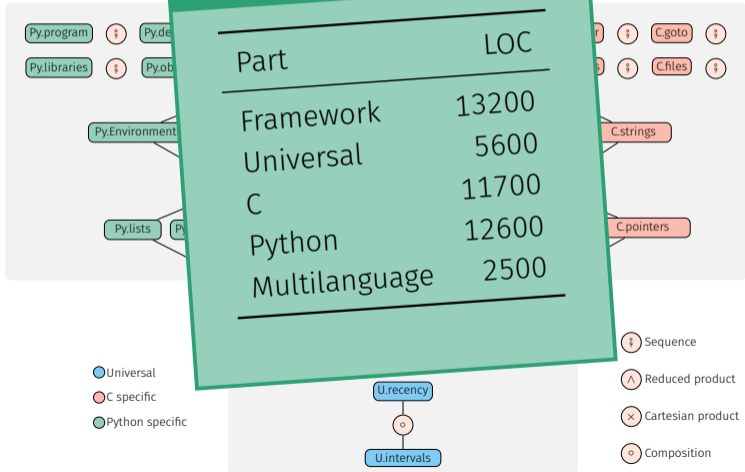


... to a multilanguage analysis!

CPython API

Implementation LOC

Part	LOC
Framework	13200
Universal	5600
C	11700
Python	12600
Multilanguage	2500



Corpus selection

- ▶ Popular, real-world libraries available on GitHub, averaging 412 stars.
- ▶ Whole-program analysis: we use the tests provided by the libraries.

Library	C + Py. Loc	Tests	🕒/test	$\frac{\# \text{ proved checks}}{\# \text{ checks}} \%$	# checks
noise	1397	15/15	1.2s	99.7%	6690
cdistance	2345	28/28	4.1s	98.0%	13716
l1ist	4515	167/194	1.5s	98.8%	36255
ahocorasick	4877	46/92	1.2s	96.7%	6722
levenshtein	5798	17/17	5.3s	84.6%	4825
bitarray	5841	159/216	1.6s	94.9%	25566

Conclusion

Difficulties

- ▶ Concrete semantics
- ▶ Memory interaction

Monat, Ouadjaout, and Miné. “A Multilanguage Static Analysis of Python Programs with Native C Extensions”. SAS 2021

Difficulties

- ▶ Concrete semantics
- ▶ Memory interaction

Previous works

- ▶ Type/exceptions analyses for the JNI
- ▶ No detection of runtime errors in C

Monat, Ouadjaout, and Miné. “A Multilanguage Static Analysis of Python Programs with Native C Extensions”. SAS 2021

Conclusion

Difficulties

- ▶ Concrete semantics
- ▶ Memory interaction

Previous works

- ▶ Type/exceptions analyses for the JNI
- ▶ No detection of runtime errors in C

Results

- ▶ Careful separation of the states and modelization of the API
- ▶ Lightweight domain on top of off-the-shelf C and Python analyses
- ▶ Shared underlying abstractions (numeric, recency)
- ▶ Scale to small, real-world libraries (using client code)

Monat, Ouadjaout, and Miné. “A Multilanguage Static Analysis of Python Programs with Native C Extensions”. SAS 2021

I'll be a permanent researcher at Inria Lille from September 2022

I'll be a permanent researcher at Inria Lille from September 2022

Considered topics:

- ▶ Ease the sound analysis of new languages
- ▶ Make static analysis more user-friendly
- ▶ Apply formal methods to legal expert systems

I'll be a permanent researcher at Inria Lille from September 2022

Considered topics:

- ▶ Ease the sound analysis of new languages
- ▶ Make static analysis more user-friendly
- ▶ Apply formal methods to legal expert systems

↪ I would be glad to set up new collaborations!

A Multilanguage Static Analysis of Python/C Programs with Mopsa

Questions

Raphaël Monat, Abdelraouf Ouadjaout, Antoine Miné

Software Languages Lab, VUB
7 July 2022

rmonat.fr

