

Static Analysis of Python Programs

A Type Abstract Domain for Python

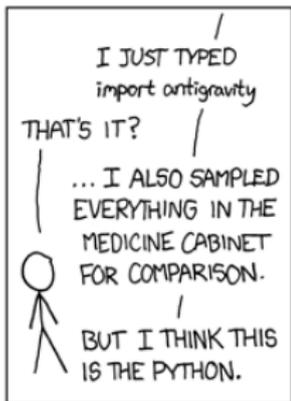
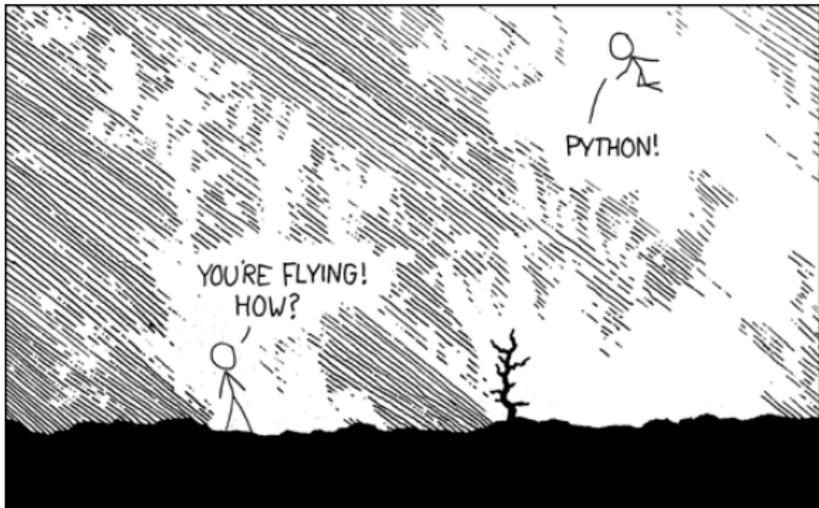
Raphaël Monat, Abdelraouf Ouadjaout, Antoine Miné

7th October 2019



**SORBONNE
UNIVERSITÉ**

Introduction



Python is a Dynamic Programming Language

It features:

- ▶ A concise and efficient syntax,
- ▶ Class-based objects,
- ▶ A large standard library,
- ▶ Dynamic typing: types are only known at runtime,
- ▶ Metaprogramming features:
 - Introspection,
 - Self-modification,
 - Metaclasses,
 - `eval`.

Python Example from the “os” Library

```
1 def fspath(p):
2     if isinstance(p, (str, bytes)):
3         return p
4     elif hasattr(p, "__fspath__"):
5         res = p.__fspath__()
6         if isinstance(res, (str, bytes)):
7             return res
8         else:
9             raise TypeError("...")
10    else:
11        raise TypeError("...")
```

Python Example from the “os” Library

Introspection, nominal types.

```
1 def fspath(p):
2     if isinstance(p, (str, bytes)):
3         return p
4     elif hasattr(p, "__fspath__"):
5         res = p.__fspath__()
6         if isinstance(res, (str, bytes)):
7             return res
8         else:
9             raise TypeError("...")
10    else:
11        raise TypeError("...")
```

Python Example from the “os” Library

Introspection, structural types.

```
1 def fspath(p):
2     if isinstance(p, (str, bytes)):
3         return p
4     elif hasattr(p, "__fspath__"):
5         res = p.__fspath__()
6         if isinstance(res, (str, bytes)):
7             return res
8         else:
9             raise TypeError("...")
10    else:
11        raise TypeError("...")
```

Python Example from the “os” Library

```
1 def fspath(p):
2     if isinstance(p, (str, bytes)):
3         return p
4     elif hasattr(p, "__fspath__"):
5         res = p.__fspath__()
6         if isinstance(res, (str, bytes)):
7             return res
8     else:
9         raise TypeError("...")
10 else:
11     raise TypeError("...")
```

Two notions of typing:

- ▶ Nominal, based on classes.
- ▶ Structural, based on attributes.

Python Example from the “os” Library

```
1 def fspath(p):
2     if isinstance(p, (str, bytes)):
3         return p
4     elif hasattr(p, "__fspath__"):
5         res = p.__fspath__()
6         if isinstance(res, (str, bytes)):
7             return res
8         else:
9             raise TypeError("...")
10    else:
11        raise TypeError("...")
```

Two notions of typing:

- ▶ Nominal, based on classes.
- ▶ Structural, based on attributes.

Type of fspath?

Python Example from the “os” Library

```
1 def fspath(p):
2     if isinstance(p, (str, bytes)):
3         return p
4     elif hasattr(p, "__fspath__"):
5         res = p.__fspath__()
6         if isinstance(res, (str, bytes)):
7             return res
8         else:
9             raise TypeError("...")
10    else:
11        raise TypeError("...")
```

Two notions of typing:

- ▶ Nominal, based on classes.
- ▶ Structural, based on attributes.

Type of fspath?

str → str; bytes → bytes and any object having a method `__fspath__` returning str or bytes.

Motivation

- ▶ Detect all runtime errors,
- ▶ At “compile time”, before the program execution.

Motivation

- ▶ Detect all runtime errors,
- ▶ At “compile time”, before the program execution.

A theoretical hurdle: Rice's theorem

“any non-trivial semantic property of programs is undecidable”.

Motivation

- ▶ Detect all runtime errors,
- ▶ At “compile time”, before the program execution.

A theoretical hurdle: Rice's theorem

“any non-trivial semantic property of programs is undecidable”.

⇒ We will compute approximate results;

Motivation

- ▶ Detect all runtime errors,
- ▶ At “compile time”, before the program execution.

A theoretical hurdle: Rice's theorem

“any non-trivial semantic property of programs is undecidable”.

⇒ We will compute approximate results;

⇒ our approximate, pessimistic approach may yield false alarms.

Motivation

- ▶ Detect all runtime errors,
- ▶ At “compile time”, before the program execution.

A theoretical hurdle: Rice's theorem

“any non-trivial semantic property of programs is undecidable”.

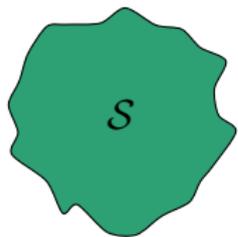
⇒ We will compute approximate results;

⇒ our approximate, pessimistic approach may yield false alarms.

Goals

- ▶ Automatic analysis: no expert knowledge required.
- ▶ Sound analysis: if no bug is detected, none will occur.

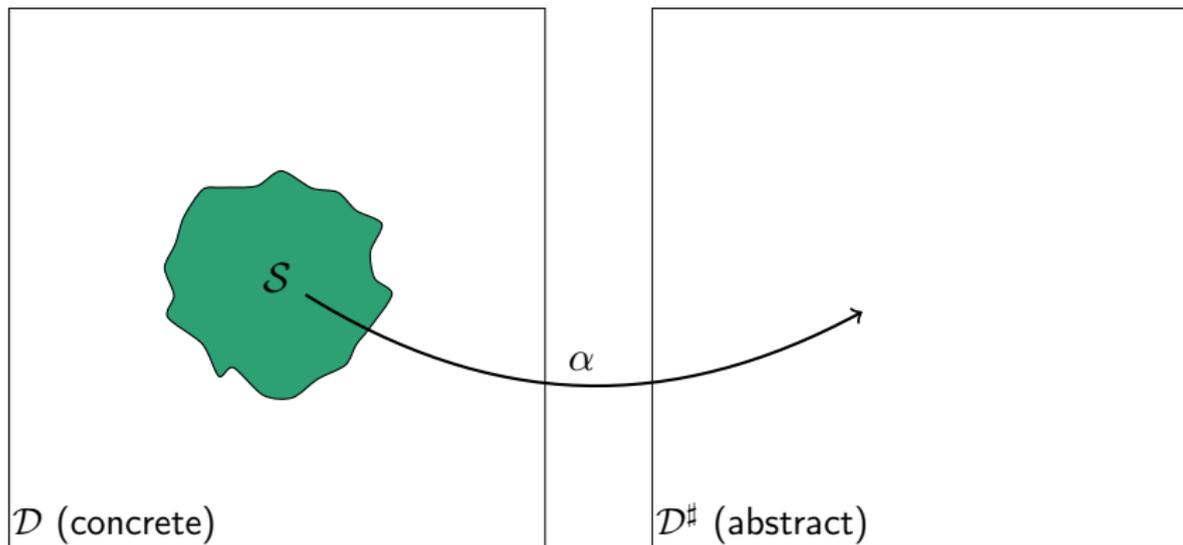
Static Analysis by Abstract Interpretation 101



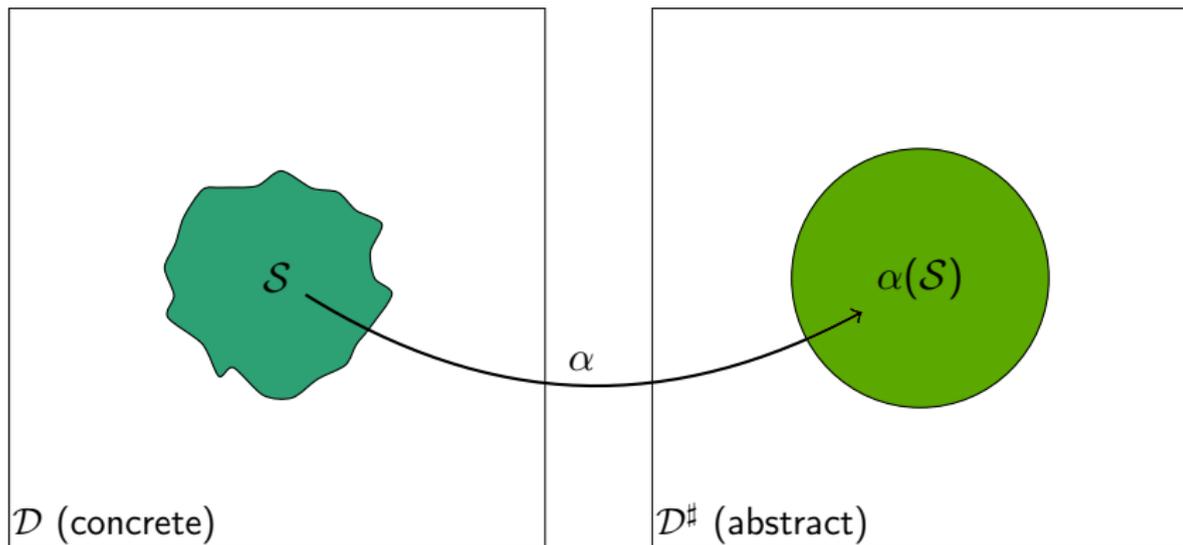
\mathcal{D} (concrete)

$\mathcal{D}^\#$ (abstract)

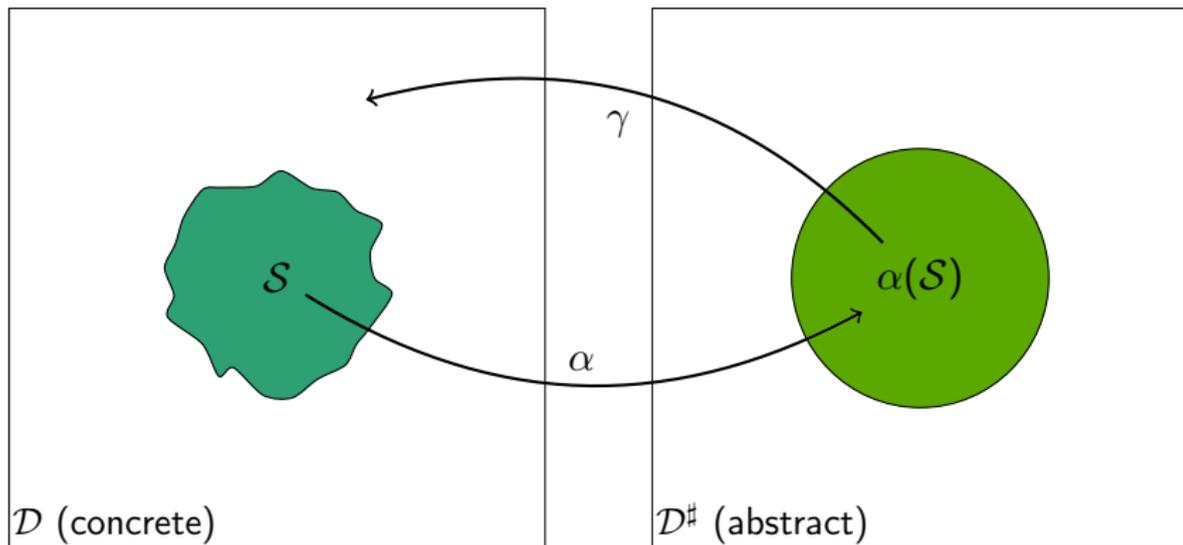
Static Analysis by Abstract Interpretation 101



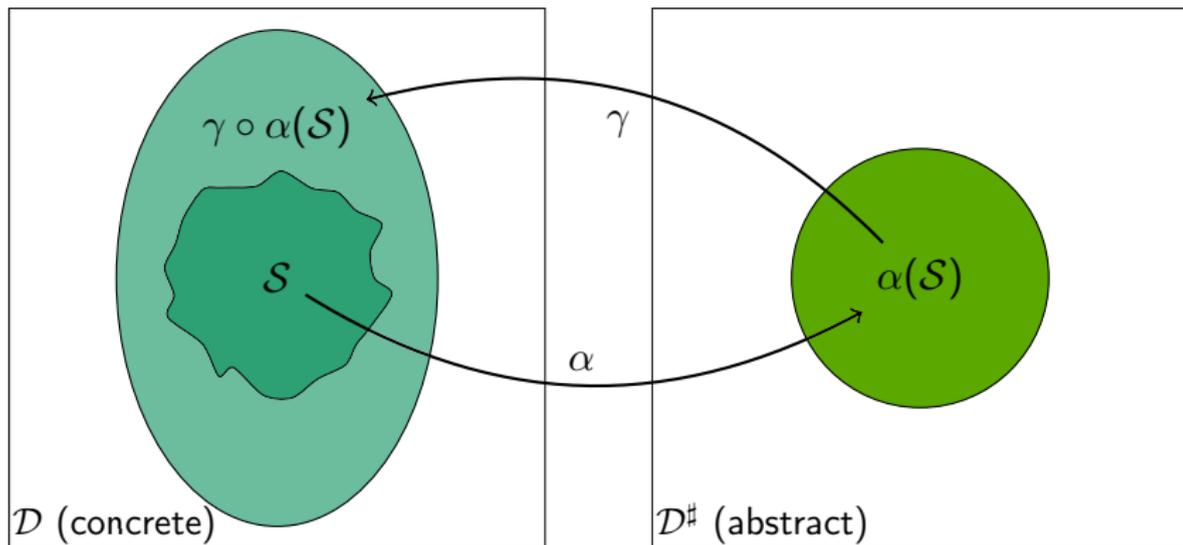
Static Analysis by Abstract Interpretation 101



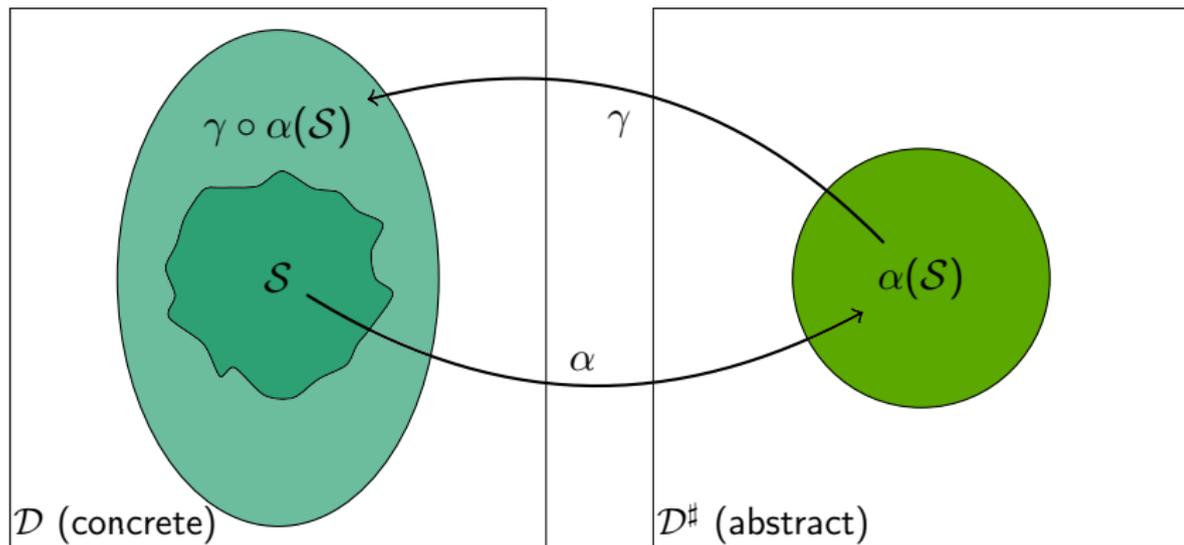
Static Analysis by Abstract Interpretation 101



Static Analysis by Abstract Interpretation 101

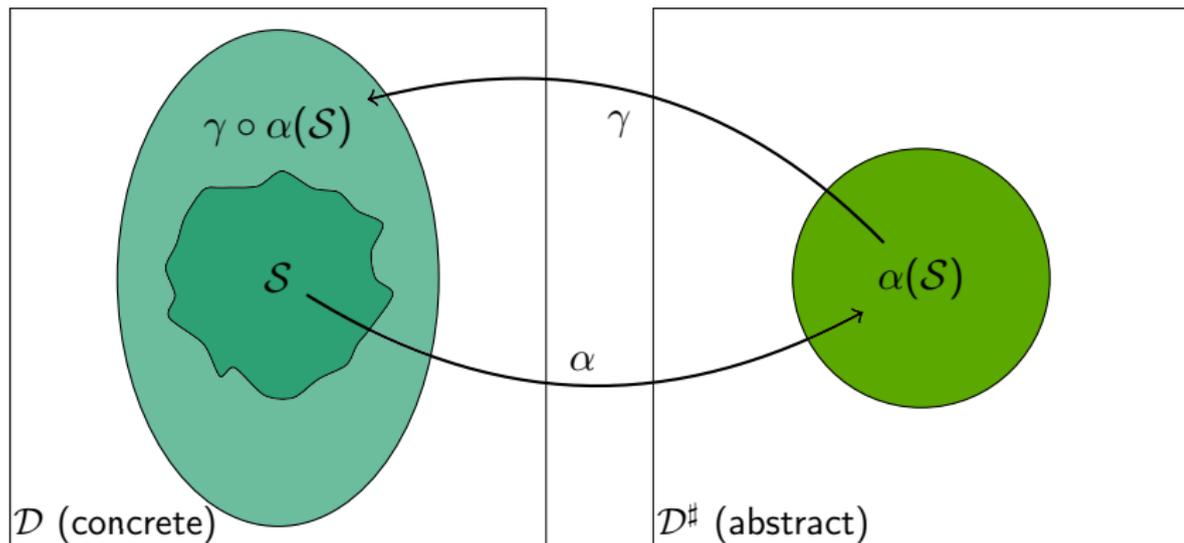


Static Analysis by Abstract Interpretation 101



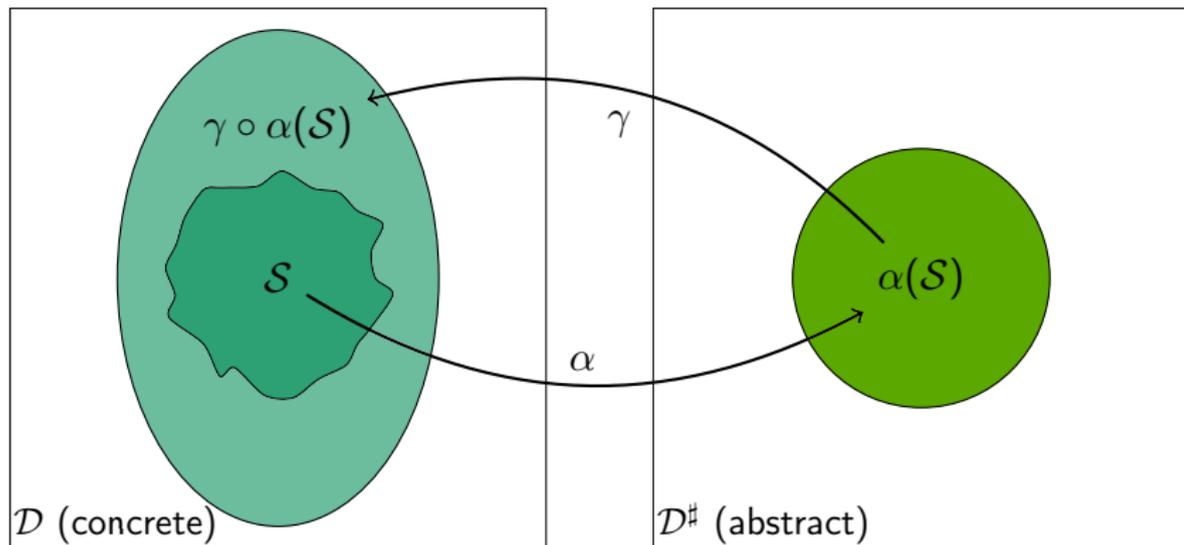
► $\mathcal{D} = \mathcal{P}(\mathbb{Z})$

Static Analysis by Abstract Interpretation 101



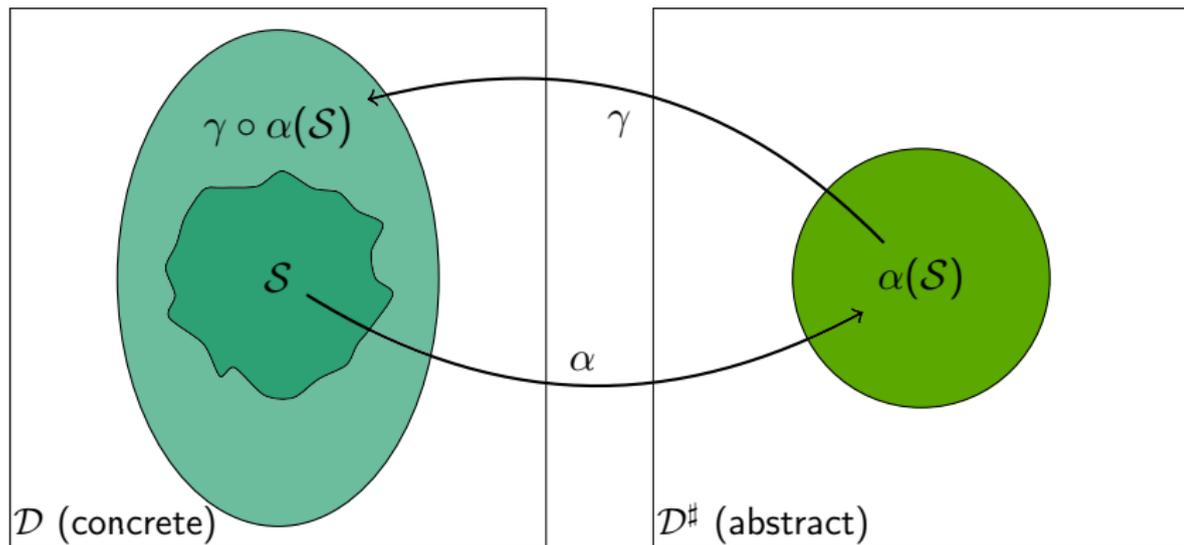
- ▶ $\mathcal{D} = \mathcal{P}(\mathbb{Z})$
- ▶ $\mathcal{D}^\sharp = \{ [a, b] \mid a \in \mathbb{Z}, b \in \overline{\mathbb{Z}}, a \leq b \}$

Static Analysis by Abstract Interpretation 101



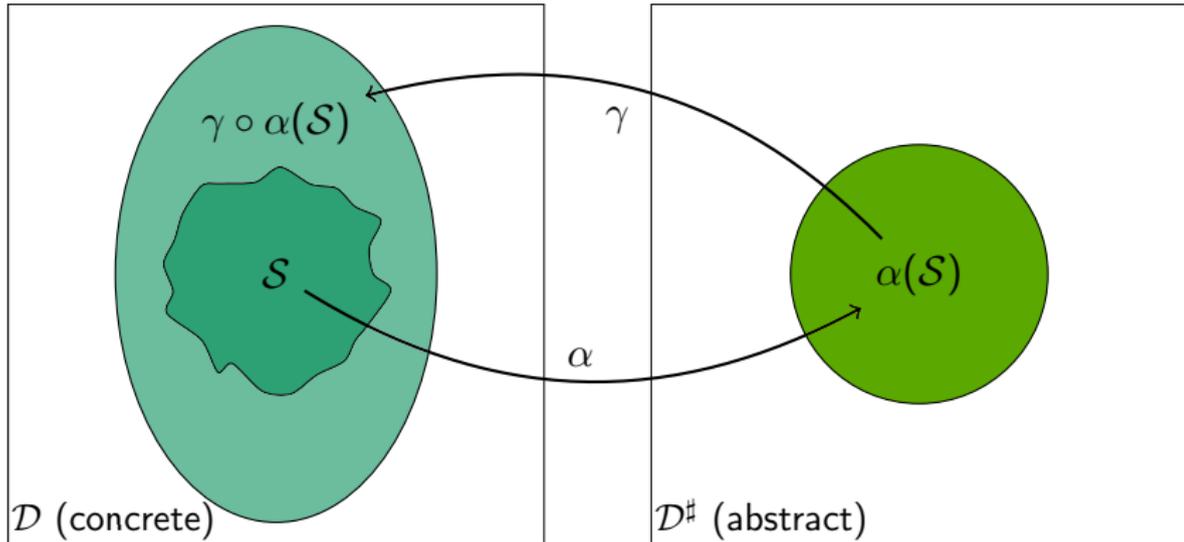
- ▶ $\mathcal{D} = \mathcal{P}(\mathbb{Z})$
- ▶ $\mathcal{D}^\sharp = \{ [a, b] \mid a \in \mathbb{Z}, b \in \overline{\mathbb{Z}}, a \leq b \}$
- ▶ $S = \{ 1, 3, 5 \}$

Static Analysis by Abstract Interpretation 101



- ▶ $\mathcal{D} = \mathcal{P}(\mathbb{Z})$
- ▶ $\mathcal{D}^\sharp = \{ [a, b] \mid a \in \mathbb{Z}, b \in \overline{\mathbb{Z}}, a \leq b \}$
- ▶ $S = \{ 1, 3, 5 \}$
- ▶ $\alpha(S) = [1, 5]$

Static Analysis by Abstract Interpretation 101



- ▶ $\mathcal{D} = \mathcal{P}(\mathbb{Z})$
- ▶ $\mathcal{D}^\# = \{ [a, b] \mid a \in \mathbb{Z}, b \in \overline{\mathbb{Z}}, a \leq b \}$
- ▶ $S = \{ 1, 3, 5 \}$
- ▶ $\alpha(S) = [1, 5]$
- ▶ $\gamma \circ \alpha(S) = \{ 1, 2, 3, 4, 5 \}$

Static analyses successfully work on critical embedded C software.
Contrary to C, Python leaves less information in the syntax.

Static analyses successfully work on critical embedded C software.
Contrary to C, Python leaves less information in the syntax.

**Static analyses could be especially helpful – though difficult –
on dynamic programming languages.**

Static analyses successfully work on critical embedded C software.
Contrary to C, Python leaves less information in the syntax.

**Static analyses could be especially helpful – though difficult –
on dynamic programming languages.**

We present a type abstract domain for Python...

Static analyses successfully work on critical embedded C software.
Contrary to C, Python leaves less information in the syntax.

**Static analyses could be especially helpful – though difficult –
on dynamic programming languages.**

We present a type abstract domain for Python...
but first let us take a look at Python's semantics.

Concrete Semantics of Python

Concrete Semantics – Motivation

Semantics? A mathematical description of the behavior of Python operators, acting like a collecting interpreter over program states.

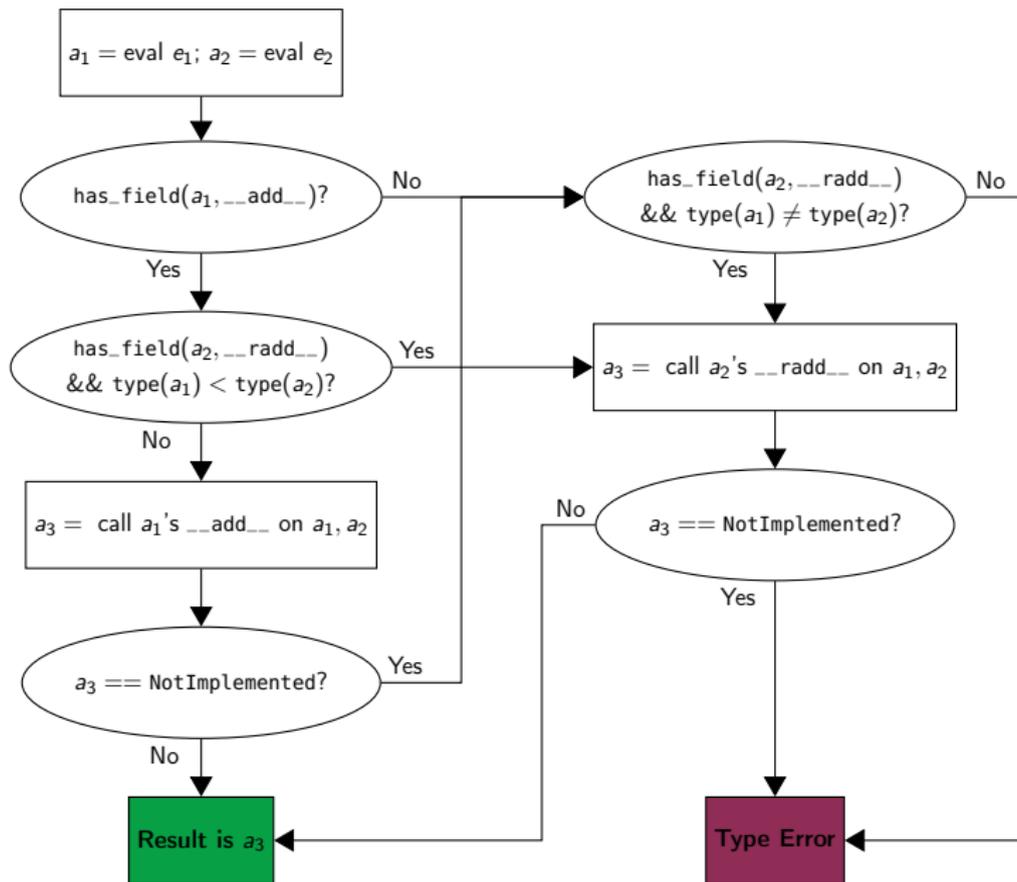
Why? To relate static analyses with the actual program behavior, and prove that our static analyses are sound.

Issues

- ▶ No standard.
- ▶ Size: Python is a huge language.
- ▶ Some parts are not formalized yet (`eval`, ...).
- ▶ Correctness: concrete semantics not implemented. But we can run our analysis on CPython's unittests.

⇒ We need a concrete semantics, but this is not our endgoal.

Semantics – Example: $e_1 + e_2$



Semantics – Example: $e_1 + e_2$

```
 $\mathbb{E}[e_1 + e_2](f, \epsilon, \sigma) \stackrel{\text{def}}{=} \\ \text{if } f \neq \text{cur} \text{ then } (f, \epsilon, \sigma) \text{ else} \\ \text{letif } (f, \epsilon, \sigma, a_1) = \mathbb{E}[e_1](f, \epsilon, \sigma) \text{ in} \\ \text{letif } (f, \epsilon, \sigma, a_2) = \mathbb{E}[e_2](f, \epsilon, \sigma) \text{ in} \\ \text{if } \text{hasattr}(\sigma(a_1), \_ \_ \text{add} \_ \_) \text{ then} \\ \quad \text{if } \text{hasattr}(\sigma(a_2), \_ \_ \text{radd} \_ \_) \wedge \text{type}(a_1) < \text{type}(a_2) \text{ then} \\ \quad \quad \text{letif } (f, \epsilon, \sigma, a_r) = \mathbb{E}[a_2.\_ \_ \text{radd} \_ \_ (a_1)] \text{ in} \\ \quad \quad \text{if } \sigma(a_r) = \text{NotImpl} \text{ then } \text{empty\_addr} \circ \mathbb{S}[\text{raise TypeError}](f, \epsilon, \sigma) \\ \quad \quad \text{else } (f, \epsilon, \sigma, a_r) \\ \quad \text{else letif } (f, \epsilon, \sigma, a_r) = \mathbb{E}[a_1.\_ \_ \text{add} \_ \_ (a_2)] \text{ in} \\ \quad \quad \text{if } \sigma(a_r) = \text{NotImpl} \text{ then} \\ \quad \quad \quad \text{if } \text{hasattr}(\sigma(a_2), \_ \_ \text{radd} \_ \_) \wedge \text{type}(a_1) \neq \text{type}(a_2) \text{ then} \\ \quad \quad \quad \quad \text{letif } (f, \epsilon, \sigma, a_r) = \mathbb{E}[a_2.\_ \_ \text{radd} \_ \_ (a_1)] \text{ in} \\ \quad \quad \quad \quad \text{if } \sigma(a_r) = \text{NotImpl} \text{ then } \text{empty\_addr} \circ \mathbb{S}[\text{raise TypeError}](f, \epsilon, \sigma) \\ \quad \quad \quad \quad \text{else } (f, \epsilon, \sigma, a_r) \\ \quad \quad \quad \text{else } (f, \epsilon, \sigma, a_r) \\ \text{else if } \text{hasattr}(\sigma(a_2), \_ \_ \text{radd} \_ \_) \wedge \text{type}(a_1) \neq \text{type}(a_2) \text{ then} \\ \quad \text{letif } (f, \epsilon, \sigma, a_r) = \mathbb{E}[a_2.\_ \_ \text{radd} \_ \_ (a_1)] \text{ in} \\ \quad \text{if } \sigma(a_r) = \text{NotImpl} \text{ then } \text{empty\_addr} \circ \mathbb{S}[\text{raise TypeError}](f, \epsilon, \sigma) \\ \quad \text{else } (f, \epsilon, \sigma, a_r) \\ \text{else } \text{empty\_addr} \circ \mathbb{S}[\text{raise TypeError}](f, \epsilon, \sigma)$ 
```

Type Abstract Domain

Features of the Analysis

Our goal: Detect uncaught exceptions, such as `TypeError`, `AttributeError`. Have a sound analysis.

Features of the Analysis

Our goal: Detect uncaught exceptions, such as `TypeError`, `AttributeError`. Have a sound analysis.

```
def dint(x):  
    if isinstance(x, int): return x*2  
    else: raise TypeError
```

```
try: z2 = dint('a')  
except TypeError: z2 = dint(1)  
# z2: int
```

Features of the Analysis

Our goal: Detect uncaught exceptions, such as `TypeError`, `AttributeError`. Have a sound analysis.

```
def dint(x):  
    if isinstance(x, int): return x*2  
    else: raise TypeError
```

```
try: z2 = dint('a')  
except TypeError: z2 = dint(1)  
# z2: int
```

⇒ Flow-sensitive analysis
including exceptions.

Features of the Analysis

Our goal: Detect uncaught exceptions, such as `TypeError`, `AttributeError`. Have a sound analysis.

```
def dint(x):  
    if isinstance(x, int): return x*2  
    else: raise TypeError
```

```
try: z2 = dint('a')  
except TypeError: z2 = dint(1)  
# z2: int
```

⇒ Flow-sensitive analysis
including exceptions.

```
class A:  
    def __init__(self):  
        self.val = 0  
    def update(self, x):  
        self.val = x  
  
x = A()  
c = x.val # c: int  
y = x  
# x, y point to the same address  
y.update('a')  
z = x.val # z: str
```

Features of the Analysis

Our goal: Detect uncaught exceptions, such as `TypeError`, `AttributeError`. Have a sound analysis.

```
def dint(x):  
    if isinstance(x, int): return x*2  
    else: raise TypeError
```

```
try: z2 = dint('a')  
except TypeError: z2 = dint(1)  
# z2: int
```

⇒ Flow-sensitive analysis
including exceptions.

```
class A:  
    def __init__(self):  
        self.val = 0  
    def update(self, x):  
        self.val = x  
  
x = A()  
c = x.val # c: int  
y = x  
# x, y point to the same address  
y.update('a')  
z = x.val # z: str
```

⇒ Handle addresses and
aliasing.

Static Type Analysis Example

```
class Path:
    def __fspath__(self): return 42

p = "/dev" if random() else Path()

def fspath(p):
    if isinstance(p, (str, bytes)):
        return p
    elif hasattr(p, "__fspath__"):
        r = p.__fspath__()
        if isinstance(r, (str, bytes)):
            return r
        else: raise TypeError
    else: raise TypeError

r = fspath(p)
```

Static Type Analysis Example

```
class Path:  
    def __fspath__(self): return 42
```

```
p = "/dev" if random() else Path()
```

```
def fspath(p):  
    if isinstance(p, (str, bytes)):  
        return p  
    elif hasattr(p, "__fspath__"):  
        r = p.__fspath__()  
        if isinstance(r, (str, bytes)):  
            return r  
        else: raise TypeError  
    else: raise TypeError
```

```
r = fspath(p)
```

Static Type Analysis Example

```
class Path:  
    def __fspath__(self): return 42
```

```
p = "/dev" if random() else Path()
```

```
def fspath(p):  
    if isinstance(p, (str, bytes)):  
        return p  
    elif hasattr(p, "__fspath__"):  
        r = p.__fspath__()  
        if isinstance(r, (str, bytes)):  
            return r  
        else: raise TypeError  
    else: raise TypeError
```

```
r = fspath(p)
```

$$\left\{ \begin{array}{l} p \mapsto \{ @str, @Path \} \\ @Path \mapsto \{ __fspath__ \} \\ @Path.__fspath__ \mapsto \{ @fun \} \end{array} \right.$$

Static Type Analysis Example

```
class Path:  
    def __fspath__(self): return 42
```

```
p = "/dev" if random() else Path()
```

```
def fspath(p):  
    if isinstance(p, (str, bytes)):  
        return p  
    elif hasattr(p, "__fspath__"):  
        r = p.__fspath__()  
        if isinstance(r, (str, bytes)):  
            return r  
        else: raise TypeError  
    else: raise TypeError
```

```
r = fspath(p)
```

$$\left\{ \begin{array}{l} p \mapsto \{ @str, @Path \} \\ @Path \mapsto \{ __fspath__ \} \\ @Path.__fspath__ \mapsto \{ @fun \} \end{array} \right.$$

Static Type Analysis Example

```
class Path:  
    def __fspath__(self): return 42
```

```
p = "/dev" if random() else Path()
```

```
def fspath(p):  
    if isinstance(p, (str, bytes)):  
        return p  
    elif hasattr(p, "__fspath__"):  
        r = p.__fspath__()  
        if isinstance(r, (str, bytes)):  
            return r  
        else: raise TypeError  
    else: raise TypeError
```

```
r = fspath(p)
```

$$\left\{ \begin{array}{l} p \mapsto \{ @str, @Path \} \\ @Path \mapsto \{ __fspath__ \} \\ @Path.__fspath__ \mapsto \{ @fun \} \\ p \mapsto \{ @str \} \end{array} \right.$$

Static Type Analysis Example

```
class Path:  
    def __fspath__(self): return 42
```

```
p = "/dev" if random() else Path() ●
```

```
def fspath(p):  
    if isinstance(p, (str, bytes)): ●  
        return p  
    elif hasattr(p, "__fspath__"):  
        r = p.__fspath__() ●  
        if isinstance(r, (str, bytes)):  
            return r  
        else: raise TypeError  
    else: raise TypeError
```

```
r = fspath(p)
```

$$\left\{ \begin{array}{l} p \mapsto \{ @str, @Path \} \\ @Path \mapsto \{ __fspath__ \} \\ @Path.__fspath__ \mapsto \{ @fun \} \\ p \mapsto \{ @str \} \end{array} \right.$$

Static Type Analysis Example

```
class Path:  
    def __fspath__(self): return 42
```

```
p = "/dev" if random() else Path() ●
```

```
def fspath(p):  
    if isinstance(p, (str, bytes)): ●  
        return p  
    elif hasattr(p, "__fspath__"):  
        r = p.__fspath__() ●  
        if isinstance(r, (str, bytes)):  
            return r  
        else: raise TypeError  
    else: raise TypeError
```

```
r = fspath(p)
```

$$\left\{ \begin{array}{l} p \mapsto \{ @str, @Path \} \\ @Path \mapsto \{ __fspath__ \} \\ @Path.__fspath__ \mapsto \{ @fun \} \end{array} \right.$$
$$p \mapsto \{ @str \}$$
$$\left\{ \begin{array}{l} p \mapsto \{ @Path \}; r \mapsto \{ @int \} \\ @Path \mapsto \{ __fspath__ \} \\ @Path.__fspath__ \mapsto \{ @fun \} \end{array} \right.$$

Static Type Analysis Example

```
class Path:  
    def __fspath__(self): return 42
```

```
p = "/dev" if random() else Path() ●
```

```
def fspath(p):  
    if isinstance(p, (str, bytes)): ●  
        return p  
    elif hasattr(p, "__fspath__"):  
        r = p.__fspath__() ●  
        if isinstance(r, (str, bytes)):  
            return r  
        else: raise TypeError  
    else: raise TypeError
```

```
r = fspath(p) ●
```

$$\left\{ \begin{array}{l} p \mapsto \{ @str, @Path \} \\ @Path \mapsto \{ __fspath__ \} \\ @Path.__fspath__ \mapsto \{ @fun \} \end{array} \right.$$
$$p \mapsto \{ @str \}$$
$$\left\{ \begin{array}{l} p \mapsto \{ @Path \}; r \mapsto \{ @int \} \\ @Path \mapsto \{ __fspath__ \} \\ @Path.__fspath__ \mapsto \{ @fun \} \end{array} \right.$$

Static Type Analysis Example

```
class Path:
```

```
    def __fspath__(self): return 42
```

```
p = "/dev" if random() else Path() ●
```

```
def fspath(p):
```

```
    if isinstance(p, (str, bytes)): ●  
        return p
```

```
    elif hasattr(p, "__fspath__"):
```

```
        r = p.__fspath__() ●
```

```
        if isinstance(r, (str, bytes)):  
            return r
```

```
        else: raise TypeError
```

```
    else: raise TypeError
```

```
r = fspath(p) ●
```

$$\left\{ \begin{array}{l} p \mapsto \{ @str, @Path \} \\ @Path \mapsto \{ __fspath__ \} \\ @Path.__fspath__ \mapsto \{ @fun \} \end{array} \right.$$

$$p \mapsto \{ @str \}$$

$$\left\{ \begin{array}{l} p \mapsto \{ @Path \}; r \mapsto \{ @int \} \\ @Path \mapsto \{ __fspath__ \} \\ @Path.__fspath__ \mapsto \{ @fun \} \end{array} \right.$$

$$\text{TypeError} \vee r \mapsto \{ @str \}$$

Classical Typing vs Static Type Analysis

“Classical” Typing

Valid programs may be **rejected**

Our Approach

No restriction on the language

Classical Typing vs Static Type Analysis

"Classical" Typing

Valid programs may be **rejected**

Type errors are **fatal**

Our Approach

No restriction on the language

Type errors are **catchable exceptions**

Classical Typing vs Static Type Analysis

"Classical" Typing

Valid programs may be **rejected**

Type errors are **fatal**

Flow-insensitive analysis

Our Approach

No restriction on the language

Type errors are **catchable exceptions**

Flow-sensitive analysis (dynamic typing & exceptions)

Classical Typing vs Static Type Analysis

“Classical” Typing

Valid programs may be **rejected**

Type errors are **fatal**

Flow-insensitive analysis

Immutable types

Our Approach

No restriction on the language

Type errors are **catchable exceptions**

Flow-sensitive analysis (dynamic typing & exceptions)

Dynamic attribute addition **changes types**

Classical Typing vs Static Type Analysis

“Classical” Typing

Valid programs may be **rejected**

Type errors are **fatal**

Flow-insensitive analysis

Immutable types

Parametric polymorphism

Our Approach

No restriction on the language

Type errors are **catchable exceptions**

Flow-sensitive analysis (dynamic typing & exceptions)

Dynamic attribute addition **changes types**

Similar, relational domain

Classical Typing vs Static Type Analysis

“Classical” Typing

Valid programs may be **rejected**

Type errors are **fatal**

Flow-insensitive analysis

Immutable types

Parametric polymorphism

Functions are analyzed in **isolation**

Top-down analysis

Our Approach

No restriction on the language

Type errors are **catchable exceptions**

Flow-sensitive analysis (dynamic typing & exceptions)

Dynamic attribute addition **changes types**

Similar, relational domain

More costly, context-sensitive interprocedural analysis

Bottom-up analysis

Classical Typing vs Static Type Analysis

"Classical" Typing	Our Approach
Valid programs may be rejected	No restriction on the language
Type errors are fatal	Type errors are catchable exceptions
Flow-insensitive analysis	Flow-sensitive analysis (dynamic typing & exceptions)
Immutable types	Dynamic attribute addition changes types
Parametric polymorphism	Similar, relational domain
Functions are analyzed in isolation	More costly, context-sensitive interprocedural analysis
Top-down analysis	Bottom-up analysis

In addition, we can easily combine our analysis with other analyses (value analysis, ...).

Description of the Type Abstract Domain

Abstract environment (nominal types) $\mathcal{E}^\# = \mathbf{Id} \rightarrow \mathcal{P}(\mathbf{Addr}^\#)$

Description of the Type Abstract Domain

Abstract environment (nominal types) $\mathcal{E}^\# = \mathbf{Id} \rightarrow \mathcal{P}(\mathbf{Addr}^\#)$

Abstract Addresses ($\mathbf{Addr}^\#$) have:

- ▶ A kind, representing the nominal type.
- ▶ A mode:
 - **weak**, meaning it summarizes multiple concrete addresses (only weak updates are possible on those addresses),
 - **strong**, abstracting one concrete address.

Description of the Type Abstract Domain

Abstract environment (nominal types) $\mathcal{E}^\# = \mathbf{Id} \rightarrow \mathcal{P}(\mathbf{Addr}^\#)$

Abstract Addresses ($\mathbf{Addr}^\#$) have:

- ▶ A kind, representing the nominal type.
- ▶ A mode:
 - **weak**, meaning it summarizes multiple concrete addresses (only weak updates are possible on those addresses),
 - **strong**, abstracting one concrete address.

We use the recency abstraction (Balakrishnan and Reps, “Recency-Abstraction for Heap-Allocated Storage”).

Description of the Type Abstract Domain

Attribute abstraction (structural types) $\mathcal{S}^\# = \mathbf{Addr}^\# \rightarrow \mathbf{Attr}^\#$

Description of the Type Abstract Domain

Attribute abstraction (structural types) $S^\# = \mathbf{Addr}^\# \rightarrow \mathbf{Attr}^\#$

Attribute under/over-approximation

$$\mathbf{Attr}^\# = \mathcal{P}(\mathit{string})^2$$

$$\gamma_{\mathbf{Attr}}^\# : \begin{cases} \mathbf{Attr}^\# & \rightarrow \mathcal{P}(\mathcal{P}(\mathit{string})) \\ (u, o) & \mapsto \{s \mid u \subseteq s \subseteq o\} \end{cases}$$

$$\gamma_{\mathbf{Attr}}^\#(\{a\}, \{a, b, c\}) = \{\{a\}, \{a, b\}, \{a, c\}, \{a, b, c\}\}$$

Description of the Type Abstract Domain

Abstract environment (nominal types)

$$\mathcal{E}^\# = \text{Id} \rightarrow \mathcal{P}(\text{Addr}^\#)$$

Attribute abstraction (structural types)

$$\mathcal{S}^\# = \text{Addr}^\# \rightarrow \text{Attr}^\#$$

$$\mathcal{S}^\# \llbracket x = e \rrbracket (\epsilon \in \mathcal{E}^\#, \sigma \in \mathcal{S}^\#) \stackrel{\text{def}}{=} \llbracket x = e \rrbracket (\epsilon, \sigma)$$

let **Obj** \mathcal{C} , $(\epsilon, \sigma) = \mathbb{E}^\# \llbracket e \rrbracket (\epsilon, \sigma)$ in

$\epsilon[x \mapsto \{ \mathcal{C} \}], \sigma$

Description of the Type Abstract Domain

Abstract environment (nominal types)

$$\mathcal{E}^\# = \text{Id} \rightarrow \mathcal{P}(\text{Addr}^\#)$$

Attribute abstraction (structural types)

$$\mathcal{S}^\# = \text{Addr}^\# \rightarrow \text{Attr}^\#$$

$$\mathcal{S}^\# \llbracket x = e \rrbracket (\epsilon \in \mathcal{E}^\#, \sigma \in \mathcal{S}^\#) \stackrel{\text{def}}{=}$$

let **Obj** \mathcal{O} , $(\epsilon, \sigma) = \mathbb{E}^\# \llbracket e \rrbracket (\epsilon, \sigma)$ in

$\epsilon[x \mapsto \{ \mathcal{O} \}], \sigma$

$$\mathbb{E}^\# \llbracket \text{object} _ _ \text{new} _ _ (e) \rrbracket (\epsilon, \sigma) \stackrel{\text{def}}{=}$$

let ee , $(\epsilon, \sigma) = \mathbb{E}^\# \llbracket e \rrbracket (\epsilon, \sigma)$ in

if $ee = \text{Class } c$ then

let \mathcal{O} , $(\epsilon, \sigma) = \mathbb{E}^\# \llbracket \text{alloc}(\text{Inst } c) \rrbracket (\epsilon, \sigma)$ in **Obj** \mathcal{O} , (ϵ, σ)

else $\text{empty_eval} \circ \mathcal{S}^\# \llbracket \text{raise TypeError} \rrbracket (\epsilon, \sigma)$

Description of the Type Abstract Domain

Abstract environment (nominal types) $\mathcal{E}^\# = \mathbf{Id} \rightarrow \mathcal{P}(\mathbf{Addr}^\#)$

Attribute abstraction (structural types) $\mathcal{S}^\# = \mathbf{Addr}^\# \rightarrow \mathbf{Attr}^\#$

```
 $S^\# \llbracket \text{object.}\_\_\text{setattribute}\_\_ (x \in \mathbf{Id}, \text{attr} \in \text{string}, e) \rrbracket (\epsilon, \sigma) \stackrel{\text{def}}{=} \\ \text{let } \mathbf{Obj} \ @_x, (\epsilon, \sigma) = \mathbb{E}^\# \llbracket x \rrbracket (\epsilon, \sigma) \text{ in} \\ \text{let attr\_var} = \mathbf{Var} (@_x, \text{attr}) \text{ in} \\ \text{if strong\_addr } @_x \text{ then} \\ \quad S^\# \llbracket \text{attr\_var} = e \rrbracket (\epsilon, \text{add\_under}(\sigma, @_x, \text{attr})) \\ \text{else} \\ \quad S^\# \llbracket \text{attr\_var} \stackrel{\text{weak}}{=} e \rrbracket (\epsilon, \text{add\_over}(\sigma, @_x, \text{attr}))$ 
```

Description of the Type Abstract Domain – Polymorphism

```
1 def get_sep(s):
2     if isinstance(s, str): return '/'
3     elif isinstance(s, bytes): return b'/'
4     else: raise TypeError
5
6     if *: r = '/dev/null'
7     else: r = b'/dev/null'
8
9     sep = get_sep(r)
```

$$\underbrace{\begin{cases} r \in \{ @_{str} \} \\ sep \in \{ @_{str} \} \end{cases}}_{\text{if branch}}$$

Description of the Type Abstract Domain – Polymorphism

```
1 def get_sep(s):
2     if isinstance(s, str): return '/'
3     elif isinstance(s, bytes): return b'/'
4     else: raise TypeError
5
6     if *: r = '/dev/null'
7     else: r = b'/dev/null'
8
9     sep = get_sep(r)
```

$$\underbrace{\begin{cases} r \in \{ @_{str} \} \\ sep \in \{ @_{str} \} \end{cases}}_{\text{if branch}}$$

$$\underbrace{\begin{cases} r \in \{ @_{bytes} \} \\ sep \in \{ @_{bytes} \} \end{cases}}_{\text{else branch}}$$

Description of the Type Abstract Domain – Polymorphism

```
1 def get_sep(s):
2     if isinstance(s, str): return '/'
3     elif isinstance(s, bytes): return b'/'
4     else: raise TypeError
5
6     if *: r = '/dev/null'
7     else: r = b'/dev/null'
8
9     sep = get_sep(r)
```

$$\underbrace{\left\{ \begin{array}{l} r \in \{ @_{str} \} \\ sep \in \{ @_{str} \} \end{array} \right\}}_{\text{if branch}} \sqcup \underbrace{\left\{ \begin{array}{l} r \in \{ @_{bytes} \} \\ sep \in \{ @_{bytes} \} \end{array} \right\}}_{\text{else branch}}$$

Description of the Type Abstract Domain – Polymorphism

```
1 def get_sep(s):
2     if isinstance(s, str): return '/'
3     elif isinstance(s, bytes): return b'/'
4     else: raise TypeError
5
6     if *: r = '/dev/null'
7     else: r = b'/dev/null'
8
9     sep = get_sep(r)
```

$$\underbrace{\begin{cases} r \in \{\textcircled{str}\} \\ sep \in \{\textcircled{str}\} \end{cases}}_{\text{if branch}} \sqcup \underbrace{\begin{cases} r \in \{\textcircled{bytes}\} \\ sep \in \{\textcircled{bytes}\} \end{cases}}_{\text{else branch}} = \begin{cases} r \in \{\textcircled{str}, \textcircled{bytes}\} \\ sep \in \{\textcircled{str}, \textcircled{bytes}\} \end{cases}$$

Description of the Type Abstract Domain – Polymorphism

```
1 def get_sep(s):
2     if isinstance(s, str): return '/'
3     elif isinstance(s, bytes): return b'/'
4     else: raise TypeError
5
6     if *: r = '/dev/null'
7     else: r = b'/dev/null'
8
9     sep = get_sep(r)
```

$$\underbrace{\left\{ \begin{array}{l} r \in \{ @_{str} \} \\ sep \in \{ @_{str} \} \end{array} \right\}}_{\text{if branch}} \sqcup \underbrace{\left\{ \begin{array}{l} r \in \{ @_{bytes} \} \\ sep \in \{ @_{bytes} \} \end{array} \right\}}_{\text{else branch}} = \left\{ \begin{array}{l} r \in \{ @_{str}, @_{bytes} \} \\ sep \in \{ @_{str}, @_{bytes} \} \end{array} \right\} \wedge r \equiv sep$$

Description of the Type Abstract Domain – Polymorphism

```
1 def get_sep(s):
2     if isinstance(s, str): return '/'
3     elif isinstance(s, bytes): return b'/'
4     else: raise TypeError
5
6     if *: r = '/dev/null'
7     else: r = b'/dev/null'
8
9     sep = get_sep(r)
```

$$\underbrace{\begin{cases} r \in \{\textcircled{str}\} \\ sep \in \{\textcircled{str}\} \end{cases}}_{\text{if branch}} \sqcup \underbrace{\begin{cases} r \in \{\textcircled{bytes}\} \\ sep \in \{\textcircled{bytes}\} \end{cases}}_{\text{else branch}} = \begin{cases} r \in \{\textcircled{str}, \textcircled{bytes}\} \\ sep \in \{\textcircled{str}, \textcircled{bytes}\} \end{cases} \wedge r \equiv sep$$

Then, we can analyze $res = r + sep$ without false alarms.

The polymorphism improves the precision.

A Modular List Domain

Ideas:

- ▶ Smash each list into one weak, abstract contents variable.
- ▶ The contents variable is built upon the list's abstract address.
- ▶ Delegate most of the work to the other domains.

A Modular List Domain

Ideas:

- ▶ Smash each list into one weak, abstract contents variable.
- ▶ The contents variable is built upon the list's abstract address.
- ▶ Delegate most of the work to the other domains.

$$\mathbb{E}^\# \llbracket [e_1, \dots, e_n]^{loc} \rrbracket s \stackrel{\text{def}}{=} \\ \text{let } @, s = \mathbb{E}^\# \llbracket \text{alloc}(\mathbf{List} \text{ } loc) \rrbracket s \text{ in} \\ \text{let contents} = \text{Var}(@, \text{"contents"}) \text{ in} \\ \text{let } s = \mathbb{S}^\# \llbracket \text{contents} \stackrel{\text{weak}}{=} e_n \rrbracket \circ \dots \circ \mathbb{S}^\# \llbracket \text{contents} \stackrel{\text{weak}}{=} e_1 \rrbracket s \text{ in} \\ \text{Obj } @, s$$

Lists and Polymorphisms

```
1 if *: l = [1, 2, 3]
2 else: l = [0.1, 0.2, 0.3]
3 x = l[0]
```

After line 1:

$$l \mapsto \{ @list1 \}$$
$$@list1.content \mapsto \{ @int \}$$

After line 2:

$$l \mapsto \{ @list2 \}$$
$$@list2.content \mapsto \{ @float \}$$

Lists and Polymorphisms

```
1 if *: l = [1, 2, 3]
2 else: l = [0.1, 0.2, 0.3]
3 x = l[0]
```

After line 1:

$$l \mapsto \{ @list1 \}$$
$$@list1.content \mapsto \{ @int \}$$

After line 2:

$$l \mapsto \{ @list2 \}$$
$$@list2.content \mapsto \{ @float \}$$

After line 3 (join):

$$\left\{ \begin{array}{l} l \mapsto \{ @list1, @list2 \} \\ @list1.content \mapsto \{ @int \} \\ @list2.content \mapsto \{ @float \} \\ x \mapsto \{ @int, @float \} \end{array} \right.$$

Lists and Polymorphisms

```
1 if *: l = [1, 2, 3]
2 else: l = [0.1, 0.2, 0.3]
3 x = l[0]
```

After line 1:

$$l \mapsto \{ @list1 \}$$
$$@list1.content \mapsto \{ @int \}$$

After line 2:

$$l \mapsto \{ @list2 \}$$
$$@list2.content \mapsto \{ @float \}$$

After line 3 (join):

$$\left\{ \begin{array}{l} l \mapsto \{ @list1, @list2 \} \\ @list1.content \mapsto \{ @int \} \\ @list2.content \mapsto \{ @float \} \\ x \mapsto \{ @int, @float \} \end{array} \right.$$

After unification and join:

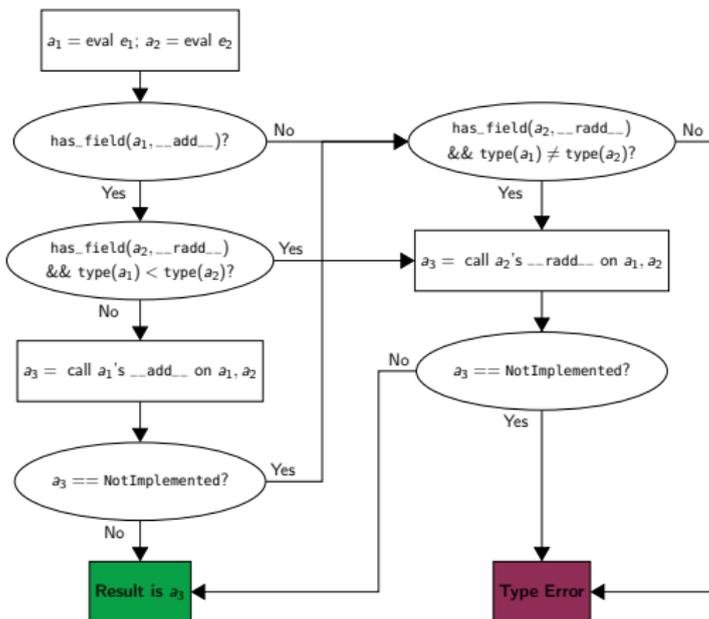
$$\left\{ \begin{array}{l} l \mapsto \{ @list_{\{1,2\}} \} \\ @list_{\{1,2\}}.content \mapsto \{ @int, @float \} \\ x \mapsto \{ @int, @float \} \\ x \equiv @list_{\{1,2\}}.content \end{array} \right.$$

Interprocedural Analysis

What about a context-insensitive analysis?

```
def f(e1, e2): return e1 + e2
```

f is valid whenever the evaluation hits the green box:



What about a context-insensitive analysis?

```
def f(e1, e2): return e1 + e2
```

⇒ We focus on a context-sensitive analysis.

What about a context-insensitive analysis?

```
def f(e1, e2): return e1 + e2
```

⇒ We focus on a context-sensitive analysis.

Inlining most precise, but costly.

What about a context-insensitive analysis?

```
def f(e1, e2): return e1 + e2
```

⇒ We focus on a context-sensitive analysis.

Inlining most precise, but costly.

Towards function summaries a simple cache keeping the previous function analyses achieves up to 25x speedup on our benchmarks.

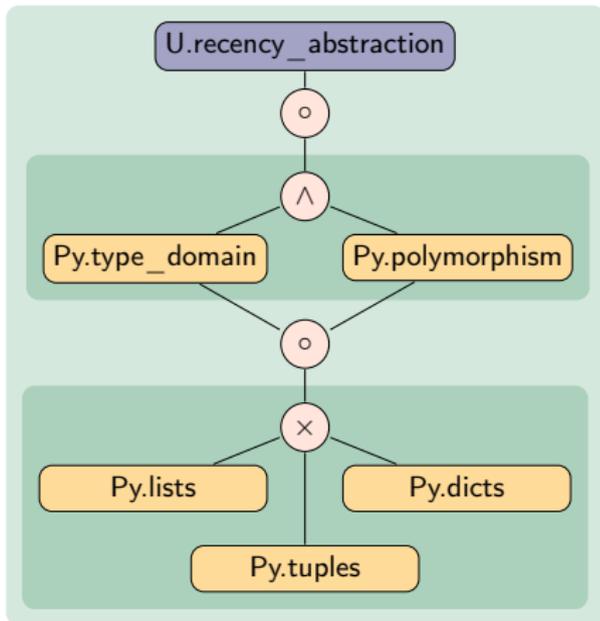
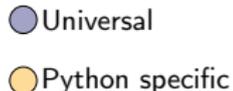
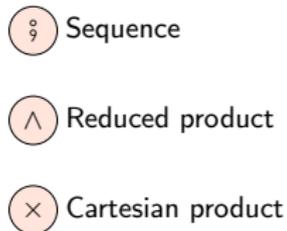
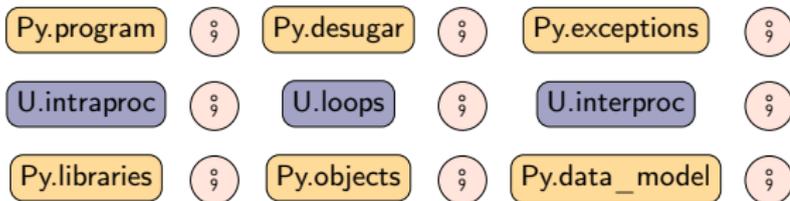
Experimental Evaluation

Modular Open Platform for Static Analysis¹

- ▶ Modular abstract domains are small “blocks”, handling everything from: abstract values to control-flow statements.
- ▶ Statements flow through these domains until one answers.
- ▶ The user can select the combination of abstract domains.
- ▶ Supports Python and C analysis
(some parts are shared in a “universal” language).

¹Journault et al. “Combinations of reusable abstract domains for a multilingual static analyzer”. VSTTE19 Proceedings.

Implementation into MOPSA



Official Python Benchmarks:

Name	LOC	Time (inlining)	Time (cache)	# Alarms	# False Alarms
fannkuch.py	59	0.07s	0.07s	0	0
float.py	63	0.10s	0.06s	0	0
spectral_n.py	74	3.9 s	0.33s	0	1
nbody.py	157	2.6s	1.5s	0	1
chaos.py	324	19s	5.9s	1	0
unpack_seq.py	458	5.6s	5.4s	0	0
hexiom.py	674	61.7m	2.2m	0	52

Official Python Benchmarks:

Name	LOC	Time (inlining)	Time (cache)	# Alarms	# False Alarms
fannkuch.py	59	0.07s	0.07s	0	0
float.py	63	0.10s	0.06s	0	0
spectral_n.py	74	3.9 s	0.33s	0	1
nbody.py	157	2.6s	1.5s	0	1
chaos.py	324	19s	5.9s	1 ²	0
unpack_seq.py	458	5.6s	5.4s	0	0
hexiom.py	674	61.7m	2.2m	0	52

²A real bug was found: a piece of currently unused code was working in Python 2.x, but not in Python 3.x.

Related Work

JavaScript: more work has been done on JavaScript, including:

- ▶ Bodin et al.: semantics in Coq,
- ▶ Jensen, Møller, and Thiemann: Value Analysis for JavaScript.

Dynamic Analysis: Pyannotate and MonkeyType collect the types of a program at runtime. The types are valid for the explored trace.

Gradual Typing: annotated program parts are typechecked, while other parts have an unknown “top” type, from which any static type can be cast to and from. In Python: Mypy, Pyre.

Static Analysis of Python

- ▶ Fromherz et al³: a value analysis for Python (more costly). We have started a modular implementation of this analysis using our type domain.
- ▶ Typpete⁴: encodes type inference into a MaxSMT instance.
- ▶ Fritz & Hage⁵: performs a dataflow analysis.
- ▶ Pytype, tool from Google, performing a dataflow analysis.

³Fromherz, Ouadjaout, and Miné. “Static Value Analysis of Python Programs by Abstract Interpretation”. NFM 2018 Proceedings.

⁴Hassan et al. “MaxSMT-Based Type Inference for Python 3”. CAV 2018 Proceedings.

⁵Fritz and Hage. “Cost versus precision for approximate typing for Python”. PEPM 2017 Proceedings.

Static Analysis of Python

- ▶ Fromherz et al³: a value analysis for Python (more costly). We have started a modular implementation of this analysis using our type domain.
- ▶ Typpete⁴: encodes type inference into a MaxSMT instance.
- ▶ Fritz & Hage⁵: performs a dataflow analysis.
- ▶ Pytype, tool from Google, performing a dataflow analysis.

Our type analysis uniquely takes into account object mutability and control-flow.

³Fromherz, Ouadjaout, and Miné. "Static Value Analysis of Python Programs by Abstract Interpretation". NFM 2018 Proceedings.

⁴Hassan et al. "MaxSMT-Based Type Inference for Python 3". CAV 2018 Proceedings.

⁵Fritz and Hage. "Cost versus precision for approximate typing for Python". PEPM 2017 Proceedings.

Conclusion

Conclusion & Future Work

We have developed a static type analysis for Python.

- ▶ More precise than the state-of-the-art static type analyses.
- ▶ Analyzes real-world benchmarks!
- ▶ Sound: it relates the analysis and the concrete semantics.

Future Work

- ▶ Support official type annotations (PEP 484).
- ▶ More efficient, summary-based function analysis.
- ▶ Handle libraries.
- ▶ Analyze real-world programs.

Appendix

Implementation into MOPSA

Implementation size:

- ▶ 5500 lines of OCaml for Python's semantics,
- ▶ 2500 for Python's type abstract domain,
- ▶ 2100 for Python's containers abstractions,
- ▶ 1800 for the universal language (loop & function analysis),
- ▶ 15000 for the modular framework.

Flow-sensitive analysis

Flow-sensitive Analysis Performed by induction over the syntax.

Using flow tokens to label continuation-passed states.

$$\mathcal{F}^\# = \{ \text{cur}, \text{ret}, \text{brk}, \text{exn } a, a \in \mathbf{Addr}^\# \}$$

- ▶ Abstract environment (nominal types) $\mathcal{E}^\# = \mathbf{Id} \rightarrow \mathcal{P}(\mathbf{Addr}^\#)$
- ▶ Attribute abstraction (structural types) $\mathcal{S}^\# = \mathbf{Addr}^\# \rightarrow \mathbf{Attr}^\#$

$$\mathbb{E}^\# \llbracket x \rrbracket (f \in \mathcal{F}^\#, \epsilon \in \mathcal{E}^\#, \sigma \in \mathcal{S}^\#) \stackrel{\text{def}}{=} \bigcup_{a^\# \in \epsilon(x)} \mathbf{Obj } a^\#, (f, \epsilon[x \mapsto \{ a^\# \}], \sigma)$$

$$\mathbb{S}^\# \llbracket x = e \rrbracket (f \in \mathcal{F}^\#, \epsilon \in \mathcal{E}^\#, \sigma \in \mathcal{S}^\#) \stackrel{\text{def}}{=} \\ \text{letif } (\mathbf{Obj } @, (f, \epsilon, \sigma)) = \mathbb{E}^\# \llbracket e \rrbracket (f, \epsilon, \sigma) \text{ in} \\ (f, \epsilon[x \mapsto @], \sigma)$$

Description of the Type Abstract Domain – Polymorphism

```
1 def get_sep(s):
2     if isinstance(s, str): return '/'
3     elif isinstance(s, bytes): return b'/'
4     else: raise TypeError
5
6     if *: r = '/dev/null'
7     else: r = b'/dev/null'
8
9     sep = get_sep(r)
```

$$\underbrace{\begin{cases} r \in \{ @_{str}^w \} \\ sep \in \{ @_{str}^s \} \end{cases}}_{\text{if branch}}$$

Description of the Type Abstract Domain – Polymorphism

```
1 def get_sep(s):
2     if isinstance(s, str): return '/'
3     elif isinstance(s, bytes): return b'/'
4     else: raise TypeError
5
6     if *: r = '/dev/null'
7     else: r = b'/dev/null'
8
9     sep = get_sep(r)
```

$$\underbrace{\begin{cases} r \in \{ @_{str}^w \} \\ sep \in \{ @_{str}^s \} \end{cases}}_{\text{if branch}}$$

$$\underbrace{\begin{cases} r \in \{ @_{str}^w \} \\ sep \in \{ @_{str}^s \} \end{cases}}_{\text{else branch}}$$

Description of the Type Abstract Domain – Polymorphism

```
1 def get_sep(s):
2     if isinstance(s, str): return '/'
3     elif isinstance(s, bytes): return b'/'
4     else: raise TypeError
5
6     if *: r = '/dev/null'
7     else: r = b'/dev/null'
8
9     sep = get_sep(r)
```

$$\underbrace{\left\{ \begin{array}{l} r \in \{ @_{str}^w \} \\ sep \in \{ @_{str}^s \} \end{array} \right\}}_{\text{if branch}} \sqcup \underbrace{\left\{ \begin{array}{l} r \in \{ @_{str}^w \} \\ sep \in \{ @_{str}^s \} \end{array} \right\}}_{\text{else branch}}$$

Description of the Type Abstract Domain – Polymorphism

```
1 def get_sep(s):
2     if isinstance(s, str): return '/'
3     elif isinstance(s, bytes): return b'/'
4     else: raise TypeError
5
6     if *: r = '/dev/null'
7     else: r = b'/dev/null'
8
9     sep = get_sep(r)
```

$$\underbrace{\left\{ \begin{array}{l} r \in \{ @_{str}^w \} \\ sep \in \{ @_{str}^s \} \end{array} \right\}}_{\text{if branch}} \sqcup \underbrace{\left\{ \begin{array}{l} r \in \{ @_{str}^w \} \\ sep \in \{ @_{str}^s \} \end{array} \right\}}_{\text{else branch}} = \left\{ \begin{array}{l} r \in \{ @_{str}^w, @_{bytes}^w \} \\ sep \in \{ @_{str}^s, @_{bytes}^s \} \end{array} \right\}$$

Description of the Type Abstract Domain – Polymorphism

```
1 def get_sep(s):
2     if isinstance(s, str): return '/'
3     elif isinstance(s, bytes): return b'/'
4     else: raise TypeError
5
6     if *: r = '/dev/null'
7     else: r = b'/dev/null'
8
9     sep = get_sep(r)
```

$$\underbrace{\left\{ \begin{array}{l} r \in \{ @_{str}^w \} \\ sep \in \{ @_{str}^s \} \end{array} \right\}}_{\text{if branch}} \sqcup \underbrace{\left\{ \begin{array}{l} r \in \{ @_{str}^w \} \\ sep \in \{ @_{str}^s \} \end{array} \right\}}_{\text{else branch}} = \left\{ \begin{array}{l} r \in \{ @_{str}^w, @_{bytes}^w \} \\ sep \in \{ @_{str}^s, @_{bytes}^s \} \end{array} \right\} \wedge r \equiv sep$$

Description of the Type Abstract Domain – Polymorphism

```
1 def get_sep(s):
2     if isinstance(s, str): return '/'
3     elif isinstance(s, bytes): return b'/'
4     else: raise TypeError
5
6     if *: r = '/dev/null'
7     else: r = b'/dev/null'
8
9     sep = get_sep(r)
```

$$\underbrace{\left\{ \begin{array}{l} r \in \{ @_{str}^w \\ sep \in \{ @_{str}^s \} \end{array} \right\}}_{\text{if branch}} \sqcup \underbrace{\left\{ \begin{array}{l} r \in \{ @_{str}^w \\ sep \in \{ @_{str}^s \} \end{array} \right\}}_{\text{else branch}} = \left\{ \begin{array}{l} r \in \{ @_{str}^w, @_{bytes}^w \\ sep \in \{ @_{str}^s, @_{bytes}^s \} \end{array} \right\} \wedge r \equiv sep$$

Then, we can analyze $res = r + sep$ without false alarms.
The polymorphism improves the precision.