

RIGHTTYPER: Effective and Efficient Type Annotation for Python

Juan Altmayer Pizzorno

University of Massachusetts Amherst
Amherst, MA, USA
jpizzorno@cs.umass.edu

Emery D. Berger[†]

University of Massachusetts Amherst /
Amazon Web Services
Amherst, MA, USA
emery@cs.umass.edu

Abstract

Python type annotations bring the benefits of static type checking to the language. However, manually writing annotations can be time-consuming and tedious. The result is that most real-world Python code remains largely untyped. Past approaches to annotating types in Python code fall short in a number of ways. Static approaches struggle with dynamic features and infer overly broad types. AI-based methods are inherently unsound and can miss rare or user-defined types. Dynamic methods can impose extreme runtime overheads, degrading performance by up to 270×, abort execution as they exhaust resources, and even infer incorrect types that lead to runtime errors. Crucially, all prior work assumes implicitly *that the code to be annotated is already correct*. This assumption is generally unwarranted, especially for large codebases that have been untyped.

This paper presents RIGHTTYPER, a novel approach for Python that overcomes these disadvantages. RIGHTTYPER not only generates precise type annotations based on actual program behavior, improving recall in type checking relative to prior approaches. It also turns *type checking into anomaly detection*, allowing the type checker to identify corner cases that the programmer can audit for unintended behavior. RIGHTTYPER is also fast and space-efficient, imposing just 30% performance overhead on average. RIGHTTYPER achieves these characteristics by a principled yet pervasive use of sampling—guided by self-profiling—along with statistical filtering and careful resolution and aggregation of type information.

Keywords

Type Inference, Type Annotations, Gradual Typing, Python, Dynamic Analysis

1 Introduction

Python is now firmly established as one of the most popular programming languages. Despite being dynamically typed, Python has supported static type annotations since version 3.5, released in 2015 [36]. While the Python runtime does not enforce these type annotations, separate tools such as mypy and pyright [14, 33] use them to perform static type checking. Beyond static verification, type annotations serve multiple purposes across the Python ecosystem. They help clarify developer intent and improve code readability, making codebases easier to understand and maintain. Integrated development environments and editors such as PyCharm and VSCode leverage annotations to provide better code completion and refactoring support. Additionally, annotations support test case generation, runtime validation, and serialization frameworks (e.g., pydantic and beartype).

Despite these many uses and a decade of language support, Python codebases remain largely unannotated. A recent study of almost ten thousand popular Python project repositories finds that only 7% of them use type annotations at all. Among those with annotations, only about 8% of function arguments and return types are actually annotated [6]. This gap highlights the practical challenges developers face in manually adding annotations, especially to large codebases. To address this, *type inference* tools aim to automatically infer and insert type annotations, reducing manual effort and making static typing more accessible.

Past approaches to Python type inference fall into three broad categories: *static*, *AI-based*, and *dynamic*. Static approaches reason about code without executing it, using information such as program structure, operations, and control flow to infer types [20]. However, they are limited by dynamic language features, either supporting only subsets of the language or leaving portions of the code unannotated [11, 28]. They also tend to be imprecise, as they conservatively over-approximate program behavior [3]. For example, if a function could theoretically operate on either an integer or a string, a static tool may reflect both in its inference, even if only one is ever actually used. While technically correct, these annotations can be too permissive, diminishing the effectiveness of static checkers and diluting the effectiveness of annotations for catching bugs and guiding development.

AI-based approaches operate on source code to estimate the likely types of program elements, most commonly by applying machine learning models. These methods are better equipped than static approaches to handle the complexities of dynamic languages, but they sacrifice soundness: predicted types are not guaranteed to reflect actual program behavior, undermining the effectiveness of type checking through false positives or false negatives. Some hybrid systems attempt to address this unsoundness by using static analysis to validate model predictions [22, 26]. However, this validation can only discard incorrect type predictions, potentially leaving more of the code unannotated. Moreover, these methods struggle with rare types, and many support only a limited type vocabulary, making them ill-equipped to handle user-defined or evolving types [19].

In contrast, dynamic approaches observe types at runtime by executing the program and recording the actual values passed to and returned from functions [26]. This strategy trades static completeness for runtime fidelity: it requires code execution, but yields more precise and realistic types that reflect actual program behavior.

Unfortunately, existing dynamic type inference tools for Python—MonkeyType and PyAnnotate—fall short in key ways [17, 35]. MonkeyType logs all function calls to a SQLite database, resulting in inordinate overhead: execution can slow down by up to 270×, and logging can consume gigabytes of disk space per second [10]. Worse,

[†]Work done at Amazon Web Services.

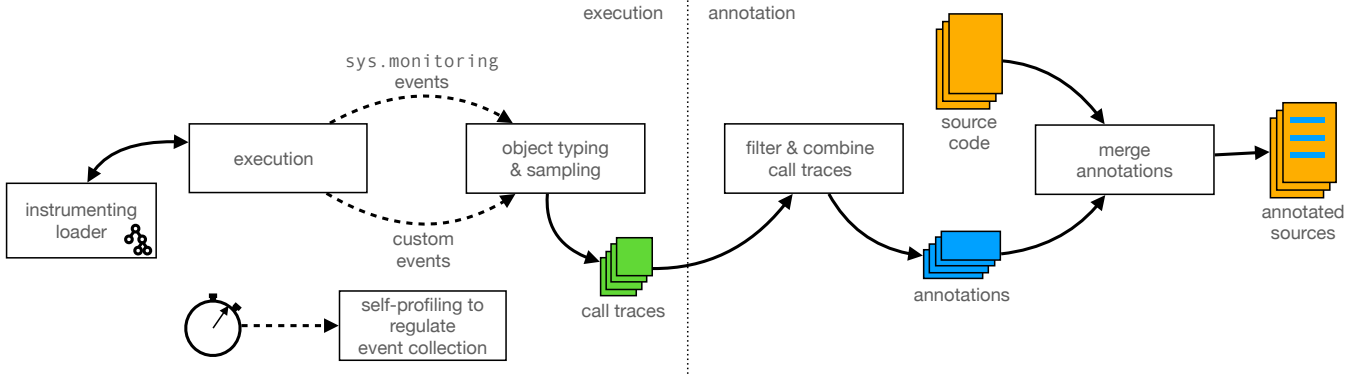


Figure 1: Overview: RIGHTTYPER enables anomaly detection through precise type annotations that capture typical program behavior. It executes the target program under instrumentation (§ 2.1, 2.2) to gather *call traces* containing the types observed at runtime. To control overhead, RIGHTTYPER profiles itself and dynamically adjusts event collection. After execution, it filters out uncommon traces and merges the remaining ones to generate type annotations (§ 2.3).

its high memory usage can also exhaust system resources and cause the program being analyzed to abort when applied to large codebases. PyAnnotate mitigates its runtime overhead through sampling, but its strategy results in biased samples. It records only a limited number of calls per function and inspects just the first four elements of data structures, leading to annotations that may not reflect typical program behavior. Moreover, it trusts Python’s often misleading type names, sometimes producing annotations that cause runtime errors.

Beyond the above drawbacks, all prior type inference systems make a strong assumption: *that the code they annotate is correct*, with approaches even reinforcing this assumption by explicitly measuring the quality of inferred types on whether they avoid static type checking errors [2, 19, 45]. This assumption is generally unwarranted, especially in large or previously untyped codebases, and as Section 3.1 shows, can lead to inferred annotations that mask bugs. For instance, a list intended to contain only numbers might erroneously include a string. An inferred type that includes both numbers and strings (e.g., `int | str`) would accurately reflect the observed data but would silently conceal the underlying issue. An ideal system should infer types that reflect *typical* behavior while surfacing anomalous cases to the developer, especially when first annotating previously untyped code.

This paper proposes RIGHTTYPER, a novel dynamic type inference approach for Python. RIGHTTYPER sidesteps the pitfalls of static and AI-based methods by sampling types observed during actual program execution. It limits overhead by profiling its execution to control instrumentation in a principled way, ensuring uniform sampling across the target program’s execution. It also efficiently samples elements from data structures such as lists and dictionaries uniformly at random (§ 2.2). RIGHTTYPER resolves type names to correct for Python’s frequently misleading introspection, and filters out types that appear only infrequently in a given context. By analyzing usage patterns, it emits annotations that reflect relationships between types. This combination of sampling and filtering turns type checking into a form of anomaly detection: rare type occurrences are excluded from annotations, prompting static checkers

to flag them for developer review as potential bugs or unexpected behavior. This paper makes the following contributions:

- It introduces RIGHTTYPER, a novel sampling-based type inference approach for Python that reframes type checking as anomaly detection;
- It presents a prototype implementation of RIGHTTYPER has been released as open-source software on GitHub [25];
- It empirically compares the types inferred by RIGHTTYPER to those from prior work, demonstrating a marked improvement in the quality of annotations;
- It evaluates RIGHTTYPER’s runtime overhead against other dynamic inference tools, showing that it is significantly more efficient.

2 Approach

RIGHTTYPER enables anomaly detection by generating precise type annotations that capture normal program behavior. It runs the target program using a combination of custom instrumentation and standard Python monitoring mechanisms. During execution, it inspects objects—such as function parameters and return values—to determine their types. Rather than tracking every event exhaustively, RIGHTTYPER employs sampling, not only to determine the complete types of container objects like lists and dictionaries, but also to regulate event collection. To manage overhead, it profiles its execution and dynamically enables or disables monitoring. It resolves type names to correct for Python’s often misleading introspection. After execution, it filters out call traces that occur only rarely in each context and merges the remaining traces to produce precise, representative type annotations. Figure 1 provides an overview.

2.1 Instrumentation and Event Sampling

Most of RIGHTTYPER’s instrumentation utilizes `sys.monitoring`, a low-overhead, fine-grained mechanism for tracking various events during execution introduced in Python 3.12 [29]. This feature allows selective activation and deactivation of monitoring for specific portions of user code. RIGHTTYPER leverages its `PY_START` and

PY_RETURN events to let it capture function arguments and return values. Additionally, it uses PY_YIELD to monitor values emitted by generator functions, which temporarily pause execution after yielding a value and resume only when the program requests the next value.

To ensure comprehensive coverage, RIGHTTYPER initially enables these events across all code, guaranteeing that every executed function is sampled at least once. As it processes events, RIGHTTYPER selectively disables monitoring for the specific code segments that generated them, reducing overhead. At regular intervals, RIGHTTYPER interrupts execution to profile itself, updating a running counter that tracks the proportion of times found executing its own code. If this percentage falls below a predefined threshold—5% by default—RIGHTTYPER re-enables monitoring for all code locations. This mechanism allows continued sampling of previously executed code portions while maintaining low overhead.

Python coroutines, like generator functions, also yield values. However, unlike generators, coroutines can also receive values during execution. Since `sys.monitoring` does not provide built-in support for intercepting these received values, RIGHTTYPER instruments the code to capture them. It intercepts code loading, injects instrumentation into the abstract syntax tree (AST), and recompiles it before execution.

2.2 Object Typing and Sampling

As it processes monitoring events, RIGHTTYPER examines the objects passed into and returned from each function to infer their types. Given the numerous extensions introduced since type annotations were first added to Python [7], RIGHTTYPER adjusts how it types objects based on the target Python version to make the best use of the typing system.

Import Paths. Because most types—aside from a few built-ins—must be imported for use in annotations, RIGHTTYPER must determine the module and name under which each type is available; together, we refer to these here as the *import path*. Python exposes an object’s type via the `type()` function, whose `__module__` and `__qualname__` (“qualified” name indicating its path within the module) attributes indicate where they are defined. In many simple cases, these attributes are sufficient to generate an import path. However, in numerous other cases—particularly for types implemented in C, such as Python’s built-ins—`__module__` may contain only a package name, or `__qualname__` may in fact provide an unqualified name or even expose an internal name not accessible from Python code. Such cases do not necessarily imply that the type is unavailable: many are accessible under different names or from entirely different modules. For example, when called on a function, `type()` returns a type named `function` from the module `builtins`. However, in Python code, that type is actually available as `FunctionType` from the `types` module. To address these inconsistencies, RIGHTTYPER builds a map from type objects to import paths by scanning Python modules for type definitions.

Most Python types are accessible through multiple names. Such definitions can be intentional or incidental, as modules often re-export symbols defined in their submodules. While not strictly necessary for typing, RIGHTTYPER attempts to identify a canonical

import path for each type. RIGHTTYPER prioritizes paths explicitly included in a module’s public interface—namely, those listed in the module’s `__all__` attribute—over other definitions. Paths containing identifiers that begin with an underscore are deprioritized, as such names conventionally indicate internal or non-public use [37]. To choose among the remaining options, RIGHTTYPER applies heuristics informed by common Python practices. If a candidate path originates from a different package, it is treated as an alias and disregarded: packages must import a type before re-exporting it, causing such aliases to appear earlier in the scan. When two paths originate from the same package, RIGHTTYPER selects the shorter one (i.e., with fewer identifiers), as shorter aliases are often introduced to simplify access. Selecting paths in this manner is intended to produce more natural and consistent annotations across the codebase.

While this approach works well for most objects whose types are visible from Python, certain commonly used built-in objects, such as iterators, lack importable type definitions altogether. For these, RIGHTTYPER falls back on structural subtyping, using *protocols* such as `Iterator` [36].

Generic Types. If the type is a generic, RIGHTTYPER also attempts to discover its parameters. Instances of generic classes, defined in Python source code, contain an `__orig_class__` attribute that describes the class and the types taken by each parameter. However, certain standard library generic objects, such as containers, lack this attribute. To determine their type parameters, RIGHTTYPER samples their contents. Since most standard containers do not support efficient random sampling, where necessary, RIGHTTYPER uses an iterator-based approach to extract samples. To mitigate overhead, it also limits the number of elements to sample—by default, to 1,000.

Given the ubiquity of dictionaries in Python, RIGHTTYPER replaces the standard `dict` with a custom implementation that also maintains its elements in a list, enabling statistically correct sampling in $O(1)$ time complexity thanks to constant-time list access. This replacement is performed using the AST-based mechanism described in Section 2.1. However, since many existing Python programs depend on functionality that explicitly requires the standard `dict`, RIGHTTYPER provides this replacement as an optional feature.

When encountering an empty container, RIGHTTYPER specifies its element type as `Never`, provided the target Python version supports it. This special type allows static checkers to detect and flag code behavior that is inconsistent with an empty container.

Iterators pose special challenges for typing because their interface does not allow inspection of return values without destructively updating their state. For iterators implemented in Python, whose iteration method calls are visible through `sys.monitoring`, RIGHTTYPER begins with a partial type and updates it to a complete one once an iteration has been observed. Built-in iterators, however, do not trigger monitoring events. For these, RIGHTTYPER uses Python’s standard `gc` module to locate the underlying object(s) and infers types based on them, falling back to a generic `Iterator[Any]` when necessary.

Call Traces. RIGHTTYPER combines the types of a function’s arguments and return value into a *call trace*, from which it later generates annotations. By capturing all types involved in a call, call

```
def add(a, b):
    return a + b

...
add(10, 20)
add("foo", "bar")

def add(a: float|str, b: float|str) -> float|str:
    return a + b

error: Unsupported operand types for + ("float" and "str") [
    operator]
```

Figure 2: Naïve typing with unions: This example illustrates a function with interdependent argument and return types. A union of observed types (bottom) is overly permissive, incorrectly allowing mixed-type inputs that result in a type error. MonkeyType and PyAnnotate emit such unions (see Section 3.1.1).

traces enable RIGHTTYPER to identify type patterns and produce more precise annotations (see Section 2.3).

2.3 Filtering and Combining Call Traces

After the target program finishes executing, RIGHTTYPER uses the collected call traces to generate type annotations for each function. It first filters out rare traces—by default, keeping those that account for 80% of calls—then infers annotations from the remainder. When the type of an argument or return value varies across traces, the simplest annotation is a *union* of the observed types. For instance, if a function receives a string argument in one sample and an integer in another, it may be annotated as `f(a: str | int)`. While unions offer a flexible way to describe multiple observed types, they can also be overly permissive, potentially leading to errors or false negatives during static analysis. In particular, they may fail to capture essential relationships between arguments. For example, Figure 2 illustrates the function `add`, which returns the sum of two numbers or the concatenation of two strings. Naïvely annotating both arguments as `float | str` permits invalid combinations—such as one argument being a float and the other a string—that would result in either a static type error or, if undetected, a runtime error. In this example, because the arguments are used directly with the addition operator, mypy can detect the error, issuing the message shown in Figure 2.

For this reason, rather than immediately constructing unions, RIGHTTYPER first searches for patterns in the observed types. When it detects consistent variability across argument or return value types, RIGHTTYPER introduces a *type argument* to capture this shared variability, binding it to the concrete types encountered during execution. The center portion of Figure 3 illustrates the resulting annotation. The search operates recursively, allowing it to detect patterns nested within type arguments. For instance, a function that returns one of the elements of a list might be annotated as `f(l: list[T]) -> T`. Algorithm 1 describes the process in more detail. If the target Python version does not support type arguments, RIGHTTYPER instead defines a *type variable* to achieve the same effect, albeit in a more verbose manner, as the bottom section of Figure 3 shows.

When RIGHTTYPER cannot identify a consistent pattern across argument or return value types, it constructs a union of the observed

Algorithm 1 Call trace generalization: the algorithm examines the types observed at each *position* (argument or return value) in a collection of call traces, generating a type for each corresponding position in the function signature. It transposes the call traces to group types by position and operates recursively, enabling it to identify patterns in the arguments of generic types. The symbol \oplus denotes concatenation.

```
function GENERALIZE(traces)
    function REBUILD(pos_types)
        if all pos_types are specializations of a generic G then
            arg_positions ← TRANSPOSE(ARGUMENTS(pos_types))
            arguments ← [ ]
            for each p in arg_positions do
                arguments ← arguments  $\oplus$  REBUILD(p)
            return G[arguments]
        else if pos_types occurs more than once then
            T ← assign or retrieve typevar for pos_types
            return T
        else
            return UNION(pos_types)
    signature ← [ ]
    for each p in TRANSPOSE(traces) do
        signature ← signature  $\oplus$  REBUILD(p)
    return signature
```

```
def add[T: (float, str)](a: T, b: T) -> T:
    return a + b

from typing import TypeVar

rt_T1 = TypeVar("rt_T1", int, str)
def add(a: rt_T1, b: rt_T1) -> rt_T1:
    return a + b
```

Figure 3: RIGHTTYPER recognizes type patterns: recognizing that `add` (Figure 2) is consistently called with either strings or numbers, RIGHTTYPER annotates the function using either a type argument (if supported by the target Python version; top portion) or a generated type variable (`rt_T1`; bottom portion).

types at that position. Before forming the union, RIGHTTYPER attempts to simplify this set by replacing groups of individual types with a suitable common supertype. To ensure semantic correctness, it verifies that the chosen supertype defines all methods and attributes shared by the types it replaces. For numeric types, it follows Python’s “numeric tower”; for instance, if both `int` and `float` are observed, the type is simplified to `float` [36].

2.4 Annotating Numerical Array Shapes

Numerical array shape mismatches are a common source of errors in machine learning programs [9]. While numpy does not support array shape annotations at the time of writing, RIGHTTYPER can optionally include shape information using the format supported by jaxtyping, a package that enables runtime shape checking. For instance, a 2×3 numpy array would be annotated as `Float64[ndarray, "2 3"]`. When emitting such annotations for


```

class Shape:
    ...
    def scale(self, factor):
        ...

class Line(Shape):
    ...

line = Line()
...
line = line.scale(1.10)

class Shape:
    ...
    def scale(self: Line, factor: float):
        ...

error: The erased type of self "line.Line" is not a supertype of
its class "shape.Shape" [misc]

```

Figure 4: Inheritance complicates typing: in the top section, class `Line` inherits the `scale()` method from its parent class, `Shape`. When `scale()` is invoked on a `Line` object, its `self` parameter refers to that object—an instance of `Line`. However, annotating `self: Line` (bottom section) would be incorrect, as the mypy error demonstrates.

arrays, `RIGHTTYPER` also attempts to identify patterns in the shapes observed across traces, replacing them with variables.

2.5 Typing Methods

Methods require special handling for type annotation. One reason is that their arguments and return values may reference the current instance, which could be a subtype of the class where the method is defined. Consider the example in Figure 4: when `scale()` is invoked, its `self` parameter holds a reference to a `Line` object. Annotating it as `self: Line` would be incorrect, as demonstrated by the mypy error in the bottom section. Instead, `RIGHTTYPER` infers the type of the current instance and annotates such objects using the class that defines the method. For class instances, it similarly uses the defining class (`type[Shape]` in this example). When generating annotations for a Python version that supports it, `RIGHTTYPER` uses `Self` (or `type[Self]`) from the typing module, which explicitly refers to the current instance or class, respectively. While using `Self` improves readability, it is particularly useful for methods that return their instance [30].

Additional handling is also needed to avoid overly narrow annotations on method arguments. If a method overrides one from a superclass, using concrete types observed may improperly restrict the original type, violating the Liskov Substitution Principle (LSP) [15]. Consider the example in Figure 5, where `Line` overrides the method `draw` from `Shape`. Even though `Line.draw` is only ever called with `Screen` instances as the `medium` argument, annotating it as `medium: Screen` would be incorrect. Doing so would violate the LSP by excluding other valid media accepted by the base class.

To address this issue, `RIGHTTYPER` inspects parent classes to determine whether a method overrides one from a superclass. If so, it annotates the method’s arguments using the union of the types observed in the call traces and those specified in the overridden method’s signature. The bottom section of Figure 5 illustrates this approach.

```

class Shape:
    ...
    def draw(self: Self, medium: Medium) -> None:
        ...

class Line(Shape):
    ...
    def draw(self, medium):
        ...

line = Line()
...
line.draw(Screen())

class Line(Shape):
    ...
    def draw(self: Self, medium: Medium|Screen) -> None:
        ...

```

Figure 5: Method overriding complicates typing: in the top section, class `Line` overrides the `draw()` method from its parent class, `Shape`. Although `Line.draw()` is only ever called with `Screen` as its `medium`, the method must remain compatible with all argument types accepted by `Shape.draw()`. Accordingly, `RIGHTTYPER` types the parameter using the union of the observed type(s) and those accepted by the parent method (bottom section). Note: if `Screen` were a subtype of `Medium`, `RIGHTTYPER` would reduce the union to simply `Medium`.

To determine parent class annotations, `RIGHTTYPER` consults `typeshed` when necessary—a collection of *type stubs* that provide type annotations for Python’s standard library and many third-party packages [12]. Consulting `typeshed` is essential because the vast majority of the standard library, as well as numerous popular packages, lack inline type annotations of the kind shown in the examples.

2.6 Implementation

Our implementation of `RIGHTTYPER` comprises approximately 5,000 lines of Python code. While it relies on `sys.monitoring` and therefore requires Python 3.12 or later to run, it can emit annotations targeting Python versions 3.9 through 3.13. It uses the standard `ast` module for AST-based instrumentation, but leverages the `libcst` concrete syntax tree library to parse and modify Python source code [32]. It uses `typeshed_client` to access the `typeshed` type stub repository [12, 49].

3 Evaluation

Our evaluation investigates the following questions:

- RQ1:** Do type annotations generated by `RIGHTTYPER` improve upon those from prior approaches? (Section 3.1)
- RQ2:** How does `RIGHTTYPER`’s runtime overhead compare to that of other dynamic approaches? (Section 3.2)

3.1 [RQ1] Type Comparisons

Evaluating type inference is a challenging task. One of the difficulties lies in establishing a reliable ground truth. Some prior work derive it from developer-written annotations, sometimes supplemented with types inferred by static analyzers such as `pytype` or `pyre` [2, 19]. However, many projects contain type errors, as type

Table 1: Litmus tests used to investigate RQ1: these tests reveal shortcomings of type inference tools. **R** = causes runtime error, **T** = causes spurious type-checking error, **F** = fails to type arg. or return, **M** = masks existing error, **P** = potentially masks errors, \checkmark = none apply. \dagger : using Python code manually created from tool output.

Test	Description	MonkeyType	PyAnnotate	pytype	QuAC \dagger	TypeT5 \dagger	RIGHTTYPER
Interdependent Types	Typing a function where arguments and/or return value depend on each other	T	T	F	P	T	\checkmark
Unnamed Types	Annotating elements lacking user-visible type names and optional values	F	R	F,P	P	T,P	\checkmark
Inheritance and Overriding	Annotating overridden methods	T	R,T,F	R,T,F	P	T,P	\checkmark
Edge case error	An edge case error, unexposed by testing	F,M	T	F,M	\checkmark	\checkmark	\checkmark

checking is not yet a standard part of development workflows, and relying on developer-written annotations limits evaluation to already-typed codebases [6]. Moreover, types inferred by static tools are only as reliable as the tools themselves. As discussed, limited by dynamic language features, these tools may conservatively over-approximate or fail to infer types (§ 4). Available evaluation datasets may omit type parameters—e.g., reducing `list[int]` to `list`—which significantly reduces their utility for assessing inference precision [38].

Another difficulty lies in comparing inferred types. A program element containing a `list[int]`, for example, can be validly annotated in multiple ways. In addition to trivial aliases like `List[int]` and less precise forms like `list`, it may also be appropriate to use protocol-based types such as `Collection`, `Sequence`, or `Iterable`. An effective comparison metric must be sufficiently nuanced to recognize when each of these types is appropriate.

Some prior work adopts a “correctness modulo type checker” approach to evaluation, measuring the quality of inferred types on whether they pass static type checking without errors. [2, 24, 45]. However, doing so reinforces the often unwarranted assumption that the underlying code is correct. As some examples in this section illustrate, annotations can effectively silence type checkers, inadvertently concealing bugs.

Type coverage—the proportion of code elements for which types are inferred—is another common metric in prior work. However, for a dynamic system like RIGHTTYPER, type coverage is largely determined by the code coverage achieved by the driver (e.g., the test suite). As a result, it offers little insight into the system’s effectiveness and is not a meaningful metric in this context.

To sidestep these challenges, and drawing inspiration from prior work on memory consistency models, we develop a set of “litmus tests” to guide our investigation [4, 34]. Each test is a small code example designed to expose specific difficulties in type inference—such as interdependencies between argument and return types, method overriding through inheritance, or the use of unnamed types. We use these examples to compare the type annotations produced by MonkeyType, PyAnnotate, pytype, TypeT5, QuAC, and RIGHTTYPER. Table 1 provides an overview; for brevity, we omit some individual results.

Because QuAC and TypeT5 produce type annotations in textual formats rather than as Python code, we manually translate their outputs into Python for evaluation and presentation purposes. This manual conversion implicitly grants them an unfair advantage, as

it sidesteps the runtime errors that other tools encounter when generating or inserting annotations automatically.

3.1.1 Interdependent Types. We first examine typing a function with interdependent argument and return types, already described in Section 2.3. When invoked with two numbers, the function returns their sum; when called with two strings, it returns their concatenation. Calling the function with one number and one string argument yields a runtime error, as the operation is invalid.

As Figure 6b shows, MonkeyType and PyAnnotate annotate the function with a union, leading mypy to emit errors that warn about the potential runtime error. pytype fails to type the function entirely, potentially excluding it from static type checking. TypeT5 types both arguments as `int`, resulting in type errors when the function is called with strings (Figure 6c). QuAC instead uses the `SupportsAdd` protocol, indicating that the object supports the addition operation (Figure 6d). While its annotation avoids immediate mypy errors, it also prevents detection of mismatched argument types. In contrast, RIGHTTYPER (Figure 6e) annotates the arguments and return value with a type variable, allowing mypy to detect mismatches when they occur.

3.1.2 Unnamed types. Python libraries often use objects—such as container iterators or key, value, and item view objects—whose types lack user-visible names. Figure 7a illustrates an example in which a dictionary’s value view is passed to a function. Additionally, the function’s second parameter is optional, allowing us to assess how different type assistants handle both unnamed objects and optional parameters.

Figures 7b to 7d show the types inferred by the various systems: MonkeyType is unable to infer any types for the function; PyAnnotate and TypeT5 infer types, but the resulting annotations lead to errors, either at runtime or during type checking. While pytype and QuAC’s types do not cause errors, they are overly broad, reducing the recall and thus the effectiveness of type checking. In contrast, RIGHTTYPER’s annotations (Figure 7e) are precise and do not introduce errors.

3.1.3 Inheritance and Overriding. We next look into typing an overridden method. In Figure 8a, class `Line` implements `Shape`, and its `draw` method extends the original interface by also permitting instances of `NonStandardMedium` in its `medium` argument.

Litmus test: Interdependent types

```
def add(a, b):  
    return a + b  
  
add(10, 20)  
add("foo", "bar")
```

(a) **Litmus Test:** This example shows a function that, when invoked with two numbers, returns their sum, but when called with two strings, returns their concatenation.

MonkeyType, PyAnnotate

```
def add(a: Union[int, str], b: Union[int, str]) -> Union[int, str]:  
    return a + b  
  
error: Unsupported operand types for + ("int" and "str") [  
    operator]  
note: Both left and right operands are unions
```

(b) **MonkeyType and PyAnnotate cause errors:** MonkeyType and PyAnnotate use unions, which lead to mpy errors.

TypeT5

```
def add(a: int, b: int) -> int:  
    return a + b  
  
error: Argument 1 to "add" has incompatible type "str";  
    expected "int" [arg-type]  
error: Argument 2 to "add" has incompatible type "str";  
    expected "int" [arg-type]
```

(c) **TypeT5 causes errors:** TypeT5 types the arguments as int, causing mpy errors when the function is used with strings.

QuAC

```
from _typeshed import SupportsAdd  
def add(a: SupportsAdd, b: SupportsAdd) -> SupportsAdd:  
    return a + b  
  
add(10, "bar")  
  
TypeError: unsupported operand type(s) for +: 'int' and 'str'
```

(d) **QuAC silences mypy:** QuAC types “add” using the SupportsAdd protocol, which prevents mypy from detecting a mismatch error (middle section). Executing the program yields a Python runtime error.

RIGHTTYPER

```
def add[T1: (int, str)](a: T1, b: T1) -> T1:  
    return a + b
```

(e) **RIGHTTYPER recognizes the type pattern:** RIGHTTYPER uses a type variable, both avoiding immediate errors and enabling mypy to find mismatch errors.

Figure 6: Litmus test for arguments and return value that depend on each other (§ 3.1.1).

Figures 8b, 8c, 9a and 9b present annotations generated by the various baseline systems, along with errors these annotations introduce; only RIGHTTYPER (Figure 9c) correctly infers all argument and return types.

Litmus test: Unnamed types

```
def f(a, b=None):  
    m = max(a)  
    if b is not None and b > m:  
        return [m, b]  
    return [m]  
  
d = {'a': 1, 'b': 2}  
  
f(d.values())  
f(d.values(), 20)
```

(a) **Litmus test:** In this example, the function is called with a dictionary’s value view object. These objects, internally named dict_values in Python, do not have a user-visible type name.

PyAnnotate

```
def f(a: dict_values, b: Optional[int] = None) -> List[int]:  
  
NameError: name 'dict_values' is not defined
```

MonkeyType

```
def f(a, b=None): ...
```

(b) **PyAnnotate causes runtime error:** PyAnnotate correctly annotates b as an optional integer but uses the internal name dict_values for a, resulting in a runtime error. MonkeyType, in contrast, rejects the internal name as an invalid type and discards the corresponding trace—leaving the function entirely unannotated.

TypeT5

```
def f(a: list, b: int=None) -> list: ...  
  
error: Incompatible default for argument "b" (default has type  
    "None", argument has type "int") [assignment]  
error: Argument 1 to "f" has incompatible type "dict_values[str  
    , int]"; expected "list[Any]" [arg-type]
```

(c) **TypeT5 types incorrectly:** TypeT5 mistakes argument a for a list and fails to recognize that b is optional, resulting in errors. pytype does not cause errors, but fails to annotate either function argument. It correctly identifies the base return value as list but omits the element type, weakening the recall of type checking.

pytype

```
def f(a, b=None) -> list: ...
```

QuAC

```
def f(a: Iterable, b) -> list: ...
```

(d) **Partial annotations:** pytype and QuAC do not cause errors, but fail to fully annotate the function. While QuAC correctly types argument a using the Iterable protocol, pytype is unable to annotate either argument. Like also TypeT5 above, both omit the list element type, weakening type checking.

RIGHTTYPER

```
def f(a: ValuesView[int], b: int|None=None) -> list[int]: ...
```

(e) **RIGHTTYPER types correctly:** RIGHTTYPER types argument a using the ValuesView protocol, indicating its element type, correctly identifies b as an optional integer argument, as well as fully types the function’s return value.

Figure 7: Litmus test for unnamed types (§ 3.1.2).

Litmus test: Inheritance and Overriding

```
from abc import ABC, abstractmethod
from typing import Self

class Medium(ABC):
    @abstractmethod
    def plot(self, x_points, y_points): ...

class Shape(ABC):
    @abstractmethod
    def draw(self: Self, medium: Medium) -> None: ...

    @abstractmethod
    def zoom(self: Self, factor: float) -> Self: ...

class NonStandardMedium:
    def plot(self, x_points, y_points, shrinkage): ...

class Line(Shape):
    def __init__(self, x, y):
        super().__init__()
        self.x = x
        self.y = y

    def zoom(self, factor):
        self.x = (self.x[0] * factor, self.x[1] * factor)
        self.y = (self.y[0] * factor, self.y[1] * factor)
        return self

    def draw(self, medium):
        if isinstance(medium, NonStandardMedium):
            medium.plot(self.x, self.y, .10)
        else:
            medium.plot(self.x, self.y)

line = Line((0., 10.), (0., 10.))
line.zoom(1.1)
line.draw(NonStandardMedium())
```

(a) **Litmus test:** Class `Line` implements `Shape` and extends the `draw` method, also permitting `NonStandardMedium` as its `medium` argument.

MonkeyType

```
class Line:
    def __init__(self, x: Tuple[float, float],
                 y: Tuple[float, float]): ...
    def zoom(self, factor: float) -> "Line": ...
    def draw(self, medium: NonStandardMedium): ...

error: Argument 1 of "draw" is incompatible with supertype "Shape"; supertype defines the argument type as "Medium" [override]
note: This violates the Liskov substitution principle
```

(b) **MonkeyType errors:** `MonkeyType` omits `Medium` from `medium`'s annotation, causing errors.

pytype

```
_TLine = TypeVar('_TLine', bound=Line)
class Line:
    def __init__(self, x, y) -> None: ...
    def zoom(self: _TLine, factor) -> _TLine: ...
    def draw(self, medium) -> None: ...

error: Name "Line" is used before definition [used-before-def]
```

(c) **pytype errors:** `pytype` fails to infer types for all arguments except `self`, and even that leads to an error.

Figure 8: Litmus test for inheritance and overriding (§ 3.1.3).

Litmus test: Inheritance and Overriding (cont'd.)

QuAC

```
class Line:
    def __init__(self, x, y): ...
    def zoom(self, factor) -> "Line": ...
    def draw(self, medium: NonStandardMedium): ...

error: Argument 1 of "draw" is incompatible with supertype "Shape"; supertype defines the argument type as "Medium" [override]
note: This violates the Liskov substitution principle
```

(a) **QuAC errors:** `QuAC` is unable to type the arguments for `__init__` and `zoom`, and omits `Medium` from `draw`'s `medium` argument, causing errors.

TypeT5

```
class Line:
    def __init__(self, x: Self, y: Self) -> None: ...
    def zoom(self, factor: float) -> "Line": ...
    def draw(self, medium: Medium) -> None: ...

error: Argument 1 to "Line" has incompatible type "tuple[float, float]"; expected "Line" [arg-type]
error: Argument 2 to "Line" has incompatible type "tuple[float, float]"; expected "Line" [arg-type]
error: Argument 1 to "draw" of "Line" has incompatible type "NonStandardMedium"; expected "Medium" [arg-type]
```

(b) **TypeT5 errors:** Applied to the code in Figure 8a, `TypeT5` incorrectly types `x` and `y` as `Self` and omits `NonStandardMedium` from `draw`'s `medium` argument, causing errors.

RIGHTTYPER

```
class Line:
    def __init__(self, x: tuple[float, float],
                 y: tuple[float, float]) -> None: ...
    def zoom(self, factor: float) -> Self: ...
    def draw(self, self,
             medium: Medium|NonStandardMedium) -> None: ...
```

(c) **RIGHTTYPER:** `RIGHTTYPER` correctly infers all argument and return types. By including both the observed `NonStandardMedium` and the inherited `Medium` in the annotation of `medium`, it abides by the Liskov substitution principle.

Figure 9: Additional results on inheritance and overriding (§ 3.1.3).

3.1.4 Edge case error. Figure 10a presents a function containing a rarely executed bug, along with a unit test that exercises it. The example simulates a bug: failing to handle an edge case in the input domain. Due to a mistake in the range check, passing the value 9 causes the function to implicitly return `None` instead of the expected `bool`.

Figures 10b and 10c show the type annotations inferred by `pytype` and `RIGHTTYPER`, along with the outcome of static checking using `mypy`. `pytype`, using static analysis, infers an `Optional` return type due to the possibility of returning `None`. This, in turn, causes `mypy` to overlook the missing return path. In contrast, `RIGHTTYPER` infers the return type observed during test execution, enabling `mypy` to surface the error.

Litmus test: Edge case error

Original Code

```
def is_value_ok(value):
    if value >= 10:
        return True
    elif value < 9:
        return False

is_value_ok(100)
is_value_ok(11)
is_value_ok(10)
is_value_ok(0)
is_value_ok(-1)
```

(a) **Litmus test:** The function in this example fails to return a boolean when passed the value 9, instead implicitly returning None. The example exercises the function with various inputs, simulating a test.

pytype

```
def is_value_ok(value) -> Optional[bool]: ...

Success: no issues found
```

(b) **pytype silences mypy:** PyType, analyzing the source statically and detecting an execution path with no return value, annotates the function with an optional return type. As a result, mypy does not report the missing return as an error.

RIGHTTYPER

```
def is_value_ok(value: int) -> bool: ...

error: Missing return statement [return]
```

(c) **RIGHTTYPER helps reveal bug:** RIGHTTYPER produces an annotation based on the return types observed during testing, allowing mypy to detect the error.

Figure 10: Litmus test for edge case error (§ 3.1.4).

[RQ1] Summary: RIGHTTYPER’s annotations substantially improve upon the quality of those of previous approaches.

3.2 [RQ2] Runtime Performance

To evaluate RIGHTTYPER’s runtime performance, we measure the time it requires to collect type information while running a suite of benchmark applications, normalized against each application’s baseline runtime. We perform the same measurements for two existing dynamic type inference tools: MonkeyType and PyAnnotate. All tools are executed with their default configurations.

We draw benchmarks from three sources:

- To represent real-world Python applications, we run the black code formatter on the large scikit-learn package, and run a Sudoku solver [21];
- We run the test suites of three widely used Python packages: black, rich, and tornado;
- We include benchmarks drawn from Python’s benchmark suite [27].

We run all experiments using Python 3.12.8 on a 10-core 3.7GHz Core i9 system with 64GB of RAM and SSD storage, running Linux

6.5.6. With the system otherwise idle, we run each experiment three times and report the median execution time.

Across all benchmarks, RIGHTTYPER consistently outperforms MonkeyType and PyAnnotate in terms of overhead. RIGHTTYPER incurs a maximum overhead of 0.9× (mean: 0.3×), compared to PyAnnotate’s 7.7× (mean: 2.9×) and MonkeyType’s 273× (mean: 13×). We observe that RIGHTTYPER incurs a mean runtime overhead of 0.3×, which substantially exceeds its target of 0.05× (see Section 2.1). Profiling reveals that this overhead stems from its call trace assembly mechanism, which requires sys.monitoring events to remain enabled until the trace is complete. Because these events are tied to code regions rather than specific invocations, RIGHTTYPER receives additional, unrelated events whenever a traced function is re-entered, such as during recursion. While we aim to address this inefficiency in future work, we note that RIGHTTYPER still substantially outperforms both PyAnnotate and MonkeyType.

[RQ2] Summary: RIGHTTYPER is substantially faster than prior dynamic approaches.

4 Related Work

There has been substantial research on type inference for dynamically typed languages, particularly JavaScript, Ruby, and Python. This section concentrates primarily on prior work targeting Python.

Static type inference. Static approaches analyze program structure without executing it, relying on techniques such as abstract interpretation, data-flow analysis, and type constraint resolution to infer types [8, 16, 31, 39, 41, 45].

Static type *checkers* such as mypy, pyright, pyre, and pytype must also perform type inference to enable type checking, though their inference may be limited, primarily focused on how types from existing annotations propagate through the code or are affected by operations [13, 14, 33, 48]. Among these tools, only pytype and pyre support adding type annotations.

While typically sound by design, static approaches are constrained by dynamic language features such as runtime type changes, reflection, monkey patching, and dynamic code generation, and often support only a subset of their target languages [8, 11, 22, 28]. Dynamic features can also lead static systems to conservatively over-approximate, inferring overly permissive types [3, 26]. In contrast to static approaches, RIGHTTYPER produces sound, precise annotations and is able to annotate all executed code paths.

AI-based type inference. These approaches estimate the likely types of program elements from features extracted from the source code, often including natural language elements. Xu et al. apply a probabilistic graphical model to predict variable types based on *hints* such as attribute accesses and variable names. [46]. TypeWriter applies a neural network to both code and natural language information and verifies the predictions using a type checker [26]. Typilus and Type4Py leverage deep similarity learning to avoid using a fixed type vocabulary [2, 19]. PYInfer uses a neural network to infer variable types from the surrounding source code context [5]. HiTyper combines static and machine learning inference by only applying the machine learning model where its type constraints

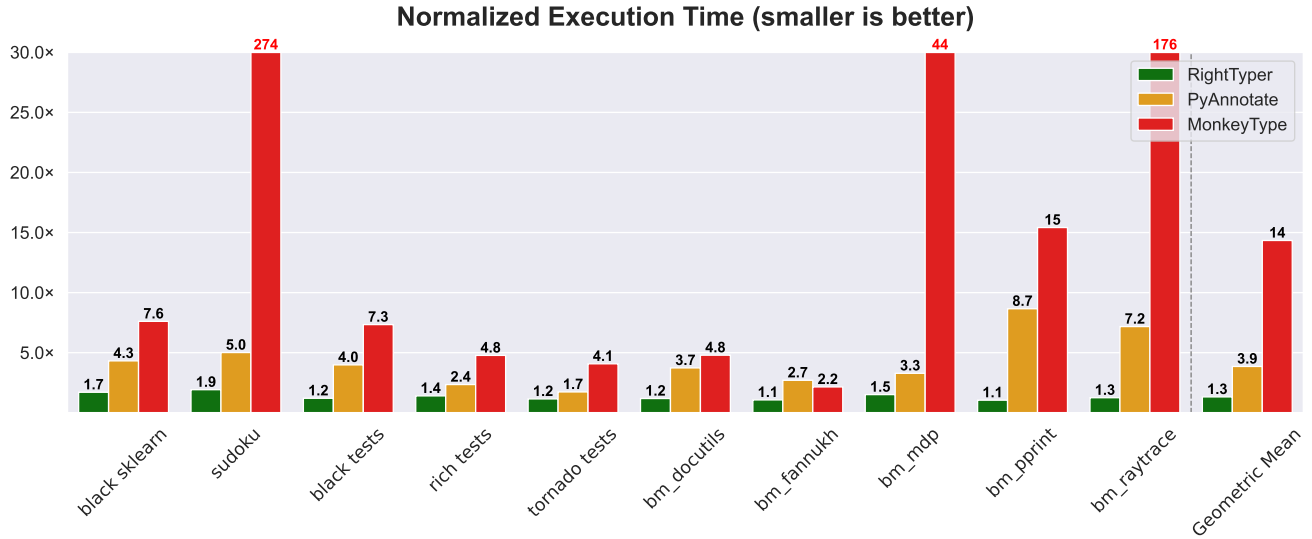


Figure 11: [RQ2] RIGHTTYPER is much faster than existing dynamic tools: across our benchmark suite, while PyAnnotate and MonkeyType incur high overheads (of 2.9× and 13×), RIGHTTYPER’s overhead remains at 0.3×.

resolution gets stuck [22]. Abdelaziz et al. combine mining docstrings and approximate duck subtyping to infer function return types [1, 18]. TypeT5 fine-tunes CodeT5 for type inference [42, 43]. DLInfer extracts program slices for variables and applies a sequence model to predict function argument types [47]. TypeGen prompts ChatGPT with chain-of-thought instructions to infer types [23, 44]. TIGER ranks generative model typing candidates with a similarity model [40].

While better able to cope with Python’s dynamic nature than purely static approaches, AI-based methods are inherently unsound: their predicted types are not guaranteed to reflect actual program behavior. By combining AI-based and static inference, approaches such as TypeWriter, HiTyper, and Typilus can reject incorrect types—but risk leaving more of the code unannotated [22, 26]. Furthermore, these approaches often only support a limited, small type vocabulary and perform poorly on rare types, making them less effective when encountering custom or evolving types [19]. Unlike AI-based approaches, RIGHTTYPER produces sound annotations that reflect actual runtime behavior.

Dynamic type inference. Dynamic approaches track types encountered during program execution. This strategy is employed by tools such as MonkeyType (developed by Instagram), PyAnnotate (developed by Python author Guido v. Rossum at Dropbox), as well as RIGHTTYPER [17, 35]. Both MonkeyType and PyAnnotate rely on Python’s runtime type attributes to identify observed types. However, these attributes are not always reliable (cf. Section 3.1). MonkeyType attempts to mitigate this by verifying type names before emitting annotations, but this pruning can result in code being left unannotated or only partially annotated. PyAnnotate performs no such validation, and as a result, may emit annotations that include unknown names, causing runtime errors. RIGHTTYPER instead builds a mapping from type objects to their names, ensuring

that the annotations are valid and, whenever possible, use public names.

By default, MonkeyType logs every function call to a SQLite database, resulting in severe performance degradation and, in some cases, causing execution to abort due to resource exhaustion [10]. To mitigate this overhead, MonkeyType can be configured to only sample one in every N function calls. However, this coarse-grained sampling risks leaving many functions unannotated, particularly during short executions. PyAnnotate also uses call sampling but follows a fixed scheme, recording only a limited number of calls per function, which favors calls early in the program’s execution. Moreover, neither MonkeyType nor PyAnnotate applies sampling to container types such as lists, sets, or dictionaries. MonkeyType exhaustively inspects every element, regardless of the size of the data structures, further increasing overhead. PyAnnotate limits overhead by inspecting the first four elements, producing biased samples. In contrast, RIGHTTYPER guarantees that every function that executes is recorded at least once. Additionally, it samples data structure contents uniformly at random, enabling more representative type inference.

5 Conclusion

This paper introduces RIGHTTYPER, a novel dynamic type inference method for Python that produces precise type annotations based on observed program behavior. By leveraging principled, self-profiled sampling, combined with statistical filtering and robust type resolution, RIGHTTYPER enables a novel interpretation of type checking as anomaly detection: rare or inconsistent behaviors are excluded from annotations, allowing static checkers to flag them for developer review. Its combination of precision, efficiency, and diagnostic power makes RIGHTTYPER a practical and effective approach for annotating Python programs.

References

- [1] Ibrahim Abdelaziz, Julian Dolby, and Kavitha Srinivas. 2022. Large Scale Generation of Labeled Type Data for Python. arXiv:2201.12242 [cs.PL] <https://arxiv.org/abs/2201.12242>
- [2] Miltiadis Allamanis, Earl T. Barr, Soline Ducouso, and Zheng Gao. 2020. Typilus: neural type hints. In *Proceedings of the 41st ACM SIGPLAN Conference on Programming Language Design and Implementation* (London, UK) (PLDI 2020). Association for Computing Machinery, New York, NY, USA, 91–105. <https://doi.org/10.1145/3385412.3385997>
- [3] Jong-hoon (David) An, Avik Chaudhuri, Jeffrey S. Foster, and Michael Hicks. 2011. Dynamic inference of static types for ruby. *SIGPLAN Not.* 46, 1 (Jan. 2011), 459–472. <https://doi.org/10.1145/1925844.1926437>
- [4] James Bornholt and Emina Torlak. 2017. Synthesizing memory models from framework sketches and Litmus tests. In *Proceedings of the 38th ACM SIGPLAN Conference on Programming Language Design and Implementation* (Barcelona, Spain) (PLDI 2017). Association for Computing Machinery, New York, NY, USA, 467–481. <https://doi.org/10.1145/3062341.3062353>
- [5] Siwei Cui, Gang Zhao, Zeyu Dai, Luochao Wang, Ruihong Huang, and Jeff Huang. 2021. PYInfer: Deep Learning Semantic Type Inference for Python Variables. *CoRR* abs/2106.14316 (2021). arXiv:2106.14316 <https://arxiv.org/abs/2106.14316>
- [6] Luca Di Grazia and Michael Pradel. 2022. The evolution of type annotations in python: an empirical study. In *Proceedings of the 30th ACM Joint European Software Engineering Conference and Symposium on the Foundations of Software Engineering* (Singapore, Singapore) (ESEC/FSE 2022). Association for Computing Machinery, New York, NY, USA, 209–220. <https://doi.org/10.1145/3540250.3549114>
- [7] Python Software Foundation. 2025. Typing PEPs. <https://peps.python.org/topic/typing/>. Accessed: 2025-05-14.
- [8] Mostafa Hassan, Caterina Urban, Marco Eilers, and Peter Müller. 2018. MaxSMT-Based Type Inference for Python 3. In *Computer Aided Verification - 30th International Conference, CAV 2018, Held as Part of the Federated Logic Conference, FloC 2018, Oxford, UK, July 14-17, 2018, Proceedings, Part II (Lecture Notes in Computer Science, Vol. 10982)*, Hana Chockler and Georg Weissenbacher (Eds.). Springer, 12–19. https://doi.org/10.1007/978-3-319-96142-2_2
- [9] Momoko Hattori, Naoki Kobayashi, and Ryosuke Sato. 2023. Gradual Tensor Shape Checking. arXiv:2203.08402 [cs.PL] <https://arxiv.org/abs/2203.08402>
- [10] Richard D Hipp. 2000. *SQLite*. <https://www.sqlite.org/>. Accessed: 2025-05-27.
- [11] Simon Holm Jensen, Anders Möller, and Peter Thiemann. 2009. Type Analysis for JavaScript. In *Proceedings of the 16th International Symposium on Static Analysis* (Los Angeles, CA) (SAS '09). Springer-Verlag, Berlin, Heidelberg, 238–255. https://doi.org/10.1007/978-3-642-03237-0_17
- [12] Matthias Kramm Jukka Lehtosalo et al. 2015. *Typeshed, a collection of library stubs for Python, with static types*. <https://github.com/python/typeshed> Accessed: 2025-05-27.
- [13] Matthias Kramm et al. 2015. *pytype*. Google. <https://github.com/google/pytype> Accessed: 2025-05-13.
- [14] Jukka Lehtosalo et al. 2012. *mypy: A static type checker for Python*. (2012). <https://github.com/python/mypy> Accessed: 2025-05-13.
- [15] Barbara Liskov. 1987. Keynote address - data abstraction and hierarchy. In *Addendum to the Proceedings on Object-Oriented Programming Systems, Languages and Applications (Addendum)* (Orlando, Florida, USA) (OOPSLA '87). Association for Computing Machinery, New York, NY, USA, 17–34. <https://doi.org/10.1145/62138.62141>
- [16] Eva Maia, Nelma Moreira, and Rogério Reis. 2012. A Static Type Inference for Python. In *Proceedings of the Workshop on Dynamic Languages and Applications (DYLA)*. ACM, 1–5.
- [17] Carl Meyer et al. 2017. *MonkeyType*. Instagram. <https://github.com/Instagram/MonkeyType> Accessed: 2025-05-13.
- [18] Nevena Milojkovic, Mohamadd Ghafari, and Oscar Nierstrasz. 2017. It's Duck (Typing) Season!. In *2017 IEEE/ACM 25th International Conference on Program Comprehension (ICPC)*. 312–315. <https://doi.org/10.1109/ICPC.2017.10>
- [19] Amir M. Mir, Evaldas Latoškinas, Sebastian Proksch, and Georgios Gousios. 2022. Type4Py: practical deep similarity learning-based type inference for Python. In *Proceedings of the 44th International Conference on Software Engineering* (Pittsburgh, Pennsylvania) (ICSE '22). Association for Computing Machinery, New York, NY, USA, 2241–2252. <https://doi.org/10.1145/3510003.3510124>
- [20] Anders Möller and Michael I. Schwartzbach. 2020. *Static Program Analysis*. <https://cs.au.dk/~amoeller/spa/>
- [21] Peter Norvig. 2006. *Solving Every Sudoku Puzzle*. <https://www.sqlite.org/> Accessed: 2025-06-13.
- [22] Yun Peng, Cuiyun Gao, Zongjie Li, Bowei Gao, David Lo, Qirun Zhang, and Michael Lyu. 2022. Static inference meets deep learning: a hybrid type inference approach for python. In *Proceedings of the 44th International Conference on Software Engineering* (Pittsburgh, Pennsylvania) (ICSE '22). Association for Computing Machinery, New York, NY, USA, 2019–2030. <https://doi.org/10.1145/3510003.3510038>
- [23] Yun Peng, Chaozheng Wang, Wenxuan Wang, Cuiyun Gao, and Michael R. Lyu. 2024. Generative Type Inference for Python. In *Proceedings of the 38th IEEE/ACM International Conference on Automated Software Engineering* (Echternach, Luxembourg) (ASE '23). IEEE Press, 988–999. <https://doi.org/10.1109/ASE56229.2023.00031>
- [24] Luna Phipps-Costin, Carolyn Jane Anderson, Michael Greenberg, and Arjun Guha. 2021. Solver-based gradual type migration. *Proc. ACM Program. Lang.* 5, OOPSLA, Article 111 (Oct. 2021), 27 pages. <https://doi.org/10.1145/3485488>
- [25] Juan Altmayer Pizzorno and Emery Berger. 2025. *RightTyper*. <https://github.com/righttyper/righttyper> Accessed: 2025-05-14.
- [26] Michael Pradel, Georgios Gousios, Jason Liu, and Satish Chandra. 2020. Type-Writer: neural type prediction with search-based validation. In *Proceedings of the 28th ACM Joint Meeting on European Software Engineering Conference and Symposium on the Foundations of Software Engineering* (Virtual Event, USA) (ESEC/FSE 2020). Association for Computing Machinery, New York, NY, USA, 209–220. <https://doi.org/10.1145/3368089.3409715>
- [27] Python Performance Team. 2008. *pyperformance*. <https://github.com/python/pyperformance> Accessed: 2025-06-13.
- [28] Gregor Richards, Sylvain Lebesne, Brian Burg, and Jan Vitek. 2010. An analysis of the dynamic behavior of JavaScript programs. In *Proceedings of the 31st ACM SIGPLAN Conference on Programming Language Design and Implementation* (Toronto, Ontario, Canada) (PLDI '10). Association for Computing Machinery, New York, NY, USA, 1–12. <https://doi.org/10.1145/1806596.1806598>
- [29] Mark Shannon. 2021. PEP 669 – Low Impact Monitoring for CPython. PEP 669. <https://peps.python.org/pep-0669/>
- [30] Pradeep Kumar Srinivasan and James Hilton-Balfe. 2021. PEP 673 – Self Type. PEP 673. <https://peps.python.org/pep-0673/>
- [31] Ke Sun, Yifan Zhao, Dan Hao, and Lu Zhang. 2023. Static Type Recommendation for Python. In *Proceedings of the 37th IEEE/ACM International Conference on Automated Software Engineering* (Rochester, MI, USA) (ASE '22). Association for Computing Machinery, New York, NY, USA, Article 98, 13 pages. <https://doi.org/10.1145/3551349.3561150>
- [32] Jennifer Taylor, Benjamin Woodruff, et al. 2019. *LibCST*. Instagram. <https://github.com/Instagram/LibCST.git> Accessed: 2025-05-27.
- [33] Eric Traut et al. 2019. *pyright: A Static Type Checker for Python*. Microsoft. <https://github.com/microsoft/pyright> Accessed: 2025-05-13.
- [34] Caroline Trippel, Daniel Lustig, and Margaret Martonosi. 2018. CheckMate: Automated Synthesis of Hardware Exploits and Security Litmus Tests. In *2018 51st Annual IEEE/ACM International Symposium on Microarchitecture (MICRO)*. 947–960. <https://doi.org/10.1109/MICRO.2018.00081>
- [35] Guido van Rossum et al. 2017. *PyAnnotate*. Dropbox. <https://github.com/dropbox/pyannotate> Accessed: 2025-05-13.
- [36] Guido van Rossum, Jukka Lehtosalo, and Łukasz Langa. 2014. *Type Hints*. PEP 484. <https://peps.python.org/pep-0484/>
- [37] Guido van Rossum, Barry Warsaw, and Alyssa Coghlan. 2001. *Style Guide for Python Code*. PEP 8. <https://peps.python.org/pep-0008/>
- [38] Ashwin Prasad Shrivastava Venkatesh, Samkuttu Sabu, Jiawei Wang, Amir M. Mir, Li Li, and Eric Bodden. 2024. TypeEvalPy: A Micro-benchmarking Framework for Python Type Inference Tools. arXiv:2312.16882 [cs.SE] <https://arxiv.org/abs/2312.16882>
- [39] Michael M. Vitousek, Andrew M. Kent, Jeremy G. Siek, and Jim Baker. 2014. Design and evaluation of gradual typing for python. In *Proceedings of the 10th ACM Symposium on Dynamic Languages* (Portland, Oregon, USA) (DLS '14). Association for Computing Machinery, New York, NY, USA, 45–56. <https://doi.org/10.1145/2661088.2661101>
- [40] Chong Wang, Jian Zhang, Yiling Lou, Mingwei Liu, Weisong Sun, Yang Liu, and Xin Peng. 2025. TIGER: A Generating-Then-Ranking Framework for Practical Python Type Inference. In *2025 IEEE/ACM 47th International Conference on Software Engineering (ICSE)*. IEEE Computer Society, Los Alamitos, CA, USA, 229–241. <https://doi.org/10.1109/ICSE55347.2025.00019>
- [41] Yin Wang et al. 2013. *PySonar2*. <https://github.com/yinwang0/pysonar2> Accessed: 2025-05-14.
- [42] Yue Wang, Weishi Wang, Shafiq Joty, and Steven C. H. Hoi. 2021. CodeT5: Identifier-aware Unified Pre-trained Encoder-Decoder Models for Code Understanding and Generation. arXiv:2109.00859 [cs.CL] <https://arxiv.org/abs/2109.00859>
- [43] Jiayi Wei, Greg Durrett, and Isil Dillig. 2023. TypeT5: Seq2seq Type Inference using Static Analysis. <https://doi.org/10.48550/arXiv.2303.09564> arXiv:2303.09564 [cs.SE]
- [44] Jason Wei, Xuezhi Wang, Dale Schuurmans, Maarten Bosma, Brian Ichter, Fei Xia, Ed H. Chi, Quoc V. Le, and Denny Zhou. 2022. Chain-of-thought prompting elicits reasoning in large language models. In *Proceedings of the 36th International Conference on Neural Information Processing Systems* (New Orleans, LA, USA) (NIPS '22). Curran Associates Inc., Red Hook, NY, USA, Article 1800, 14 pages.
- [45] Jifeng Wu and Caroline Lemieux. 2024. QuAC: Quick Attribute-Centric Type Inference for Python. *Proc. ACM Program. Lang.* 8, OOPSLA2, Article 343 (October 2024), 30 pages. <https://doi.org/10.1145/3689783>
- [46] Zhaogui Xu, Xiangyu Zhang, Lin Chen, Kexin Pei, and Baowen Xu. 2016. Python probabilistic type inference with natural language support. In *Proceedings of the 2016 24th ACM SIGSOFT International Symposium on Foundations of Software*

- Engineering* (Seattle, WA, USA) (*FSE 2016*). Association for Computing Machinery, New York, NY, USA, 607–618. <https://doi.org/10.1145/2950290.2950343>
- [47] Yanyan Yan, Yang Feng, Hongcheng Fan, and Baowen Xu. 2023. DLInfer: Deep Learning with Static Slicing for Python Type Inference. In *Proceedings of the 45th International Conference on Software Engineering* (Melbourne, Victoria, Australia) (*ICSE '23*). IEEE Press, 2009–2021. <https://doi.org/10.1109/ICSE48619.2023.00170>
- [48] Shannon Zhu et al. 2018. *pyre, a performant type checker for Python*. Facebook. <https://github.com/facebook/pyre-check> Accessed: 2025-05-13.
- [49] Jelle Zijlstra, Nicolas, Alex Waygood, and Bartosz Ślawecki. 2017. *typeshed_client*. https://github.com/JelleZijlstra/typeshed_client.git Accessed: 2025-05-27.