

Modular, parsing-based, flow analysis of dictionary data structures in scripting languages

David A. Schmidt*

Kansas State University, Manhattan, Kansas, USA

Abstract. We design and implement a modular, constant-propagation-like forwards flow analysis for a Python subset containing strings and dictionaries (hash tables). The analysis infers types of dictionaries and the functions and modules that use them. Unlike records and class-based objects, dictionaries are wholly dynamic, and we employ a domain of dictionary types that delineate which fields a dictionary must have. We have deliberately omitted unification-based inference and row variables to obtain the benefits of a forwards analysis that matches a programmer’s intuitions. Nonetheless, to accommodate a modular analysis, the values of parameters and free (global) variables are represented by tokens to which are attached constraints. At link- and function-call-time, the constraints are matched against the actual values of arguments and global variables.

Finally, programmers are encouraged to use a BNF-like syntax to define the forms of data types employed in their scripts. The analysis uses the programmer-written BNF rules to “abstractly parse” program phrases and associate them with derivations possible from the programmer-defined grammars. A prototype of the system is under construction.

1 Introduction

Programmers like the dynamic data structures of scripting languages like Perl, Ruby, and Python. The dictionary (hash table) structure is particularly useful for data modelling. For example, one can model an empty, heterogenous binary tree by the string, `'nil'`, and a nonempty tree by the dictionary, `{'val': V, 'left': T1, 'right': T2}`, where `V` can be any phrase whatsoever and each `Ti` can be a subtree. By omitting the `'left'` and `'right'` fields, we model partial trees as well. Functions that process partial trees are simply written; see Figure 1.

The partial-tree data structure is difficult to define simply in a typed language. But such heterogeneous data structures are common in practice and are handled well by dynamically-typed scripting languages. Unfortunately, when a single script or a collection of script-coded modules grows above 5000 lines, loss of comprehension cancels the gain of simplicity. Further, the time spent debugging data-structure mismatch errors (“typing errors”) becomes significant.

* das@ksu.edu. Supported by NSF CNS-0939431.

```

def count(t): # counts the values in tree t
  if t == 'nil': # empty tree ?
    return 0
  elif isinstance(t, dict): # a nonempty tree ?
    ans = 1
    if 'left' in t: ans = ans + count(t['left'])
    if 'right' in t: ans = ans + count(t['right'])
  return ans

```

Fig. 1. Function that counts nodes in a partial binary tree

As researchers have noted [2, 3, 8–11], an ideal solution would be a minimally programmer-assisted, type-inference tool to check script modules for data-structure mismatches. This paper describes an approach that does dictionary checking and complements the work on record and class-based object checking [1, 4, 6, 12]

2 Related efforts

Two key previous efforts relevant to our work are *soft typing* [2] and *gradual typing* [8, 9].

Soft typing is an ML-like type inference which infers types from the type domain of primitive types, function types, true unions, and fixed points. Phrases that cannot be assigned a type are wrapped with a run-time cast. The standard implementation of soft typing encodes a union type as a record whose fields enumerate all the type constructors in the programming language; a field’s value is set to *must* if its constructor is a component of the union type, and *must-not*, when it is certainly excluded. To define subtyping of unions, a field can also take on a *may* value. Remy’s inference algorithm for records, paired with Huet’s unification algorithm, are used to perform the inference and synthesis of recursive defined union types [6]. Soft typing is elegant but is not modular and does not scale well [3].

Gradual typing is a type inference for a type calculus that is augmented by a type *dynamic* (“?”) or “unknown.” The typing is named “gradual” because the programmer is encouraged to refine the program with more and more types to reduce occurrences of ?. Checking is performed up to the occurrences of ?, which are checked at run-time. This induces a type compatibility relationship, \sim , such that $int \sim ?$, $(int \rightarrow ?) \sim (? \rightarrow bool)$, etc. The definition of \sim is compatible with type refinement (replacement of occurrences of ? by concrete types). Gradual typing has been defined for ML and Java subsets (whose class definitions aid type inference) [8, 9].

Two other recent efforts deserve mention: Furr, et al. [4] have defined a Ruby subset, DRuby, along with a type-inference algorithm. DRuby does not support dictionaries, and it depends on class definitions to check object creation and message passing. Jensen, et al. [5] have designed a precise flow analysis to detect

$C \in Command$	$B \in BoolExpression$
$L \in LeftHandSide$	$I \in Identifier$
$E \in Expression$	$S \in StringConstant$
$C ::= L = E \mid \text{if } B: C_1 \text{ else: } C_2 \mid \text{while } B: C$ $\mid \text{def } I_1(I_2): C \mid \text{global } I \mid \text{return } E \mid C_1; C_2$	
$L ::= I \mid E[I]$	
$E ::= S \mid \text{stringop } E \mid \{E_i : E'_i\}_{0 \leq i < m} \mid L \mid L(E)$	
$B ::= E_1 == E_2 \mid E \text{ in } L \mid \text{isinstance}(L, \text{str})$ $\mid \text{isinstance}(L, \text{dict}) \mid \text{callable}(L)$	

Fig. 2. Syntax of Python subset

errors in Javascript scripts, which use prototype objects that provide information similar to a class definition’s.

Over a number of years, Felleisen and his colleagues have implemented and profiled several soft-typing and set-based analyses for module-based Scheme [3]. Their conclusions are that the methods do not scale well, are “brittle” (type errors easily arise) and do not provide good error diagnostics, and must be augmented by additional theory to accommodate modules [3]. Their latest efforts have moved towards type checking of Scheme augmented by programmer-written ML-like `datatype` definitions and annotated case-selector functions [3, 10, 11].

3 Script modules with dictionaries

Based on our experiences working with scripting languages and based on the research documented above, we learned there is value in having a tool assistant that can warn a programmer of mismatches in dictionary usage. We realize that is impossible to impose a static typing discipline on a script, so we search for a middle ground, where a modular, forwards, data-flow-analysis-based inference algorithm provides directives to the programmer where repairs might be necessary. For the reasons outlined by Felleisen [3], we employ an object-based, forwards, iterative data-flow constant-propagation algorithm as our tool’s software architecture. The tool accepts programmer assistance supplied as generalized `datatype` definitions and parameter annotations.

We start our work with a Python subset, consisting of strings, dictionaries, conditionals, loops, and functions. Every Python script is itself a module, which is executed (*imported*) and then can be linked to other modules that use the script’s public namespace (variable and function definitions). The source syntax is defined in Figure 2.

Figure 3 shows a small script with strings, dictionaries, and nested scopes. Python employs static scoping on its namespaces (environments), but declarations are unneeded — namespaces are dynamically filled, as when `f` is called in the example: `f`’s local namespace holds a binding for `z`, and its global names-

```

def f(z):
    global x
    if isinstance(y, dict): y[x] = z
    x = y
def g():
    y = 'cd'; f(y)
x = 'ab'; y = {}; g()
while x in y: # is x a field in dictionary y ?
    print y[x]; y = y[x]

```

Fig. 3. Sample script/module

pace, whose identity was fixed at `f`'s definition, is consulted for `x`. Since `y` cannot be found in `f`'s local namespace, the global namespace is searched for it. Note that global `y`'s value is *not* changed by calling `g`. Finally, there is a potential for generating an exception by evaluating the test, `x in y`, at the head of the while-loop, because `y` might take on a non-dictionary value. (This possibility will be detected by our analysis.) The example shows that Python's static scoping and its dictionaries pose problems in that both are modelled by namespaces that are dynamically filled with fieldname, value pairs. Also, fieldnames for dictionary lookups are computed at run-time, like array indexes, and dictionaries can be aliased.

Figures 4 and 5 give the language's denotational semantics. The semantics shows that a command's environment is a stack (list) of locations (addresses) of namespaces in the heap. The environment is statically calculated, but each namespace referenced by an environment is dynamically filled. The "*global*" token is a "patch" that allows a function to assign to a top-level variable. Note that primitive values (here, strings), dictionaries (namespaces) and closures can be stored in the heap and assigned to variable names.

4 Modular forwards constant-propagation analysis

The error-detection analysis is an abstract interpretation of the concrete semantics, where the main abstraction is upon namespaces:

$$\begin{aligned}
 AbsNS &= \{(AbsStr : AbsDenotable)^*\}[\!] \\
 AbsStr &= String \cup \{Str_{any}\} \\
 AbsDenotable &= SetOf(AbsStr + NSLocation + ClLocation)
 \end{aligned}$$

That is, an abstract namespace is a dictionary of abstract-string, denotable-set, mappings with an optional `!` suffix, which denotes "exactly". For example, the abstract namespace, `{ "a" : {"mm", ℓ_0 }, "b" : Strany }!` denotes all namespaces that have exactly two fields, "`a`" and "`b`", where "`a`" denotes either string "`mm`" or location ℓ_0 and "`b`" maps to some string (represented by `Strany`). Since "`b`"'s value is a singleton set, we omit the set braces to ease eye strain. If we remove the `!` suffix, the abstract namespace denotes those namespaces that have *at least* the

Semantic domains:

$\rho \in \text{Environment} = \text{NSLocation}^*$ (stack of namespaces' locations)
 $\sigma \in \text{Heap} = (\text{NSLocation} \rightarrow \text{Namespace}) \times (\text{CILocation} \rightarrow \text{Closure})$
 $\text{Namespace} = \text{Iden} \rightarrow \text{Denotable}$, where $\text{Iden} \subseteq \text{String}$
 $\text{Denotable} = \text{String} + \text{NSLocation} + \text{CILocation} + \text{"global"}$
 $\text{Expressible} = \text{Denotable} + \text{"error"}$
 $\text{Closure} = \text{Iden} \times \text{Command} \times \text{Environment}$
 $\text{LValue} = \text{NSLocation} \times \text{Iden}$
 $\text{Outcome} = \text{Heap} + (\text{Expressible} \times \text{Heap})$

Auxiliary functions:

$\text{allocateNS} : \text{Denotable} \times \text{Heap} \rightarrow \text{NSLocation} \times \text{Heap}$
 $\text{allocateNS}(v, \sigma) = (\ell_0, \sigma + [\ell_0 \mapsto v])$, where ℓ_0 is a fresh location
 $\text{allocateCl} : \text{Denotable} \times \text{Heap} \rightarrow \text{CILocation} \times \text{Heap}$ (similar)

$\text{findNSof} : \text{Iden} \times \text{Environment} \times \text{Heap} \rightarrow \text{NSLocation} + \text{"none"}$
 $\text{findNSof}(i, \text{nil}, \sigma) = \text{"none"}$
 $\text{findNSof}(i, \ell :: \rho', \sigma) = \text{if } i \in \text{dom}(\sigma(\ell))$
 $\quad \text{then if } \sigma(\ell) = \text{"global" then last}(\ell :: \rho') \text{ else } \ell$
 $\quad \text{where last}(\ell_0 :: \ell_1 :: \dots :: \ell_n) = \ell_n$
 $\quad \text{else findNSof}(i, \rho', \sigma)$

$\text{lookup} : \text{LValue} \times \text{Heap} \rightarrow \text{Denotable} + \text{"error"}$
 $\text{lookup}((\ell, i), \sigma) = \text{if } i \in \sigma(\ell) \text{ then } \sigma(\ell)(i) \text{ else "error"}$

$\text{update} : \text{LValue} \times \text{Denotable} \times \text{Heap} \rightarrow \text{Heap}$
 $\text{update}((\ell, i), v, \sigma) = \sigma + [\sigma(\ell) \mapsto \sigma(\ell) + [i : v]]$

$\text{apply} : (\text{Heap} \rightarrow \text{Outcome}) \times \text{Outcome} \rightarrow \text{Outcome}$
 $\text{apply}(f, \psi) = \text{if } \psi \in \text{Heap} \text{ then } f(\psi) \text{ else } \psi$

$C : \text{Command} \rightarrow \text{Environment} \rightarrow \text{Heap} \rightarrow \text{Outcome}$

$\text{let } \theta, \sigma' = \mathcal{L}[L]\rho\sigma$
 $C[L = E]\rho\sigma = \text{let } v, \sigma'' = \mathcal{E}[E]\rho\sigma'$
 $\quad \text{in update}(\theta, v, \sigma'')$

$C[\text{if } B : C_1 \text{ else } : C_2]\rho\sigma = \text{let } b, \sigma' = \mathcal{B}[B]\rho\sigma$
 $\quad \text{if } b \text{ then } C[C_1]\rho\sigma' \text{ else } C[C_2]\rho\sigma'$

$C[\text{while } B \text{ do } : C]\rho = f,$
 $\quad \text{where } f\sigma = \text{let } b, \sigma' = \mathcal{B}[B]\rho\sigma$
 $\quad \text{if } b \text{ then apply}(f, C[C]\rho\sigma') \text{ else } \sigma'$

$\text{let } \theta, \sigma_1 = \mathcal{L}[I_1]\rho\sigma$
 $C[\text{def } I_1(I_2) : C]\rho\sigma = \text{let } \ell_2, \sigma_2 = \text{allocateNS}(\{\}, \sigma_1)$
 $\quad \text{let } \ell_3, \sigma_3 = \text{allocateCl}((I_2, C, (\ell_0 :: \rho)), \sigma_2)$
 $\quad \text{in update}(\theta, \ell_3, \sigma_3)$

$C[\text{global } I]\rho\sigma = \text{update}((\text{head}(\rho), I), \text{"global"}, \sigma)$

$C[\text{return } E]\rho\sigma = \mathcal{E}[E]\rho\sigma$

$C[C_1 ; C_2]\rho\sigma = \text{apply}(C[C_2]\rho, C[C_1]\rho\sigma)$

Fig. 4. Semantics of Python subset

$$\begin{aligned}
\mathcal{L} &: \text{LeftHandSide} \rightarrow \text{Environment} \rightarrow \text{Heap} \rightarrow ((\text{LValue} + \text{"error"}) \times \text{Heap}) \\
\mathcal{L}[I]\rho\sigma &= \text{let } \ell = \text{findNSof}(I, \rho, \sigma) \\
&\quad \text{if } \ell = \text{"none"} \text{ then } ((\text{head}(\rho), I), \sigma) \text{ else } ((\ell, I), \sigma) \\
\mathcal{L}[E[I]]\rho\sigma &= \text{let } \ell, \sigma' = \mathcal{E}[E]\rho\sigma \\
&\quad \text{if } v \in \text{NSLocation} \text{ then } ((v, I), \sigma') \text{ else } (\text{"error"}, \sigma') \\
\mathcal{E} &: \text{Expression} \rightarrow \text{Environment} \rightarrow (\text{Expressible} \times \text{Heap}) \\
\mathcal{E}[\text{"abc"}]\rho\sigma &= \text{"abc"}, \sigma \\
\mathcal{E}[L]\rho\sigma &= \text{let } \theta, \sigma' = \mathcal{L}[L]\rho\sigma \text{ in } (\text{lookup}(\theta, \sigma'), \sigma') \\
&\quad \text{let } \ell, \sigma' = \mathcal{E}[L]\rho\sigma \\
&\quad \text{if } \ell \notin \text{CLocation} \text{ then } (\text{"error"}, \sigma') \\
\mathcal{E}[L(E)]\rho\sigma &= \text{else let } I, C, (\ell_0 :: \rho_0) = \sigma(\ell) \\
&\quad \text{let } v, \sigma'' = \mathcal{E}[E]\rho\sigma' \\
&\quad \text{in if } v = \text{"error"} \text{ then } (\text{"error"}, \sigma'') \\
&\quad \quad \text{else } \mathcal{C}[C](\ell_0 :: \rho)(\text{update}((\ell_0, I), v, \sigma'')) \\
&\quad \quad \text{let } \ell, \sigma_0 = \text{allocateNS}(\sigma) \\
&\quad \quad \text{foreach } E_i, E'_i, 0 \leq i < m, \\
\mathcal{E}[\{E_i : E'_i\}_{0 \leq i < m}]\rho\sigma &= \quad \text{let } \phi_i, \sigma'_i = \mathcal{E}[E_i]\rho\sigma_i \\
&\quad \quad \text{let } v_i, \sigma''_i = \mathcal{E}[E'_i]\rho\sigma'_i \\
&\quad \quad \text{let } \sigma_{i+1} = \text{update}((\ell, \phi_i), v_i, \sigma''_i) \\
&\quad \quad \text{in } (\ell, \sigma_m) \\
\mathcal{B} &: \text{BoolExpr} \rightarrow \text{Environment} \rightarrow \text{Heap} \rightarrow (\text{Boolean} + \text{"error"}) \times \text{Heap} \\
&\quad \text{let } v_1, \sigma_1 = \mathcal{E}[E]\rho\sigma \\
&\quad \text{let } v_2, \sigma_2 = \mathcal{E}[L]\rho\sigma_1 \\
\mathcal{B}[E \text{ in } L]\rho\sigma &= \text{if } v_1 \notin \text{String} \text{ or } v_2 \notin \text{NSLocation} \text{ then } (\text{"error"}, \sigma_2) \\
&\quad \text{else } (v_1 \in \text{dom}(\sigma_2(v_2)), \sigma_2) \\
&\quad \quad \text{let } v, \sigma' = \mathcal{E}[L]\rho\sigma \\
\mathcal{B}[\text{instanceof}(L, \text{dict})]\rho\sigma &= \text{if } v = \text{"error"} \text{ then } (\text{"error"}, \sigma') \\
&\quad \text{else } (v \in \text{NSLocation}, \sigma') \\
\mathcal{B}[\text{callable}(L)]\rho\sigma &= \text{let } v, \sigma' = \mathcal{E}[L]\rho\sigma \\
&\quad \text{if } v = \text{"error"} \text{ then } (\text{"error"}, \sigma') \text{ else } (v \in \text{CLocation}, \sigma') \\
&\quad \text{let } v_1, \sigma_1 = \mathcal{E}[E_1]\rho\sigma \\
&\quad \text{let } v_2, \sigma_2 = \mathcal{E}[E_2]\rho\sigma_1 \\
\mathcal{B}[E_1 == E_2]\rho\sigma &= \text{if } v_1 \text{ or } v_2 = \text{"error"} \text{ then } (\text{"error"}, \sigma_2) \text{ else } (v_1 \equiv v_2, \sigma_2) \\
&\quad \text{where } \equiv \text{ defines equality on strings, "deep equality"} \\
&\quad \quad \text{on namespaces, and CLocation equality on closures}
\end{aligned}$$

Fig. 5. Semantics of Python subset, concl.

"a" and "b" fields with the values just described. Figure 6 defines the abstract patterns, their concretizations, and their partial orderings.

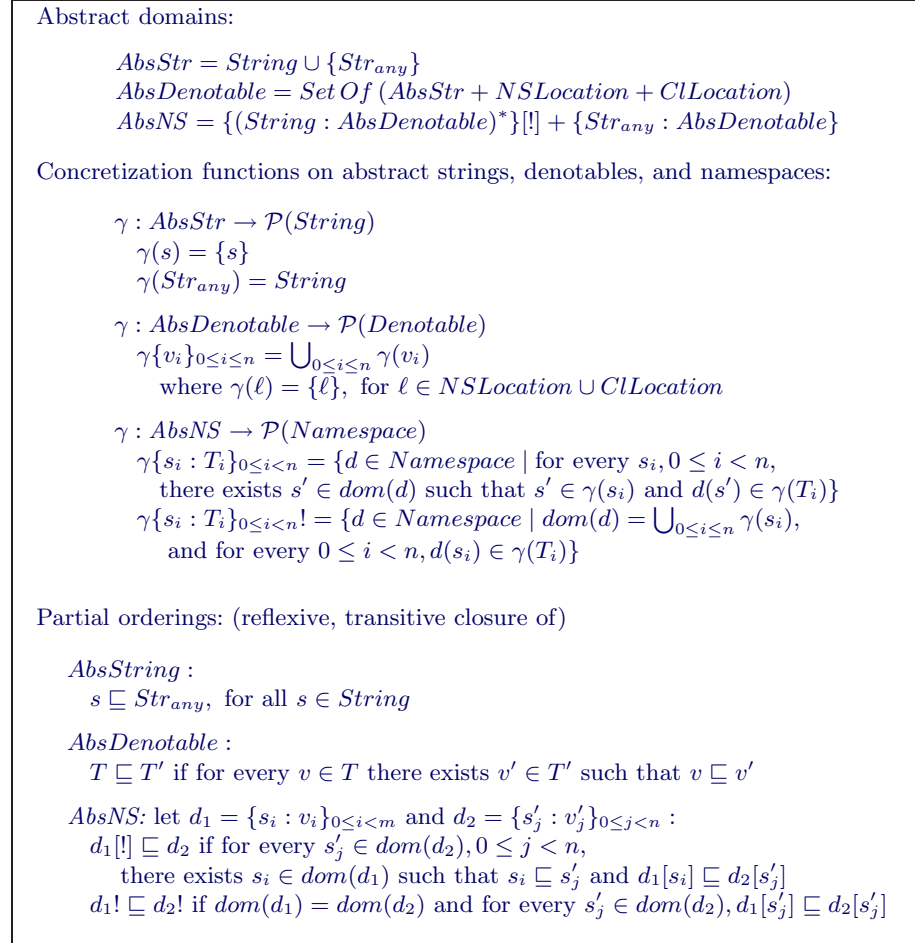


Fig. 6. Abstract domains of strings and namespaces

The partial orderings are used by the flow-analysis algorithm to compute least fixed points — the key property is that $P_1 \sqsubseteq P_2$ implies $\gamma(P_1) \subseteq \gamma(P_2)$. Since the token, Str_{any} , denotes “any string at all”, we have that “abc” $\sqsubseteq Str_{any}$, etc. An abstract denotable value is now a set of possible denotables, and two abstract denotable sets are compared by using the lower powerset ordering, e.g., $\{\text{“abc”}, \ell_2\} \sqsubseteq \{Str_{any}, \ell_2, \ell_1\}$. Dictionaries are ordered by contravariance on the field names: $\{\text{“a”} : \{\text{“m”}, \text{“n”}\}, \text{“b”} : \{\ell_0\}\}! \sqsubseteq \{\text{“a”} : \{Str_{any}\}, \text{“b”} : \{\ell_0, \text{“n”}\}\}!$ $\sqsubseteq \{\text{“a”} : \{Str_{any}\}, \text{“b”} : \{\ell_0, \text{“n”}\}\} \sqsubseteq \{\text{“b”} : \{\ell_0, \text{“n”}\}\} \sqsubseteq \{\text{“b”} : \{\ell_0, Str_{any}\}\} \sqsubseteq \{Str_{any} : \{\ell_0, Str_{any}\}\} \sqsubseteq \{\}$.

Program and flow equations:

<pre style="margin: 0;"> x = {}_α while ...: x['a'] = {}_β x = x['a'] ... x['a'] ... </pre>	<pre style="margin: 0;"> P₀ = initialHeap_{n₀} P₁ = assign x {}_α P₀ P₂ = P₁ ⊔ P₄ P₃ = assign (index x "a") {}_β P₂ P₄ = assign x (index x "a") P₃ P₅ = ... (index x "a") ... P₂ </pre>
---	--

At all program points, P_i , the environment has value, $[n_0]$, where n_0 is the address of the global namespace. The initial heap is $[n_0 \mapsto \{\}] \times []$ — there are no variable bindings yet in n_0 's namespace. The second heap component, $[]$ (closure bindings) stays empty and is omitted hereon. The worklist solution to the flow equations goes as follows:

```

P0 = [n0 ↦ {}!]
P1 = [n0 ↦ {x : lα}!, lα ↦ {}!]
P2 = [n0 ↦ {x : lα}!, lα ↦ {}!]
P3 = [n0 ↦ {x : lα}!, lα ↦ {"a" : lβ}!, lβ ↦ {}!]
P4 = [n0 ↦ {x : lβ}!, lα ↦ {"a" : lβ}!, lβ ↦ {}!]

P2 = [n0 ↦ {x : {lα, lβ}}, lα ↦ {}, lβ ↦ {}!]
P3 = [n0 ↦ {x : {lα, lβ}}, lα ↦ {"a" : lβ}, lβ ↦ {"a" : lβ}!]
P4 = [n0 ↦ {x : lβ}!, lα ↦ {"a" : lβ}, lβ ↦ {"a" : lβ}!]

P2 = [n0 ↦ {x : {lα, lβ}}, lα ↦ {}, lβ ↦ {}]
P3 = [n0 ↦ {x : {lα, lβ}}, lα ↦ {"a" : lβ}, lβ ↦ {"a" : lβ}]
P4 = [n0 ↦ {x : lβ}!, lα ↦ {"a" : lβ}, lβ ↦ {"a" : lβ}]

P2 = [n0 ↦ {x : {lα, lβ}}, lα ↦ {}, lβ ↦ {}]
(converges)

```

(To reduce eye strain, the brackets around singleton-set abstract denotables are omitted.)

Fig. 7. Sample program and its flow analysis

Using the abstract domains, we compute a forwards data-flow analysis, where each line in the source program generates one or more flow equations. The family of simultaneously defined flow equations are solved so that each equation computes to an environment (stack of *NSlocation*) and abstract heap of form $(NSlocation \rightarrow AbsNS) \times (CLlocation \rightarrow Closure)$. Since Python is statically scoped, the computation of the environments can be undertaken prior to the computation of the abstract heaps.

Figure 8 shows a small loop program, its flow equations, and the solution to the equations. Because the loop's iterations are unknown, so is the quantity of dictionaries allocated within the loop. For simplicity, we employ *summary nodes* [7, 13], l_α and l_β , to denote the dictionary objects allocated during the analysis. The heap at loop exit is summarized by $[n_0 \mapsto \{x : \{l_\alpha, l_\beta\}\}, l_\alpha \mapsto \{\}, l_\beta \mapsto \{\}]$. That is, at loop exit, x (and no other variable) is definitely

defined in the global namespace and denotes a dictionary. The dictionary’s value is overapproximated by $\{\}$, which means the dictionary’s fields and values are uncertain; the dictionary might be empty. (Indeed, this is the case at run time.) The heap lets us analyze the reference to $x['a']$ at point P_5 — a missing-field error might occur, since there is no guarantee that the field labelled 'a' exists.

The primary purpose of the analysis is to detect operator-operand incompatibility errors and missing-field errors (*key errors*) in dictionary usage. The analysis looks superficially like type analyses for class-based object-oriented languages, but as the previous example shows, dictionaries pose new challenges in that fields can be added dynamically and field names for lookups can be computed at run-time.

5 Function-definition analysis

To make the analysis modular, we analyze a function definition with an abstract heap that contains symbolic constants, called *tokens*, for the values of the function’s free variables, namely, the function’s formal parameters and nonlocal variables.

While the function body is analysed, references to the free variables generate *conditional constraints* on the corresponding tokens. The constraints are included along with the function’s output image of the abstract heap, as the function’s abstract denotation — a summary template. At function-call time, the template is compared to the actual-parameter arguments and global-variable values in the call-time heap for a compatibility check (“type check”). If the constraints are satisfied, the call-time actual-parameter arguments and global-variable values are substituted for the tokens in the function’s template, producing the function call’s output heap, which is then reconciled with the input heap to the call. (This last step is discussed in the next section.)

Figure 8 gives an example of a template resulting from an analyzed function. The value of P_1 shows how free variables d and g reside in the local and global namespaces, respectively. The conditional’s test generates the constraint that t_d is $\{\}$ — d ’s value is a dictionary — in the then-arm. Since d is dereferenced within the then-arm, the constraint is expanded into a conditional constraint, t_d is $\{\} \Rightarrow (t_d$ is $\{ "b" : t_d \})$ — if d ’s initial value is a dictionary, then it must possess a "b" field. The updates to $d["c"]$ and $d["b"]$ are expressed *in the heap, on fresh tokens*, as seen in P_5 .

P_6 ’s value shows how g is updated in the else arm; note the constraint that d is “not a dictionary” within the else arm. The negation can be expressed as a disjunction of all the type structures that are not dictionary structures (here, “string or closure”) or can be left as is; see Section 8.

Finally, the function’s output is defined by P_7 . The output heap collects updates to both local and free variables. Precision is lost by joining the heaps at P_5 and P_6 . (In Section 9 we see how to regain precision by annotating d with a user-defined type that lets us generate multiple templates for a function definition.) The output heap is used at function-call time to update the call-time

Function and flow equations:

<pre> def f(d): global g if isinstance(d,dict): d['c'] = g g = d['b'] d['b'] = {} else: g = d </pre>	$ \begin{aligned} P_0 &= \text{assign } d \ t_d \ \text{initialHeap}_{n_1, n_0} \\ P_1 &= \text{assign global } g \ t_g \ P_0 \\ P_{2t} &= \text{makeTrue } (\text{isinstance } d \ \text{dict}) \ P_1 \\ P_3 &= \text{assign } (\text{index } d \ \text{"c"}) \ g \ P_{2t} \\ P_4 &= \text{assign } g \ (\text{index } d \ \text{"b"}) \ P_3 \\ P_5 &= \text{assign } (\text{index } d \ \text{"b"}) \ \{\} \ P_4 \\ P_{2f} &= \text{makeFalse } (\text{isinstance } d \ \text{dict}) \ P_1 \\ P_6 &= \text{assign } g \ d \ P_{2f} \\ P_7 &= P_5 \sqcup P_6 \end{aligned} $
--	---

The environment is $[n_1 :: n_0]$, where n_0 is the address of the global namespace. The initial heap is $[n_1 \mapsto \{\}!, n_0 \mapsto \{\}] \times [\dots] \times ()$, where the second component maps defined functions to their templates and the third component collects constraints on the function's free variables. The worklist solution is

$$\begin{aligned}
P_0 &= [n_1 \mapsto \{d : t_d\}!, n_0 \mapsto \{\}], [\dots], () \\
P_1 &= [n_1 \mapsto \{d : t_d\}!, n_0 \mapsto \{g : t_g\}], [\dots], () \\
P_{2t} &= [n_1 \mapsto \{d : t_d\}!, n_0 \mapsto \{g : t_g\}, t_d \mapsto \{\}], [\dots], t_d \text{ is } \{\} \Rightarrow () \\
P_3 &= [n_1 \mapsto \{d : t_d\}!, n_0 \mapsto \{g : t_g\}, t_d \mapsto \{\text{"c"} : t_g\}], [\dots], \\
&\quad t_d \text{ is } \{\} \Rightarrow () \\
P_4 &= [n_1 \mapsto \{d : t_d\}!, n_0 \mapsto \{g : t_{d'}\}, t_d \mapsto \{\text{"c"} : t_g, \text{"b"} : t_{d'}\}], [\dots], \\
&\quad t_d \text{ is } \{\} \Rightarrow (t_d \text{ is } \{\text{"b"} : t_{d'}\}) \\
P_5 &= [n_1 \mapsto \{d : t_d\}!, n_0 \mapsto \{g : t_{d'}\}, t_d \mapsto \{\text{"b"} : \ell_0, \text{"c"} : t_g\}, \ell_0 \mapsto \{\}!], [\dots], \\
&\quad t_d \text{ is } \{\} \Rightarrow (t_d \text{ is } \{\text{"b"} : t_{d'}\}) \\
P_{2f} &= [n_1 \mapsto \{d : t_d\}!, n_0 \mapsto \{g : t_g\}], [\dots], t_d \text{ is } \sim \{\} \Rightarrow () \\
P_6 &= [n_1 \mapsto \{d : t_d\}!, n_0 \mapsto \{g : t_d\}, t_d \text{ is } \sim \{\} \Rightarrow () \\
P_7 &= [n_1 \mapsto \{d : t_d\}!, n_0 \mapsto \{g : \{t_{d'}, t_d\}\}, t_d \mapsto \{\text{"b"} : \ell_0, \text{"c"} : t_g\}, \ell_0 \mapsto \{\}!], [\dots], \\
&\quad t_d \text{ is } \{\} \Rightarrow (t_d \text{ is } \{\text{"b"} : t_{d'}\})
\end{aligned}$$

The value of P_7 , heap and constraints, constitute the function's template.

Fig. 8. Analysis of a function

heap, and the collected constraints are validated on the call-time arguments and call-time heap. We consider this next.

6 Function-call processing

A function's template is applied when the function is called. These steps are undertaken:

1. At the point of call, the arguments and global variables are checked against the constraints listed in the function's template. If a constraint cannot be validated, a mismatch error (type error) is generated. If any of the arguments or globals themselves contain tokens, the constraints are copied into the caller's constraint set.

2. The value of the function template’s output heap is used to update the call-time heap: For each of the namespaces listed in the caller’s environment, the values of the function’s free variables overwrite the values in the call-time heap. Here, locally constructed objects can “escape” from the called function’s body and can be tracked by the analysis (cf. “escape analysis [13]).

Here are two calls of the function that was analyzed in Figure 8. First, we have the call,

```

g = 'mm'
h = {'b' : 'n'}
f(h)

```

At the point of the call, $\mathbf{f}(\mathbf{h})$, the call-time heap is

$$[\{n_0 \mapsto \{g : "mm", h : \ell_0\}!, \ell_0 \mapsto \{"b" : "n"}\}!].$$

The call generates these two bindings to the function’s tokens: $t_d = \ell_0$ and $t_g = "mm"$. With the bindings, the analysis validates that the input constraint, $t_d \text{ is } \{\} \Rightarrow (t_d \text{ is } \{"b" : t_d\})$, holds true; in the process, the binding, $t_d = "n"$, is established. The three bindings are applied to the function’s output heap — recall that it is:

$$[n_1 \mapsto \{d : t_d\}!, n_0 \mapsto \{g : \{t_d, t_d\}, t_d \mapsto \{"b" : \ell'_0, "c" : t_g\}, \ell'_0 \mapsto \{\}!]$$

the above heap is reconciled with the call-time heap, generating this heap that results from the completed function call:

$$[n_0 \mapsto \{g : \{"n", \ell_0\}, h : \ell_0, \}!, \ell_0 \mapsto \{"b" : \ell'_0, "c" : "mm"}!, \ell'_0 \mapsto \{\}!$$

The example shows the loss of precision regarding \mathbf{g} ’s value. Precision can be improved if the function template would generate an output heap for each execution path (or, each generated constraint).

Here is a second, example call:

```

g = 'mm'
f(g)

```

At the point of call, the call-time heap is $[\{n_0 \mapsto \{g : "mm"}\}]$. The bindings, $t_g = "mm"$ and $t_d = "mm"$, validate the input constraint. The function’s output heap, listed above, is reconciled with the input heap, yielding $[n_0 \mapsto \{g : \{\perp, "nn"}\}]$, which equals

$$[n_0 \mapsto \{g : "nn"}]$$

In this example, there is no binding to t_d , since the token is a subpart of *dictionary* t_d , which does not appear in the output heap that results from this call. ℓ'_0 does not appear for the same reason.

<pre> ### Module M.py w = ... # global variable def set(x): global w w = x def get(): return ...w... # end module M ### Module N.py import M M.set(...) ...M.get()... # end module N </pre>	<pre> environment is [n_M] heap is [n_M ↦ {w : ...}]! function environment is [n_set :: n_M] entry heap is [n_set ↦ {x : t_x}!, n_M ↦ {w : t_w}] exit heap is [n_set ↦ {x : t_x}!, n_M ↦ {w : t_x}] function environment is [n_get :: n_M] heap is [n_get ↦ {}]!, n_M ↦ {w : t_w}] environment is [n_N] heap is [n_N ↦ {M : ℓ_0}!, ℓ_0 ↦ {set : ..., get : ...}]! heap is unchanged (update to w is hidden) </pre>
--	--

Fig. 9. Module definition, importation, and invocation

7 Module importation

Once a module is completely analyzed, a summary template must be generated from it. This is necessary because the caller of the module’s functions owns an environment whose namespaces are disjoint from the module’s. In particular, updates that a called function makes to a module’s variables are hidden from the caller. Figure 9 shows an example module, **M.py**, and its importation and use by another, **N.py**. The environment used by module **N** is disjoint from module **M**’s. This means the call, `M.set(...)`, does not record the update to `set`, since `set` has already been flow-analyzed.

But this demands that `set`’s analysis must ensure that all possible arguments for parameter `x` meet the the input constraints generated by `set`’s body. (And global variable `w` must hold a value that always satisfies the input constraints.)

For these reasons, the analysis first constructs module **M**’s summary template by unioning and solving all constraints that were generated by all **M**’s functions on all **M**’s variables. The solution to the constraints defines each module variable’s *invariant type*. Next, each function’s output heap is checked to validate that the values of the global variables in the output heap satisfy the invariant types. (If a function fails to satisfy a variable’s invariant type, that function cannot be called externally.)

The modular analysis of **M.py** generates a namespace of function templates that is imported by module **N** and used just like those in a local function call.

8 Filter functions

Within a scripting language, the tests in conditional commands play the role of classifying run-time values into types, and a useful analysis must exploit the information. To do so, the analysis interprets type-discriminating tests as

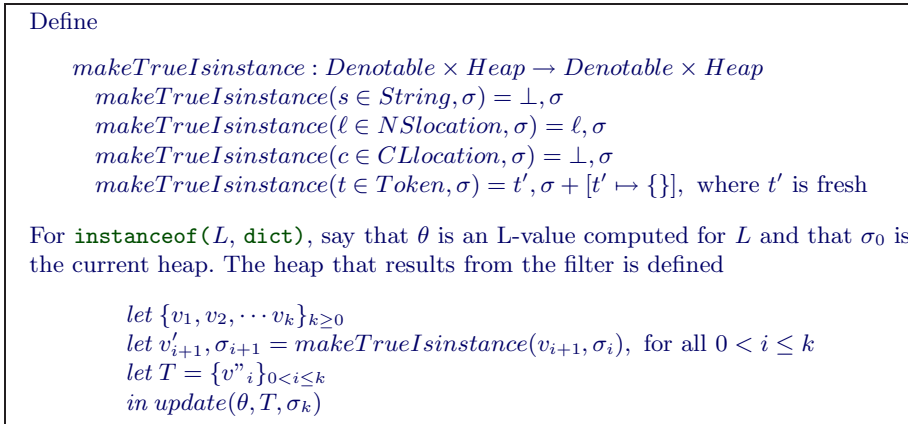


Fig. 10. Module definition, importation, and invocation

filter functions, which narrow (make more precise) the values that flow into the bodies of conditionals and loops. Figure 8 showed how the test predicate, **instanceof**(**d**, **dict**), was used to update the value of the heap that entered the conditional’s true arm to show that **d**’s value must be a dictionary (denoted, {}).

Figure 10 gives the definition used by the analysis for the **instanceof** predicate. For **instanceof**(*L*, **dict**), the definition examines each of the denotable values that *L* denotes in the abstract heap. Each denotable value is filtered through **instanceof**; some proceed, some do not. In the case of a token, the token is specialized to a dictionary value. The results are collected into a set, which updates *L*’s value.

The complement operation, *makeFalseInstance*, is defined similarly. In Figure 8, a “meta pattern”, $\sim\{\}$, was used as the output for *makeFalseInstance*. This metapattern can be instantiated to a set of denotables, as suggested by Figure 10, or can be carried along as is. In the case of test predicates too complex to filter (e.g., $E_1 == E_2$), the filter function can default to the identity filter, which overapproximates the flow into the arms of the conditional.

9 User-defined type equations and abstract parsing

User-defined types can improve the quality of an analysis; a user may write BNF-like type equations whose components are abstract-string and abstract-namespace values, e.g.,

$$N = \{ 's' : N \} ! \mid 'nil'$$

A type equation defines a named denotable value that can be calculated upon by the flow analysis. Figure 11 displays a example function whose analysis utilizes the above named type, *N*. The named type generates a finite set of denotables,

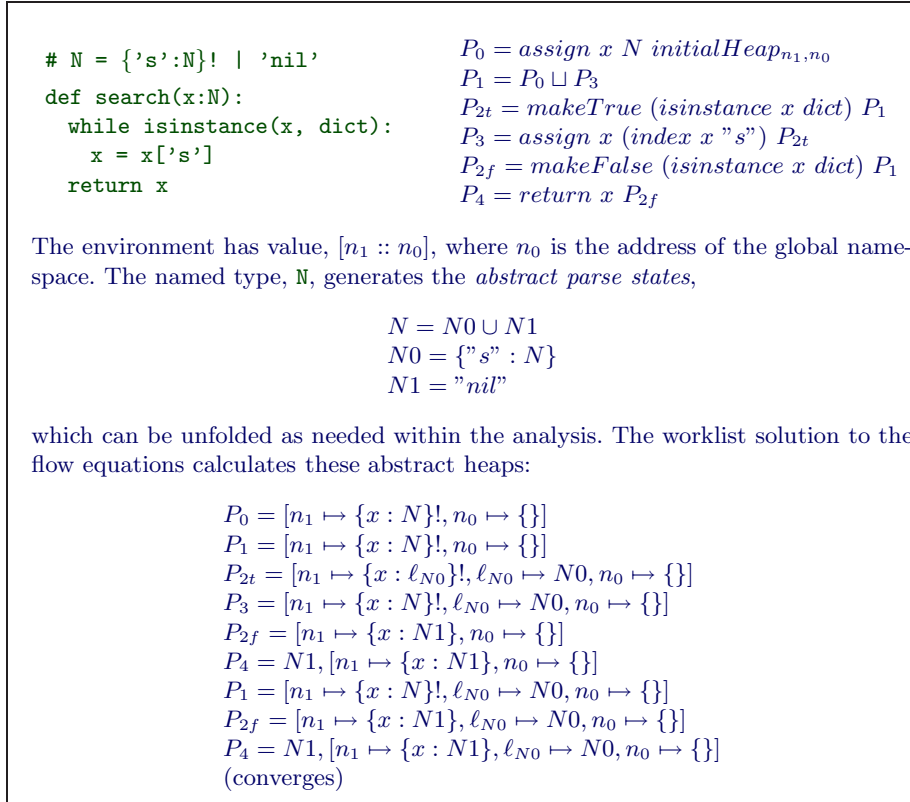


Fig. 11. Analysis of function with named parameter type

each of which is named by the parse state that is obtained by parsing a value with respect to the BNF equation. The parse states are used as denotable values by the analysis and are unfolded as needed (e.g., at the test, `isinstance(x,dict)`). In the example, parse state $N0$ acts as a “summary node,” limiting the size of the abstract heap.

The type definitions also allow the analysis to generate multiple templates per function. For example, Figure 11 could be revised to generate an analysis of `def search(x:N0)` and an analysis of `def search(x:N1)`. The separate analyses prove useful in practice when a function is written to act as a discriminator, to be used in the test of a conditional command [3].

10 Implementation status

We have implemented a prototype analyzer that incorporates much (but not all) of the developments in this paper. The analyzer has proved to be a useful development and refinement tool for the techniques outlined here. The work is

ongoing, and we expect refinements of the techniques described in this paper as a result of our profiling experiments.

References

1. M. Abadi and L. Cardelli. *A Theory of Objects*. Springer Verlag, 1996.
2. R. Cartwright and M. Fagan. Soft typing. In *Proc. Conf. on Prog. Lang. Design and Implementation*, pages 278–292. ACM Press, 1991.
3. M. Felleisen. From Soft Scheme to Typed Scheme: Experiences from 20 years of script evolution, and some ideas on what works. Technical Report www.ccs.neu.edu/home/matthias/Presentations/STOP/stop.pdf, Keynote address, Workshop on Scripts to Programs, Genova, Italy, 2009.
4. M. Furr, J.-H. An, J. Foster, and M. Hicks. Static type inference for Ruby. In *Proc. Symposium on Applied Computing*, pages 1859–1866. ACM Press, 2009.
5. S.H. Jensen, A. Möller, and P. Thiemann. Type inference for Javascript. In *Proc. Static Analysis Symposium*, pages 238–255. Springer LNCS 5673, 2009.
6. D. Remy. Typechecking records and variants in a natural extension of ML. In *Proc. Symp. on Principles of Prog. Lang.*, pages 77–88. ACM Press, 1989.
7. S. Sagiv, T. Reps, and R. Wilhelm. Parametric shape analysis via 3-valued logic. *ACM Trans. Programming Lang. and Systems*, 24:217–298, 2002.
8. J. Siek and W. Taha. Gradual typing for functional languages. In *Proc. Scheme and Functional Programming Workshop*, pages 374–388. ACM Press, 2006.
9. J. Siek and W. Taha. Gradual typing for objects. In *Proc. European Conf. on Object-Oriented Programming*, pages 2–27. Springer LNCS 4609, 2007.
10. S. Tobin-Hochstadt and M. Felleisen. Interlanguage migration: from scripts to programs. In *Proc. OOPSLA Companion*, pages 964–974. ACM Press, 2006.
11. S. Tobin-Hochstadt and M. Felleisen. The design and implementation of typed scheme. In *Proc. Symp. Principles of Prog. Languages*, pages 395–406. ACM Press, 2007.
12. M. Wand. Type inference for record concatenation and multiple inheritance. In *Proc. Symp. on Logic in Computer Science*, pages 92–97. IEEE Press, 1989.
13. J. Whaley and M. Rinard. Compositional pointer and escape analysis for Java programs. In *Proc. OOPSLA'99*, pages 187–206. ACM Press, 2009.