



Eidgenössische Technische Hochschule Zürich  
Swiss Federal Institute of Technology Zurich

# Static Program Analysis of Data Usage Properties

Master Thesis

Simon Wehrli

Chair of Programming Methodology  
Department of Computer Science  
ETH Zurich

**Supervisors:**

Dr. Caterina Urban  
Prof. Dr. Peter Müller

2017

---

## Abstract

We introduce a new property of program inputs that denotes if the input is *used* and develop a fully-automatic, static analysis within the framework of Abstract Interpretation to approximate this property for both single inputs and input lists of unknown length. The developed analysis extends the well-studied *truly-live variable analysis* to the stronger property of *usage*: It detects for input items, or more generally for any program variable, if they have any effect on the output of the program.

In Abstract Interpretation, an interpreter iterates over the code of a program without executing it on concrete inputs, and imitates parts of its behaviour with objects of an abstract universe, called abstract domain. Our analysis roots in an abstract *usage* domain, which holds information about which variables are used and how they can affect the usage property of other variables along the possible execution paths of the program. Combined with a stack-like structure, it can track the *flow of usage* through nested constructs like conditionals and loops. We can differentiate direct usage in the program output from implicit usage in conditions and transitive usage via assignments.

We then extend the usage domain to be able to track the usage property not only for variables, but also for list items. This means we can find out if the individual values of a list of inputs are used, even if the length of the list is statically unknown. We achieve this by integrating the usage domain into an abstract domain to analyze list contents. This domain is based on the idea of splitting the list into segments with dynamic bounds, each segment holding an abstraction of the values the segment is covering.

A prototype implementation of all presented domains and analysis was developed alongside this thesis. It demonstrates the potential usage scenarios in real-world programs typically written in data science applications.

---

# Contents

<b>1</b>	<b>Introduction</b>	<b>3</b>
<b>2</b>	<b>Language SimplePython</b>	<b>5</b>
2.1	Assumptions . . . . .	5
2.2	Input/Output . . . . .	5
2.3	Syntax . . . . .	5
<b>3</b>	<b>Background</b>	<b>8</b>
3.1	Abstract Interpretation . . . . .	8
3.1.1	Octagon Domain . . . . .	8
3.1.2	Interval Domain . . . . .	9
3.2	Usage Property . . . . .	9
<b>4</b>	<b>Generic Functor Abstract Domains</b>	<b>11</b>
4.1	Abstract store domain . . . . .	11
4.2	Abstract stack domain . . . . .	12
<b>5</b>	<b>Usage Analysis for Integer Variables</b>	<b>13</b>
5.1	Usage Abstract Domain . . . . .	13
5.2	Usage Stores . . . . .	13
5.3	Usage Stacks . . . . .	14
5.4	Concretization . . . . .	15
5.5	Example . . . . .	15
<b>6</b>	<b>List Start Usage Analysis</b>	<b>17</b>
6.1	Combining List Start Domain and Usage Domain . . . . .	19
6.2	Concretization . . . . .	20
6.3	Example . . . . .	21
<b>7</b>	<b>List Content Usage Analysis</b>	<b>22</b>
7.1	Ingredients . . . . .	22
7.1.1	Variable Environment Abstract Domain . . . . .	22
7.1.2	Expression Abstract Domain . . . . .	22
7.1.3	Segment Limit Abstract Domain . . . . .	22
7.1.4	Predicate Abstract Domain . . . . .	23
7.2	Generic List Segmentation Abstract Domain . . . . .	23
7.2.1	Segmentation Unification . . . . .	24
7.2.2	Segmentation Operators . . . . .	26
7.2.3	Abstract Transformers for Backward Analysis . . . . .	27
7.3	List Segmentation for Usage Property . . . . .	28
7.3.1	Integer/List stores . . . . .	28
7.3.2	Abstract Transformers . . . . .	30
7.3.3	Concretization . . . . .	30
7.3.4	Example . . . . .	31
<b>8</b>	<b>Implementation</b>	<b>32</b>
8.1	Framework for Abstract Interpretation . . . . .	32
8.2	Automated Testing . . . . .	32
8.3	Examples . . . . .	33
<b>9</b>	<b>Conclusion</b>	<b>36</b>

# 1 Introduction

In the past years, the Python programming language and associated libraries became popular for data science applications due to the simple syntax and broad integration with common libraries for data analytics. For quick data inspection, the dynamic nature of Python seems very attractive, but introduces some challenges for analysis tools. Current tools mainly aim at detecting programming errors that lead to unhandled runtime exceptions. But programming errors that do not trigger exceptions at all can still have serious consequences. Erroneous treatment or exclusion of data items can lead to wrong results while showing no indication of the faulty derivation of those results.

Consider the case when we exclude data items unintentionally. Suppose we wrote a program to calculate the number of passing grades, i.e., grades greater or equal to 4, out of a list of grades provided as the input to the program. Code 1 shows a possible implementation. The structure of this program is typical for the sort of programs we tackle. In line 1 we read a list of inputs. The function `listinput()` reads one full line with an unknown number of input grades, which then are split and converted to a list of integers. Lines 5 to 11 count the passing grades. On the last line we print out the computed count.

```
1 list_grades = listinput()
   count = 0
5 i = 1 # Bug B: should be 'i = 0'
   while i < len(list_grades):
       g = list_grades[i]
       if g >= 4:
           g += 1 # Bug A: should be 'count += 1'
10      i += 1
   print(count)
```

Code 1: A program with typical bugs we want to identify with our analysis.

The presented code contains two bugs: In line 5 we set the index  $i$  to the wrong starting value (Bug B) and in line 9 we increment the wrong variable (Bug A). Bug A let the whole while-loop be effectless on the result. On the other hand, Bug B alone (with Bug A fixed) let the first grade in the list be excluded from the counting.

The *live variable analysis* is used to separate live and dead variables in a program. A variable is dead in a line when its not live, and this is when it is assigned value is not read in all possible execution path after that line. In compiler design, this knowledge can be used for performance optimizations by omitting assignments to dead variables. An extension is the *truly-live variable analysis* [RHS95, GM81]:

**Definition 1.** Variable  $x$  is *truly-live* after a statement  $stmt$  (otherwise it is truly-dead/faint) if, and only if, there is a  $x$ -definition-free path from a successor of  $stmt$  to a statement  $stmt'$  that reads the value of  $x$  and:

- $stmt'$  is a print statement, or
- $stmt'$  is a condition, or
- $stmt'$  is an assignment statement that assigns to a variable  $y$  that is truly-live after  $stmt'$ .

Whenever a variable is truly-live, it is also live, so truly-live is the stronger property.<sup>1</sup>

If we apply the truly-live variable analysis to our example, it will classify all variables as live. In fact,  $g$  is used in a condition on line 8, so Bug A will not be revealed. Detecting Bug B is out of the scope of this analysis, since it has no possibility to track liveness properties of list items individually and will denote the whole list as live as soon as one item is accessed on line 7.

## Goals

In this thesis, we want to analyze a program for data usage properties, specifically, if an input data item is used in the calculations and has an (indirect) effect on the output. We want to find this with a static analysis of the program, so without presuming concrete input data that exposes a faulty behaviour. E.g., in Code 1 the analysis should detect that the first input grade has no effect on the printed result count, whatever the list of grades looks like in an actual execution. We target a sound approach that gives us mathematical guarantees of the specified properties about any syntactically correct program and input data. We will see that we choose our analysis to report unused inputs only if they are definitely unused in every actual execution.

We list the core goals of the thesis in detail and give an outline in which section we approach them:

- Define what is considered as an input data item to the analyzed program. We will see that we can handle inputs and variables in a similar fashion. (Section 2.2)
- Give a better intuition of the property “a data item is *used*” and formalize it. (Section 3.2)
- Design a static analysis within the framework of Abstract Interpretation [CC77] that approximates the solution to the question: Does a data item fulfill the *usage* property defined before? We build our analysis for a subset of Python (Section 2). The analysis has to be able to track the *flow of usage* through nested constructs like conditionals and loops. We should differentiate direct usage in the program output from implicit usage in conditions and transitive usage via assignments. First, we study this task on programs with scalar variables only (Section 5). We then extend the analysis to be able to track the usage property not only for variables, but also for list items. This means we can find out if the individual values of a list of inputs are used, even if the length of the list is statically unknown (Section 6 and 7).
- Implement the designed analysis (Section 8) and evaluate the analysis on several example programs (Section 8.3).

---

<sup>1</sup> The live variable property would only require  $x$  to be used in a reachable assignment *or* condition.

## 2 Language SimplePython

### 2.1 Assumptions

Some assumptions are made independently of the choice of a demonstration language. We assume that the programs analyzed

- are deterministic,
- are single-threaded,
- have a single pair of an initial and a final control point (before the initial and after the final statement, respectively),
- come type annotated<sup>2</sup>.

### 2.2 Input/Output

A fully-automatic approach, so without the need for user input, is to assume that every line read from Python's standard input `sys.stdin` is a data item. The number of inputs is not statically known and we will not suppose that it is necessarily finite. This adds a real challenge. Related analysis, such as the mentioned truly-live variable analysis, assume that the number of items of interest is statically known (e.g., all variables in the program).

It is reasonable to assume that the input of a program consists of a known number of single input values or lists (of unknown length) of inputs. Even if the structure is more complex, the input reading procedure will often store the inputs in one of these forms. This implies that we can design our analysis to determine the usage property on the level of scalar variables and list items. The result of this broader analysis then can be easily matched to actual inputs and unused ones can be reported to the user. We still include the two input expressions `input()` and `listinput()` in the syntax defined in the next section, to emphasize the goal of our analysis to determine properties about input items.

### 2.3 Syntax

For the formal derivation and the implementation prototype, we use a subset of the Python language called *SimplePython*. The purpose is to have a simple language that is sufficient to show how the explained analysis work in practice. The subset also contains some constructs for presentational purposes not available in Python, but they have equivalents in Python (see the footnotes in the syntax definition below).

Python is strongly typed, so we will assume the types can be inferred by a pre-analysis or are explicitly provided by the user. For ease of presentation, SimplePython only supports integers and single-dimension lists of integers. We do not use the type annotations syntax in our examples, but use variable names that indicate the type of the associated objects (because there are only two possible). A variable also has a fixed type all over the program.

---

<sup>2</sup> Type annotations are available in Python  $\geq 3.5$ , <https://docs.python.org/3.5/library/typing.html>

Constants:

$i \in \mathbb{Z}$  integer constants

Variables:

$x, y \in \mathbb{V}\text{ar}_{\text{int}}$  integer variables  
 $l \in \mathbb{V}\text{ar}_{\text{list}}$  list variables  
 $v \in \mathbb{V}\text{ar} = \mathbb{V}\text{ar}_{\text{int}} \cup \mathbb{V}\text{ar}_{\text{list}}$  general variables

Arithmetic expressions: (with and without list index)

$a, a_1, a_2 \in \mathbb{A}\text{expr} ::=$   $i$  integer constants  
|  $x$  integer variables  
|  $l[a]$  list index access  
|  $\text{input}()$  single integer input<sup>3</sup>  
|  $-a$  unary operator  
|  $a_1 + a_2$  |  $a_1 - a_2$  |  $a_1 * a_2$  |  $a_1 / a_2$  binary operators<sup>4</sup>  
 $t \in \mathbb{A}\text{expr}^- ::=$   $\mathbb{A}\text{expr}$  without any list index accesses

Boolean expressions:

$b, b_1, b_2 \in \mathbb{B}\text{expr} ::=$   $\text{True}$  truth  
|  $\text{False}$  falsity  
|  $\text{not } b$  negation  
|  $b_1 \text{ and } b_2$  conjunction  
|  $b_1 \text{ or } b_2$  disjunction  
|  $b_1 == b_2$  boolean comparison  
|  $a_1 == a_2$  |  $a_1 <= a_2$  |  $a_1 < a_2$  arithmetic comparison

Expressions:

$e \in \mathbb{E}\text{expr} ::= a$  |  $b$

List expressions:

$\text{disp} \in \mathbb{L}\text{istDisplaySeq} ::= a$  |  $\text{disp}, a$  list variable  
 $le \in \mathbb{L}\text{istExpr} ::= l$  list variable  
|  $[\text{disp}]$  list display  $[a_1, \dots, a_n]$   
|  $[\ ]$  empty list display  
|  $\text{listinput}()$  list input<sup>5</sup>

Statements:

$stmt \in \text{Stmt} ::=$	<b>pass</b>	no effect, identity
	$x = a$	integer assignment
	$l = le$	list assignment
	$l[a_1] = a_2$	list index update
	$\text{print}(e)$	output
	<b>if</b> $b : stmt_1$ <b>else</b> $stmt_2$	conditional
	<b>while</b> $b : stmt$	iteration
	$stmt_1; stmt_2$   $stmt_1$ <newline> $stmt_2$	sequence of statements

Program:

$prog \in \text{Prog} ::=$	$stmt$	program
----------------------------	--------	---------

For the rest of this report, we use this naming scheme when we show formal programs and single lines of code, like expressions and statements in abstract transformers. Additionally we use  $\psi \in \text{Var}_{\text{int}} \cup \text{Var}_{\text{list}} \cup \text{Index}$ , and  $\omega \in \text{Expr} \cup \text{ListDisplay}$ . All presented code samples of this report conform to this syntax specification. Reformulating the indicated input expressions into the corresponding Python functions makes all of our code samples directly runnable by the Python interpreter.

We define a convenience operator on expressions to retrieve all the operands, since in the usage analysis we are only interested in the appearing operands, not in the operators:

$\text{opd}(i) \stackrel{\text{def}}{=} i$	
$\text{opd}(x) \stackrel{\text{def}}{=} x$	
$\text{opd}(l[a]) \stackrel{\text{def}}{=} l[a]$	
$\text{opd}(-a) \stackrel{\text{def}}{=} \text{opd}(a)$	
$\text{opd}(a_1 \diamond a_2) \stackrel{\text{def}}{=} \text{opd}(a_1) \cup \text{opd}(a_2)$	$\diamond \in \{+, -, *, /\}$
$\text{opd}(a_1 \bowtie a_2) \stackrel{\text{def}}{=} \text{opd}(a_1) \cup \text{opd}(a_2)$	$\bowtie \in \{==, <=, <\}$
$\text{opd}(-b) \stackrel{\text{def}}{=} \text{opd}(b)$	
$\text{opd}(b_1 \square b_2) \stackrel{\text{def}}{=} \text{opd}(b_1) \cup \text{opd}(b_2)$	$\square \in \{\text{and}, \text{or}, ==\}$
$\text{opd}(e) \stackrel{\text{def}}{=} \emptyset$	for any other expression $e$

<sup>3</sup> The function `input()` does not return an integer in Python3, but this line can be replaced by `int(input())` to get an integer input in Python3.

<sup>4</sup> Operator `/` denotes integer floor division here, which would be `//` in Python3.

<sup>5</sup> The function `listinput()` does not belong to the builtins of Python3, but can easily be implemented by `list(map(int, input().split()))`.



## 3 Background

### 3.1 Abstract Interpretation

Our formal derivation and our implementation builds on the basis of *Abstract Interpretation*, a framework for approximating program behaviours [CC77]. Another detailed introduction can be found in the PhD thesis of a supervisor of this work [Urb15]. The key ingredients for every analysis built on Abstract Interpretation are described next:

**Abstract Domains:** One abstract domain (or several composed to one) allow representing the property of the program behaviour we are interested in. The representation is done by objects in some abstract universe, which holds partial information about the behaviour of an actual execution of the program. The concrete semantics describe this behaviour. We say *the concrete* when we are talking about actual execution behaviours. Note that, tracking all details of the concrete is not feasible for most programming languages. Therefore the decision what information to preserve in the abstract domains is crucial. A concretization function describes the meaning of the designed abstract objects in the concrete.

**Lattices:** Abstract domains should be based on a lattice structure and define a widening operator that is applied by the interpreter to guarantee termination of the analysis.

**Abstract Semantics:** Abstract semantics describe the effects of statements and expressions of the programming language on the abstract domain. The effects are described by *abstract transformers*, which are usually specific for the chosen abstract domain, and transform the abstract object, also called *state*. They imitate the effect of the concrete execution of the program on the abstract state without mapping the state into the concrete. We have to show that given abstract semantics are sound with respect to the concrete semantics.

**Interpreter:** The interpreter iterates over the program, applying the abstract transformers until a fixpoint is reached. It can be configured in several ways: some domains ask for a forward, others for a backward iteration, some for a combination. The interpreter may have to adjust the state additionally to how the abstract transformers change it, to reduce the number of iterations or prevent itself from running infinitely. This process is called *widening*.

There are two decisions to be made before designing the details of an abstract domain and its transformers, because abstract domains are usually sensitive to these decisions:

**Type of Approximation:** We can choose to under- or over-approximate the property of interest. It is the equivalent decision to prefer false-positives or false-negatives in the result set for which members we claim they satisfy the property of interest. An analysis that produces both false negatives and false positives is rarely desirable.

**Forward/Backward Analysis:** The interpreter can iterate through the program forwards or backwards. Some properties ask for a forward/backward analysis while others can be found both ways or sometimes a combination leads to the best results.

#### 3.1.1 Octagon Domain

For our more sophisticated analysis, we use well-known numerical domains, which support us with numerical information about variables. Specifically, the Octagon Domain **Oct** [Min06] provides interval bounds,  $b_l \leq x \leq b_u$ , and relational constraints on every pair of variables,  $x + y \leq c$ , for variables  $x, y \in \text{Var}_{\text{int}}, b_l, b_u, c \in \mathbb{Z}$ . These constraints are useful to compare

expressions containing variables. We need such expressions to store dynamic information about which parts of lists are used.

When octagons are modified, they need to be *closed* to infer implicit constraints. The author of [Min06] presents a *tight closure* algorithm for the case of integer variables, but updating all constraints has a worst-case runtime complexity of  $\mathcal{O}(n^4)$ , where  $n$  is the number of variables. Therefore we use the improved *tight closure* algorithm from [BHZ08], which has the runtime complexity  $\mathcal{O}(n^3)$ . Their approach is basically a combination of the classical Floyd-Warshall shortest path algorithm with a tightening step<sup>6</sup> and a step to achieve strong coherence. During the algorithm, checks for different types of consistency are done and if one fails, the octagon is marked as infeasible (set to the bottom element).

Since the octagon domain only supports scalar variables, we use *weak assignments* for the list variables. This means we treat them as a single variable. Relational constraints with lists are not possible in the octagon domain. We do not need them because we use the octagons for the purpose of supporting our list content analysis. We only implemented it for integer variables and division is not supported yet.

### 3.1.2 Interval Domain

The interval domain [CC77] is a simple numerical domain that abstracts integer variables by an interval. We do not use the interval domain directly, since it is a strictly weaker abstraction compared to the octagon domain we use. However, we use the interval calculus as a fallback whenever the octagon can not answer queries about numerical properties because the query parameters can not be put into a form usable with octagonal constraints.

When we need the interval domain fallback, the strategy is always the same: We first build an interval abstract state from the information available in the octagon domain, specifically the worst-case bounds of a variable. Then we evaluate an expression in the interval calculus, using these bounds for variables and approximating other parts in the expressions by reasonable interval abstractions. E.g., an observed input() within an expression is abstracted by  $[-\infty, \infty]$ .

## 3.2 Usage Property

In this section, we formalize our intuition of the property, that an input is *used*. We first give an informal definition of the property we want to approximate:

**Definition 2.** Variable  $x$  is *used* after a statement  $stmt$  (otherwise *unused*), if and only if, there is a  $x$ -definition-free path from a successor of  $stmt$  to a statement  $stmt'$  that reads the value of  $x$  and:

- $stmt'$  is a print statement, or
- $stmt'$  is a condition, that if satisfied, leads to a path that potentially has an effect on at least one output of the program, or
- $stmt'$  is an assignment statement that assigns to a variable  $y$  that is used after  $stmt'$ .

We suppose that each variable  $x \in \text{Var}$  has a (possibly infinite) domain  $D_x$  of valid values with an equivalence relation  $=$ . To ease explanation we assume that the fixed number of inputs and outputs are stored in specific variables  $\text{Var}_{\text{In}} \subseteq \text{Var}$  and  $\text{Var}_{\text{Out}} \subseteq \text{Var}$ , respectively. As

<sup>6</sup>The pseudo-code in Figure 2 of [BHZ08] has an error in the tightening step: After applying the floor function, the result should again be doubled, i.e., the line should read  $w[i, \bar{i}] := 2 * \text{floor}(w[i, \bar{i}]/2)$ ;

mentioned before, we can also determine the usage property for all variables, i.e.,  $\text{Var}_{\text{In}} = \text{Var}$ . The usage property is stored in a concrete environment  $\rho \in \mathcal{E} = \text{Var} \rightarrow \{U, N\}$ , which holds for every variable if it is used ( $U$ ) or unused ( $N$ ).

Let  $\Sigma$  be the set of all program states. A program *state*  $s \in \Sigma$  is a pair consisting of a program label and an environment that defines the values of the variables  $\text{Var}$  at the program control point designated by the label. We use the shorthand notation  $s(x_i)$  to refer to the value of  $x_i$  in the environment belonging to state  $s$ . Let  $\mathcal{I} \subseteq \Sigma$  be the set of *initial program states*, which are the states with a label designating the initial control point *before* the first statement of the program. Analogously,  $\Omega \subseteq \Sigma$  is the set of *final program states* with a label designating the final control point *after* the final statement. A state  $t$  is called *reachable* from state  $s$ , denoted by the predicate  $\text{reach}(s, t)$ , if there is a trace from  $s$  to  $t$ .

The variable  $x_i$  is considered used by a program if and only if there exist two initial program states  $s, s' \in \mathcal{I}$  with equal values for all input variables except the input value for  $x_i$  and there are two final program states  $t, t' \in \Omega$ , *reachable* from  $s, s'$ , respectively, with unequal values for at least one of the output variables. Formally,

$$\begin{aligned} \mathcal{E}(x_i) = U \equiv & \exists s, s' \in \mathcal{I}, t, t' \in \Omega : \text{reach}(s, t) \wedge \text{reach}(s', t') \\ & \wedge \forall x_j \in \text{Var}_{\text{In}} : (s(x_j) = s'(x_j) \leftrightarrow j \neq i) \wedge \exists o_k \in \text{Var}_{\text{Out}} : t(o_k) \neq t'(o_k). \end{aligned}$$

Similar properties have been stated and built analysis for in the context of secure information flow [BDR04].

**Backwards Analysis** As mentioned, the design of the abstract domains is sensitive to how the interpreter works. The property for data items we want to proof asks for a backward analysis, because the usage property propagates from the output statements/variables backwards to the inputs.

**Over-Approximation** We decided to target an analysis that produces an over-approximation of the set of used variables, i.e. we have no false-negatives, every variable denoted unused is definitely unused in the concrete.

## 4 Generic Functor Abstract Domains

We define reusable, composed domains in a generic way first. A functor<sup>7</sup> usually operates on at least one other domain  $\mathbf{D}$ , defined by the set of properties  $\mathcal{D}$ , induced lattice  $L_D = \langle \mathcal{D}, \sqsubseteq_D, \sqcup_D, \sqcap_D, \perp_D, \top_D \rangle$  and a (incomplete<sup>8</sup>) set of abstract transformers. We do not introduce a domain and its elements separately all the time, but use different fonts with the same letter to distinguish the domain ( $\mathbf{D}$ ), the set of abstract elements of the domain ( $\mathcal{D}$ ) and use subscripts for operators ( $\text{GET}_D, \sqsubseteq_D$ ) and top/bottom elements ( $\perp_D$ ). A functor can also take simple sets as (additional) parameters.

### 4.1 Abstract store domain

We often want to track certain values of entities (e.g., variables) in a program by an element of a lattice  $L_{\mathcal{D}}$  of some abstraction  $\mathbf{D}$ , which motivates the definition of an *abstract store domain*<sup>9</sup>. The abstract domain functor  $\mathbf{Store}(\mathcal{S}, \mathbf{D})$  takes a set  $\mathcal{S}$  of entities to be tracked and serving as keys into the store. Second, it takes any abstract domain  $\mathbf{D}$  and creates a new domain of abstract stores where each abstract store  $q$  is a mapping  $\mathcal{S} \rightarrow \mathcal{D}$ . An abstract store is also a complete lattice  $L_{\mathcal{S} \rightarrow \mathcal{D}} = \langle \mathcal{S} \rightarrow \mathcal{D}, \dot{\sqsubseteq}, \dot{\sqcup}, \dot{\sqcap}, \dot{\perp}, \dot{\top} \rangle$ , where the operators are defined by *pointwise lifting*:

$$\begin{aligned} q \dot{\sqsubseteq} q' &\stackrel{\text{def}}{=} \forall x \in \mathcal{S} : q(x) \sqsubseteq q'(x) \\ \dot{\sqcup} Q &\stackrel{\text{def}}{=} \lambda x. \sqcup \{q(x) \mid q \in Q\} \\ \dot{\sqcap} Q &\stackrel{\text{def}}{=} \lambda x. \sqcap \{q(x) \mid q \in Q\} \\ \dot{\perp} &\stackrel{\text{def}}{=} \lambda x. \perp \\ \dot{\top} &\stackrel{\text{def}}{=} \lambda x. \top \end{aligned}$$

We define a new formalism for updates on a store, extending the common store update  $q[v \leftarrow d]$ . We often want to update a store  $q$  to some value  $d$  at several keys  $\mathcal{V} \subseteq \mathcal{S}$  at once, and we may want to do it based on some condition  $c$  (which may depend on the key). Therefore we introduce

$$q \left[ v \xleftarrow{c(v)} d \mid v \in \mathcal{V} \right] \stackrel{\text{def}}{=} \dot{\sqcup} \left\{ q \left[ v \leftarrow \begin{cases} d & \text{if } c(v) \\ q(v) & \text{otherwise} \end{cases} \right] \mid v \in \mathcal{V} \right\}.$$

When we need a store for different *types* of entities  $\mathcal{S}_1, \mathcal{S}_2$ , like different types of variables (integers/lists), we may also want them to map to different domains  $\mathbf{D}_1, \mathbf{D}_2$ . Therefore we define the union of stores as:

$$\begin{aligned} \mathbf{Store}(\mathcal{S}_1, \mathbf{D}_2) \cup \mathbf{Store}(\mathcal{S}_2, \mathbf{D}_1) &\stackrel{\text{def}}{=} \mathbf{Store}(\mathcal{S}_1 \cup \mathcal{S}_2, \mathbf{D}_1 \cup \mathbf{D}_2) = Q \\ \text{s.t. for } q \in Q &: (v \in \mathcal{S}_1 \Rightarrow q(v) \in \mathbf{D}_1) \wedge (v \in \mathcal{S}_2 \Rightarrow q(v) \in \mathbf{D}_1) \end{aligned}$$

The definition of the union  $\mathbf{D}_1 \cup \mathbf{D}_2$  varies from case to case and has to provide a common partial order and a join/meet operator. By default we can assume the elements of  $\mathcal{D}_1$  and  $\mathcal{D}_2$  to be not comparable.

<sup>7</sup> The terminology of *functors* roots in the OCaml programming language and is explained in [CCL11].

<sup>8</sup> Some domains may only be useful to interpret a subset of a programs constructs.

<sup>9</sup> We call it *abstract store* as an opposite to *symbolic store*.

## 4.2 Abstract stack domain

We define a generic abstract domain functor **Stack**( $\Gamma$ ). Each instance can be described by a tuple  $\langle \Gamma, \text{push\_element}, \text{pop\_element} \rangle$ , where  $\Gamma$  is the set of valid elements on the stack (stack alphabet) and  $\text{push\_element} : \Gamma \rightarrow \Gamma$ ,  $\text{pop\_element} : \Gamma \times \Gamma \rightarrow \Gamma$  are functions invoked when we push or pop an element, respectively. These two functions allow specific behaviour to be defined for a concrete stack instance.

Formally, a stack is defined as an (ordered) sequence  $s \in \Gamma^*$  of stack elements, which form a complete lattice  $L_\Gamma$ . With  $\epsilon \in \Gamma^*$ , we denote the empty stack. A stack supports the following operations:

Push an element onto the stack:

$$\begin{aligned} \text{PUSH} : \Gamma^k &\longrightarrow \Gamma^{k+1} \\ s &\longmapsto \underbrace{\langle s_1, \dots, s_k \rangle}_s, \text{push\_element}(s_k) \end{aligned}$$

Pop an element from the stack:

$$\begin{aligned} \text{POP} : \Gamma^{k+1} &\longrightarrow \Gamma^k \quad (k \geq 1) \\ \langle s_1, \dots, s_{k-1}, s_k, s_{k+1} \rangle &\longmapsto \langle s_1, \dots, s_{k-1}, \text{pop\_element}(s_k, s_{k+1}) \rangle \end{aligned}$$

The operators are defined by element-wise lifting of the operators defined for the domain of the elements on the stack. For two stacks  $s, s' \in \Gamma^*$  with matching length, we define:

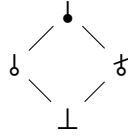
$$\begin{aligned} \langle s_1, \dots, s_k \rangle \bar{\sqsubseteq} \langle s'_1, \dots, s'_k \rangle &\stackrel{\text{def}}{=} \forall s_i. s_i \sqsubseteq s'_i \\ \langle s_1, \dots, s_k \rangle \bar{\sqcup} \langle s'_1, \dots, s'_k \rangle &\stackrel{\text{def}}{=} \langle s_1 \sqcup s'_1, \dots, s_k \sqcup s'_k \rangle \\ \langle s_1, \dots, s_k \rangle \bar{\sqcap} \langle s'_1, \dots, s'_k \rangle &\stackrel{\text{def}}{=} \langle s_1 \sqcap s'_1, \dots, s_k \sqcap s'_k \rangle \\ \bar{\perp} &\stackrel{\text{def}}{=} \langle \perp, \dots, \perp \rangle \\ \bar{\top} &\stackrel{\text{def}}{=} \langle \top, \dots, \top \rangle \end{aligned}$$

## 5 Usage Analysis for Integer Variables

### 5.1 Usage Abstract Domain

We introduce the *usage abstract domain*  $\mathbf{U}$  which tracks for a scalar variable if it is used at a specific program point. Recall the intuitive Definition 2 of the usage property. One of the criterias to let a variable become used is that the variable appears in a condition that leads to a path that potentially has an effect on an output. Some paths only traverse straight-line code, others enter and exit conditionals and loops. To (implicitly) distinguish different types of paths, we need a notion of a *scope*. Eventhough our notion of scopes is inspired by *variable scopes*, they are not the same. Our scopes track the possibility of a branching condition in the control flow graph to have an effect on the output and are derived from the structure of the control flow graph. Since conditionals can be nested in SimplePython, scopes can also be nested, but all nested scopes must be fully contained in their parent scopes.

Let  $L_{\mathcal{U}} = \langle \mathcal{U}, \sqsubseteq, \sqcup, \sqcap, \perp, \top \rangle$  be a complete lattice with the following Hasse-Diagram:



The elements  $\mathcal{U} = \{\top, \downarrow, \leftarrow, \perp\}$  of this lattice have the following intuitive meaning:

- $\top = \bullet$  : Used in this scope or in a deeper nested scope.
- $\downarrow$  : Used in an outer scope.
- $\leftarrow$  : Used in an outer scope and overridden in this scope.
- $\perp$  : Not used.

**Lattice Operators** The binary operators  $\sqsubseteq, \sqcup, \sqcap$  are defined through the Hasse-Diagram above.

**Additional Operators** Descending into a nested scope (in a backwards fashion) requires to lower the abstract elements accordingly. Let  $\text{DESCEND} : \mathcal{U} \rightarrow \mathcal{U}$ ,

$$\frac{\text{DESCEND}(u)}{u} \begin{array}{c|cccc} \downarrow & \downarrow & \leftarrow & \perp & \perp \\ \bullet & \downarrow & \leftarrow & \perp & \perp \end{array} \quad (1)$$

Exiting a nested scope (at its start because we do a backwards analysis) requires to combine the representation of the nested and the containing scope. Let  $\text{COMBINE} : \mathcal{U} \times \mathcal{U} \rightarrow \mathcal{U}$ ,

$$\frac{\text{COMBINE}(u, u')}{u} \begin{array}{c|cccccccccccccccc} \perp & \downarrow & \leftarrow & \bullet & \downarrow & \downarrow & \leftarrow & \downarrow & \leftarrow & \downarrow & \leftarrow & \downarrow & \downarrow & \leftarrow & \downarrow \\ \perp & \downarrow & \leftarrow & \bullet & \perp & \downarrow & \leftarrow & \bullet & \perp & \downarrow & \leftarrow & \bullet & \perp & \downarrow & \leftarrow \\ \perp & \perp & \perp & \perp & \downarrow & \downarrow & \downarrow & \downarrow & \leftarrow & \leftarrow & \leftarrow & \leftarrow & \bullet & \bullet & \bullet & \bullet \end{array} \quad (2)$$

### 5.2 Usage Stores

We want to track all variables in a program with a value of the complete lattice  $L_{\mathcal{U}}$ , so we use abstract stores  $q \in Q$ ,  $q : \text{Var} \rightarrow \mathcal{U}$ . This is an instantiation  $\mathbf{Q} = \text{Store}(\text{Var}, \mathbf{U})$  of the abstract store functor, hence the usage stores inherit the structure and they build also a complete lattice  $L_Q = \langle \text{Var} \rightarrow \mathcal{U}, \sqsubseteq, \sqcup, \sqcap, \perp, \top \rangle$ .

## Information Flow Operators

**Implicit Flow** Usage stores support the unary operator  $\text{FILTER}_Q[[e]] : Q \rightarrow Q$  used to update the abstract store on a branching condition in the control flow graph:

$$\text{FILTER}_Q[[e]]q \stackrel{\text{def}}{=} \bigsqcup \left\{ q \left[ x \leftarrow \begin{cases} \downarrow & \text{if } \exists y \in \text{Var} : q(y) \in \{\downarrow, \dagger\} \\ q(x) & \text{otherwise} \end{cases} \right] \mid x \in \text{opd}(e) \cap \text{Var} \right\}$$

**Explicit Flow** The operator  $\text{USE}_Q[[x = e]] : Q \rightarrow Q$  allows transitive usage of variables to be reflected in the abstract store:

$$\text{USE}_Q[[x = e]]q \stackrel{\text{def}}{=} \begin{cases} \bigsqcup \{q[y \leftarrow \downarrow] \mid y \in \text{opd}(e) \cap \text{Var}\} & \text{if } q(x) \in \{\downarrow, \downarrow\} \\ q & \text{otherwise} \end{cases}$$

**Intermitted Flow** The operator  $\text{KILL}_Q[[x = e]] : Q \rightarrow Q$  permits the potential degradation of the abstract state of a variable at an overwriting assignment:

$$\text{KILL}_Q[[x = e]]q \stackrel{\text{def}}{=} q \left[ x \leftarrow \frac{q(x) \in \{\downarrow, \downarrow\} \wedge x \notin \text{opd}(e)}{\dagger} \right]$$

## Partial Abstract Semantics

$$\begin{aligned} [[x = e]]q &\stackrel{\text{def}}{=} (\text{KILL}_Q[[x = e]] \circ \text{USE}_Q[[x = e]])q && \text{(order of application is crucial<sup>10</sup>)} \\ [[\text{print}(e)]]q &\stackrel{\text{def}}{=} \bigsqcup \{q[x \leftarrow \downarrow] \mid x \in \text{opd}(e) \cap \text{Var}\} \end{aligned} \quad (3)$$

## 5.3 Usage Stacks

We track all variables in a program with a usage property  $u \in \mathcal{U}$  of the domain  $\mathbf{U}$ . The elements  $u$  implicitly track on what type of execution path they are rooted in. Therefore, they potentially change when the interpreter enters or exits a scope. To make the usage property of variables dependent on scopes, we use a stack instantiation  $\mathbf{S} = \mathbf{Stack}(\mathbf{Q})$  with the configuration tuple  $\langle Q, \text{DESCEND}, \text{COMBINE} \rangle$ .

**Lifting Operators** We naturally lift the FILTER-operator on stores to stacks by simply modifying the top most element of a stack  $s = \langle s_1, \dots, s_k \rangle$ :

$$\text{FILTER}_S[[e]]s \stackrel{\text{def}}{=} \text{FILTER}_Q[[e]]s_k$$

**Abstract Transformers** For what follows, we use  $s, t \in S$  to denote abstract states.

<sup>10</sup> To avoid ambiguity:  $(f \circ g)x \stackrel{\text{def}}{=} f(g(x))$

$$\begin{aligned}
\llbracket \text{pass} \rrbracket s &\stackrel{\text{def}}{=} s \\
\llbracket x = e \rrbracket s &\stackrel{\text{def}}{=} \langle s_1, \dots, s_{k-1}, \llbracket x = e \rrbracket s_k \rangle \\
\llbracket \text{print}(e) \rrbracket s &\stackrel{\text{def}}{=} \langle s_1, \dots, s_{k-1}, \llbracket \text{print}(e) \rrbracket s_k \rangle \\
\llbracket \text{if } b : stmt_1 \text{ else } stmt_2 \rrbracket s &\stackrel{\text{def}}{=} \text{POP}((\text{FILTER}_S \llbracket b \rrbracket \circ \llbracket stmt_1 \rrbracket)(\text{PUSH}(s))) \\
&\quad \sqcup \text{POP}((\text{FILTER}_S \llbracket \text{not } b \rrbracket \circ \llbracket stmt_2 \rrbracket)(\text{PUSH}(s))) \\
\llbracket \text{while } b : stmt \rrbracket s &\stackrel{\text{def}}{=} \text{lfp}_s [\lambda t. \text{FILTER}_S \llbracket \text{not } b \rrbracket t \\
&\quad \sqcup \text{POP}((\text{FILTER}_S \llbracket b \rrbracket \circ \llbracket stmt \rrbracket)(\text{PUSH}(t)))]
\end{aligned} \tag{4}$$

## 5.4 Concretization

The concretization  $\gamma_S$  of a stack depends on the concretization  $\gamma_Q$  of its element domain. In our analysis the purpose of the usage stacks is to remember usage properties across scopes, so for the concrete only the top frame is relevant. For a stack  $s$  of length  $m$  we have:

$$\gamma_S(s) \stackrel{\text{def}}{=} \gamma_Q(s_m) \tag{5}$$

For the concretization of the usage store we reuse the usage property environment  $\rho \in \mathcal{E} = \text{Var} \rightarrow \{U, N\}$  from Section 3.2. We define a function that returns the set of used variables in the concrete:

$$\begin{aligned}
\text{USED} : \mathcal{E} &\longrightarrow \wp(\text{Var}) \\
\rho &\longmapsto \{o \in \text{Var} \mid \rho(o) = U\}
\end{aligned}$$

The concretization of usage stores creates an environment with the same signature as a concrete environment:

$$\begin{aligned}
\gamma_Q : Q &\longrightarrow \mathcal{E} \\
q &\longmapsto \lambda x. \begin{cases} U & \text{if } q(x) \in \{\downarrow, \delta\} \\ N & \text{otherwise} \end{cases}
\end{aligned}$$

**Theorem 1.** The usage analysis is an over-approximation of the set of definitely used variables, i.e., the abstract state  $s$  and the concrete environment  $\rho$  at every program point satisfy:

$$\text{USED}(\rho) \subseteq \text{USED}(\gamma_S(s))$$

## 5.5 Example

Code 2 shows a rewrite of our program for counting passing grades (Code 1). It uses a fixed number of input grades read with `input()`. Bug A on line 8 leads to `count` not being updated upon the check that the math grade is passing. Bug B on line 11 leads to an update of `count` but, because of an erroneous check, involving the wrong variable `physics`. On line 6, our analysis correctly detects `count`  $\rightarrow \downarrow$ , `math`  $\rightarrow \perp$ , `physics`  $\rightarrow \downarrow$ , `history`  $\rightarrow \perp$ .



```

1 math = input()
  physics = input()
  history = input()

5 count = 0

  if math >= 4:
    math += 1 # Bug A: should be 'count += 1'
  if physics >= 4:
10 count += 1
  if physics >= 4: # Bug B: should be 'history >= 4'
    count += 1

print(count)

```

Code 2: Program with two bugs detectable with the  $\mathbf{U}$  abstract domain. On line 6, the analysis outputs  $count \rightarrow \downarrow, math \rightarrow \perp, physics \rightarrow \downarrow, history \rightarrow \perp$ .

## 6 List Start Usage Analysis

In this section we explore the *list start usage domain*, denoted by  $\mathbf{H}$ , to track the content of a list with respect to the usage property. More precisely, our domain will track how much from the start of a list is used, used outside of the current scope, overwritten or definitely not used, corresponding to the elements of  $\mathcal{U}$ . Since we want to track these properties *from the start* of the list, we only have to remember an upper limit until where the content of a list satisfies a property from  $\mathcal{U}$ . If we implicitly assume that the whole list is unused if no other property is stated, we only have to store the limits of the other properties  $\downarrow, \uparrow, \cancel{\downarrow} \in \mathcal{U}$ . Thus, an element  $h \in \mathbf{H}$ , the set of all valid elements in the domain  $\mathbf{H}$ , is a triple:

$$h = \underbrace{\langle h_{\downarrow}, h_{\uparrow}, h_{\cancel{\downarrow}} \rangle}_{\text{indexed by } u \in \mathcal{U}} \in (\mathbb{N}_0 \cup \{\infty\})^3$$

The actual usage property at an index  $i$  of a list with abstraction  $h$  is the least upper bound of the properties *reaching* that index:

$$\text{GET}_{\mathbf{H}}[h](i) = \text{GET}_{\mathbf{H}}[\langle h_{\downarrow}, h_{\uparrow}, h_{\cancel{\downarrow}} \rangle](i) \stackrel{\text{def}}{=} \bigsqcup (\{u \in \mathcal{U} \mid i < h_u\} \cup \{\perp\}) \quad (6)$$

**Example 1.** Let a list  $l$  be abstracted by  $h = \langle 0, 3, 5 \rangle$ . This reads as the list items  $l[0], l[1], l[2]$  are considered used ( $\downarrow$ ) and  $l[3], l[4]$  are considered overwritten ( $\cancel{\downarrow}$ ). The rest of the list (of unknown length) is definitely unused. Figure 1 show a visualization of this example.

Whenever the limits  $h_{\uparrow}$  and  $h_{\cancel{\downarrow}}$  disagree in  $h$ , i.e., when both are greater than 0, we say that  $h$  is not *coherent*. It is also not coherent when  $h_{\uparrow}$  and  $h_{\cancel{\downarrow}}$  have no effect since they are smaller than  $h_{\downarrow}$ .

**Example 2.** Both  $h = \langle w, 3, 5 \rangle$  and  $h' = \langle 0, 10, 5 \rangle$  are not coherent.

**Definition 3.**  $h = \langle h_{\downarrow}, h_{\uparrow}, h_{\cancel{\downarrow}} \rangle$  is *coherent* if  $(h_{\downarrow} = 0 \vee h_{\cancel{\downarrow}} = 0) \wedge h_{\downarrow} > \max(h_{\uparrow}, h_{\cancel{\downarrow}})$ . An alternative definition is that we require  $\{u \in \mathcal{U} \mid i < h_u\}$  from (6) to contain at most one element.

Even though  $\text{GET}_{\mathbf{H}}$  resolves disagreeing indices to  $\downarrow$  (since  $\downarrow$  and  $\cancel{\downarrow}$  are incomparable), we would like to have the elements always in a coherent state. Therefore we define the closure

$$\overset{\circ}{h} = \text{CLOSE}[h] = \text{CLOSE}[\langle h_{\downarrow}, h_{\uparrow}, h_{\cancel{\downarrow}} \rangle] \stackrel{\text{def}}{=} \left\langle \begin{cases} h_{\downarrow} & \text{if } h_{\downarrow} > h_{\cancel{\downarrow}}, h_{\uparrow} \\ 0 & \text{otherwise} \end{cases}, \max(h_{\uparrow}, \min(h_{\downarrow}, h_{\cancel{\downarrow}})), \begin{cases} h_{\cancel{\downarrow}} & \text{if } h_{\cancel{\downarrow}} > h_{\downarrow}, h_{\uparrow} \\ 0 & \text{otherwise} \end{cases} \right\rangle$$

**Lemma 1.** The closure  $\text{CLOSE}[h]$  of any abstract element  $h \in \mathbf{H}$  is *coherent*.

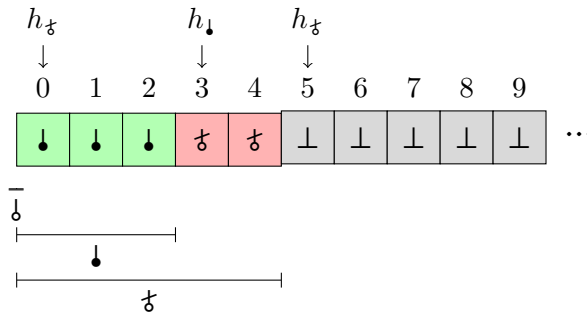


Figure 1: Visualization of the abstraction  $h = \langle 0, 3, 5 \rangle$  in the list start usage domain  $\mathbf{H}$ .

Analogously, we define an operator to represent an update of usage at a specific index with the most precise abstraction:

$$\begin{aligned} \text{SET}_{\mathbf{H}}[h](i, u) &= \text{SET}_{\mathbf{H}}[\langle h_{\downarrow}, h_{\downarrow}, h_{\ddagger} \rangle](i, u) \stackrel{\text{def}}{=} \text{CLOSE}[h'], \\ \text{where } \forall \mu \in \mathcal{U} : h'_{\mu} &= \begin{cases} \max(h_u, i) & \text{if } u = \mu \\ h_{\mu} & \text{otherwise} \end{cases} \end{aligned} \quad (7)$$

**Lattice Operators** We define the complete lattice  $L_{\mathbf{H}} = \langle \mathbf{H}, \sqsubseteq_{\mathbf{H}}, \sqcup_{\mathbf{H}}, \sqcap_{\mathbf{H}}, \perp_{\mathbf{H}}, \top_{\mathbf{H}} \rangle$ . The operators produce the same result as if we would apply them at each index individually and then find the most precise abstraction as a triple  $h \in \mathbf{H}$ :

$$\begin{aligned} h \sqsubseteq_{\mathbf{H}} h' &\stackrel{\text{def}}{=} \forall u \in \{\downarrow, \downarrow, \ddagger\} : \hat{h}_u \leq \hat{h}'_u \\ h \sqcup_{\mathbf{H}} h' &\stackrel{\text{def}}{=} \text{CLOSE} \left[ \left\langle \max(\hat{h}_{\downarrow}, \hat{h}'_{\downarrow}), \max(\hat{h}_{\downarrow}, \hat{h}'_{\downarrow}), \max(\hat{h}_{\ddagger}, \hat{h}'_{\ddagger}) \right\rangle \right] \\ h \sqcap_{\mathbf{H}} h' &\stackrel{\text{def}}{=} \text{CLOSE} \left[ \left\langle \min(h_{\downarrow}, h'_{\downarrow}), \min(h_{\downarrow}, h'_{\downarrow}), \min(h_{\ddagger}, h'_{\ddagger}) \right\rangle \right] \\ \perp_{\mathbf{H}} &\stackrel{\text{def}}{=} \langle 0, 0, 0 \rangle \\ \top_{\mathbf{H}} &\stackrel{\text{def}}{=} \langle 0, \infty, 0 \rangle \end{aligned}$$

The partial order and the meet operator require closed arguments for greatest possible precision and they all ensure that the result is again closed. The widening operator each of limits individually to infinity when their are growing and closes the result:

$$\begin{aligned} h \nabla h' &\stackrel{\text{def}}{=} \text{CLOSE}[h^{\nabla}] \\ \text{where } \forall \mu \in \mathcal{U} : h^{\nabla}_{\mu} &= \begin{cases} \infty & \text{if } h_{\mu} < h'_{\mu} \\ h'_{\mu} & \text{otherwise} \end{cases} \end{aligned}$$

**Lifting of Usage Abstract Domain Operators** We lift the operators defined on  $\mathcal{U}$  in  $\mathbf{U}$  to the domain  $\mathbf{H}$ . Recall the definition of  $\text{DESCEND}$  in (1). The result of  $\text{DESCEND}$  is bascially  $\downarrow$  if the argument was either  $\downarrow$  or  $\downarrow$ , and  $\perp$  otherwise. Thus we have to reset the limits of an abstraction  $h$  of domain  $\mathbf{H}$  at both  $\downarrow$  and  $\ddagger$ :

$$\text{DESCEND}_{\mathbf{H}}(h) \stackrel{\text{def}}{=} \langle \max(h_{\downarrow}, h_{\downarrow}), 0, 0 \rangle$$

Note that, for brevity,  $\text{DESCEND}_{\mathbf{H}}$  is given here *specifically* for the definition in (1). The lifting of the  $\text{COMBINE}$  operator to  $\text{COMBINE}_{\mathbf{H}}$  is fully generic without assuming anything about the definition of  $\text{COMBINE}$ <sup>11</sup>. Intuitively, the lifting finds all overlapping ranges and combines their properties with  $\text{COMBINE}$  and then maps this back to the most precise representation possible in  $\mathbf{H}$ . To make the presentation easier, we assume  $h_{\perp} = \infty$  is also defined for every  $h \in \mathbf{H}$ .

$$\begin{aligned} \text{COMBINE}_{\mathbf{H}}(h, h') &\stackrel{\text{def}}{=} \text{CLOSE} \left[ \left\langle \max \left\{ \min(h_u, h_{u'}) \mid u, u' \in \mathcal{U} \wedge \text{COMBINE}(u, u') = \downarrow \right\}, \right. \right. \\ &\quad \max \left\{ \min(h_u, h_{u'}) \mid u, u' \in \mathcal{U} \wedge \text{COMBINE}(u, u') = \downarrow \right\}, \\ &\quad \left. \max \left\{ \underbrace{\min(h_u, h_{u'})}_{\text{upper limit of overlapping range}} \mid \underbrace{u, u' \in \mathcal{U} \wedge \text{COMBINE}(u, u') = \ddagger}_{\text{all combinations resulting in } \ddagger} \right\} \right] \end{aligned}$$

<sup>11</sup>The implementation uses a generic way to realize this function, but is more efficient by not enumerating all the possible combinations  $u, u' \in \mathcal{U}$  but using a scanline-algorithm.

## 6.1 Combining List Start Domain and Usage Domain

We track the scalar (integer) variables of the program with members of the usage domain  $\mathbf{U}$  as in the integer-only programs from Section 5.2, and we track list variables with the list start domain  $\mathbf{H}$ . To combine these domains we use abstract stores  $q \in Q = Q_{\text{int}} \cup Q_{\text{list}}$ , where  $Q_{\text{int}} = \mathbb{V}\text{ar}_{\text{int}} \rightarrow \mathcal{U}$  and  $Q_{\text{list}} = \mathbb{V}\text{ar}_{\text{list}} \rightarrow \mathbf{H}$ .

**Information Flow Operators** We extend usage stores to be able to filter expressions containing both scalar variable accesses and list index expressions. We define the condition when the operands in a filtered expression are considered used, i.e., when a store is *dirty*, as follows:

$$\text{DIRTY}(q) \stackrel{\text{def}}{=} \exists y \in \mathbb{V}\text{ar}_{\text{int}} : q(y) \in \{\downarrow, \dagger\} \vee \exists l \in \mathbb{V}\text{ar}_{\text{list}} : q(l)_{\downarrow} > 0$$

**Implicit Flow** We account for implicit flow by checking if a store is dirty and updating the store accordingly:

$$\begin{aligned} \text{FILTER}_Q \llbracket e \rrbracket q &\stackrel{\text{def}}{=} q \left[ x \xleftarrow{\text{DIRTY}(q)} \downarrow \mid x \in \text{opd}(e) \cap \mathbb{V}\text{ar}_{\text{int}} \right] \\ &\quad \dot{\downarrow} q \left[ l \xleftarrow{\text{DIRTY}(q)} \text{SET}_H[q(l)](i; \downarrow) \mid l[i] \in \text{opd}(e) \cap \mathbb{I}\text{ndex} \right] \end{aligned}$$

**Explicit Flow** We account for explicit flow by checking if left side of assignment is currently used:

$$\begin{aligned} \text{USE}_Q \llbracket x = e \rrbracket q &\stackrel{\text{def}}{=} q \left[ y \xleftarrow{q(x) \in \{\downarrow, \downarrow\}} \downarrow \mid y \in \text{opd}(e) \cap \mathbb{V}\text{ar}_{\text{int}} \right] \\ &\quad \dot{\downarrow} q \left[ l \xleftarrow{q(x) \in \{\downarrow, \downarrow\}} \text{SET}_H[q(l)](i; \downarrow) \mid l[i] \in \text{opd}(e) \cap \mathbb{I}\text{ndex} \right] \end{aligned}$$

The transformer  $\text{USE}_Q \llbracket l[i] = e \rrbracket q$  is similar except that the condition for updating the store is  $\text{GET}_H[q(l)](i) \in \{\downarrow, \downarrow\}$ . When assigning to a list variable, we propagate the used property naturally to the right side of the assignment, whether it is another list variable or a list display expression:

$$\begin{aligned} \text{USE}_Q \llbracket l = l' \rrbracket q &\stackrel{\text{def}}{=} q \left[ l' \leftarrow q(l) \left[ p_i \xleftarrow{l \neq l' \wedge p_i \in \{\downarrow, \downarrow\}} \downarrow \mid i \in \{0, \dots, \text{len}(l)\} \right] \right] \\ \text{USE}_Q \llbracket l = [t_0, \dots, t_{m-1}] \rrbracket q &\stackrel{\text{def}}{=} q \left[ x \xleftarrow{\text{GET}_H[q(l)](k) \in \{\downarrow, \downarrow\}} \downarrow \mid x \in \text{opd}(t_k), k \in \{0, \dots, m-1\} \right] \end{aligned}$$

**Intermitted Flow** The intermitted flow for an assignment to a scalar variable is as before:

$$\text{KILL}_Q \llbracket x = e \rrbracket q \stackrel{\text{def}}{=} q \left[ x \xleftarrow{q(x) \in \{\downarrow, \downarrow\} \wedge x \notin \text{opd}(e)} \dagger \right]$$

We add the definitions for the assignment to list variables:

$$\begin{aligned} \text{KILL}_Q \llbracket l = l' \rrbracket q &\stackrel{\text{def}}{=} q \left[ l \leftarrow q(l) \left[ p_i \xleftarrow{l \neq l' \wedge p_i \in \{\downarrow, \downarrow\}} \dagger \mid i \in \{0, \dots, \text{len}(l)\} \right] \right] \\ \text{KILL}_Q \llbracket l = [t_0, \dots, t_{m-1}] \rrbracket q &\stackrel{\text{def}}{=} q \left[ l \leftarrow q(l) \left[ p_i \xleftarrow{p_i \in \{\downarrow, \downarrow\}} \dagger \mid i \in \{0, \dots, \text{len}(l)\} \right] \right] \end{aligned}$$

To determine the intermitted flow for an assignment to a list index  $l[i] = e$ , we now take advantage of the simplification that  $i$  can only be a constant. This makes the check if the index  $i$  appears in the right side as easy as checking for an occurrence of  $l[i]$ <sup>12</sup>

$$\text{KILL}_Q[\![l[i] = e]\!]q \stackrel{\text{def}}{=} q \left[ l \leftarrow \begin{cases} \text{SET}_H[q(l)](i, \sharp) & \text{if } l[i] \in \text{opd}(e) \cap \text{Var}_{\text{list}} \\ q(l) & \text{otherwise} \end{cases} \right]$$

Finally we can combine this operators to adapt the partial abstract semantics given in (3) for int/list stores:

$$\begin{aligned} \llbracket \psi = \omega \rrbracket q &\stackrel{\text{def}}{=} (\text{KILL}_Q[\![\psi = \omega]\!] \circ \text{USE}_Q[\![\psi = \omega]\!]) q \\ \llbracket \text{print}(e) \rrbracket q &\stackrel{\text{def}}{=} q [x \leftarrow \downarrow \mid x \in \text{opd}(e) \cap \text{Var}_{\text{int}}] \\ &\quad \dot{\cup} q [l \leftarrow \text{SET}_H[q(l)](i; \downarrow) \mid l[i] \in \text{opd}(e) \cap \text{Index}] \\ &\quad \dot{\cup} q [l \leftarrow \top_H \mid l[i] \in \text{opd}(e) \cap \text{Var}_{\text{list}}] \end{aligned}$$

**Abstract Transformers** The abstract transformers stay the same as in (4) except that the assignment has support for lists:

$$\llbracket \psi = \omega \rrbracket s \stackrel{\text{def}}{=} \langle s_1, \dots, s_{k-1}, \llbracket \psi = \omega \rrbracket s_k \rangle$$

## 6.2 Concretization

We focus on the concretization for the list start usage domain  $\mathbf{H}$  here, the concretization of the combination with  $\mathbf{U}$  is straight forward. For the concretization  $\gamma_{Q, \mathbf{H}}$  of the part of the store that tracks list variables  $Q_{\text{list}}$ , we assume that each list element is identified by a distinct symbol  $\sigma \in \text{Sym}$  in the concrete. Then the concrete environment is a mapping from symbols the concrete properties used/unused,  $\rho_{\text{list}} \in \mathcal{E}_{\text{sym}} = \text{Sym} \rightarrow \{U, N\}$ . The concretization is given by

$$\begin{aligned} \gamma_{Q, \mathbf{H}} : Q &\longrightarrow \mathcal{E}_{\text{sym}} \\ q &\longmapsto \lambda \sigma. \begin{cases} U & \text{if } \text{GET}_H[q(l^\sigma)](i^\sigma) \in \{\downarrow, \downarrow\}, \\ N & \text{otherwise,} \end{cases} \end{aligned}$$

where  $l^\sigma$  is the list the symbol  $\sigma$  belongs to and  $i^\sigma$  is the corresponding index in the list.

We define the concretization of the stack  $s$  of length  $m$  in the same fashion as in (5):

$$\gamma_{S, \mathbf{H}}(s) \stackrel{\text{def}}{=} \gamma_{Q, \mathbf{H}}(s_m)$$

**Theorem 2.** The list start usage analysis includes an over-approximation of the set of list items definitely used, i.e., the abstract state  $s$  and the concrete environment  $\rho_{\text{sym}}$  at every program point satisfy:

$$\text{USED}(\rho_{\text{sym}}) \subseteq \text{USED}(\gamma_{S, \mathbf{H}}(s))$$

<sup>12</sup> When we have access to an variable environment that can evaluate an expression  $a$  to an interval, we can also allow index expressions for  $l[a]$ . Refer to Section 7.3 for a more sophisticated domain accessing the variable environment.

### 6.3 Example

Code 3 shows a program to demonstrate the presented domain. The program takes two inputs, which are compiled with some constants to a list on line 4. On lines 8 to 11 we perform some random accesses to the list, using the accessed values in the printed sum. Since  $\mathbf{H}$  can only track usage properties from the start of the list, we receive on line 5 that *list1* is used up to index 4 ( $h_{\downarrow}$  is exclusive):

$$list1 \rightarrow \langle 0, 5, 0 \rangle \quad (8)$$

On line 3 we get the correct result  $x \rightarrow \downarrow, y \rightarrow \perp$ . However, if we would remove line 9 (marker A), the analysis still outputs the same result even though  $x$  is no longer used, thus is imprecise.

```
1 x = int(input())
  y = int(input())

  list1 = [1, x, 2, 3, 5, 8, y]
5
  sum = 0
  # some random accesses to list1
  sum += list1[2]
  sum += list1[1] # Marker A
10 sum += list1[4]
    sum += list1[0]
    print(sum)
```

Code 3: Program demonstrating the list start usage domain  $\mathbf{H}$ .

## 7 List Content Usage Analysis

We first have a closer look at the ingredients, before combining them to segmentations.

### 7.1 Ingredients

#### 7.1.1 Variable Environment Abstract Domain

In the context of a list content analysis the variables referring to lists are abstracted by a segmentation, so the variable environment must at least cover the abstractions for the set of scalar variables  $\mathbb{V}\text{ar}_{\text{scalar}}$  in the program. This abstraction is known under the name *variable environment abstract domain*, denoted  $\mathbf{V}$ , and should comply with the interface of a  $\mathbf{Store}(\mathbb{V}\text{ar}_{\text{scalar}}, \mathbf{D})$ . Alternatively, we can use an extension that supports at least the retrieval of some variable-specific properties, like the *Interval Domain* provides worst-case bounds for each scalar variable. We will build on a even more powerful, relational variable environment domain (*Octagon Domain*) and show some interesting modifications to operators on segmentations, as opposed to the proposal of [CCL11] to keep the segmentation fully generic and not use any relational constraints on scalar variables. We still keep it as generic as possible and only require  $\mathbf{V}$  to provide a partial order  $\sqsubseteq_V$  that is capable of answering  $v + c \sqsubseteq_V v' + c'$ , where  $v, v' \in \mathbb{V}\text{ar}$  and  $c, c' \in \mathbb{Z}$ .

#### 7.1.2 Expression Abstract Domain

The expressions used for bounds have to be in some normal form. We will use the simple normal form  $(v + c) \in E$ <sup>13</sup>,  $v \in \mathbb{V}\text{ar}_{\text{int}} \cup \{v_0\}$  where  $c \in \mathbb{Z}$  and  $v_0$  is considered to be always equal to 0 (to represent a constant with  $(v_0 + c)$ ).

**Operators** When we use expression in segmentation limits, we have to be able to decide about their equality. This can be done on a syntactic level only, by comparing the normal forms for equal variables and constants:

$$e = e' :\Leftrightarrow v = v' \wedge c = c' \quad (\text{syntactic equality})$$

But we can also use relational constraints from the variable environment abstract domain  $\mathbf{V}$  to derive expression equality:

$$e \stackrel{\Delta}{\sqsubseteq}_E e' :\Leftrightarrow e \sqsubseteq_E e' \wedge e' \sqsubseteq_E e \quad (\text{relational equality})$$

$$\text{where } e \sqsubseteq_E e' :\Leftrightarrow v - v' \sqsubseteq_V c' - c \quad (\text{partial order})$$

Together with the least upper bound  $\sqcup_E$  and greatest lower bound  $\sqcap_E$  defined as usual, a bottom element  $\perp_E$ , representing unreachability, and a top element  $\top_E$ , representing expressions that can not be put in normal form, we can define the lattice  $L_E = \langle E, \sqsubseteq_E, \sqcup_E, \sqcap_E, \perp_E, \top_E \rangle$ .

#### 7.1.3 Segment Limit Abstract Domain

Each segment limit in  $\mathbf{L}(\mathbf{E}_V)$  is a non-empty set of expressions from  $\mathbf{E}_V$ . All the bounds inside a limit set are supposed to be pairwise equals in the concrete. For limits  $l, l'$ , we lift the equality

<sup>13</sup>The parenthesis are not shown in the examples, only when needed for clarity.

definitions from the expression domain:

$$\begin{aligned}
l = l' &:\Leftrightarrow \forall b \in l : b \in l' \wedge \forall b' \in l' : b' \in l && \text{(syntactic equality } \triangleq \text{ set equality of } l, l') \\
l \sqsubseteq_L l' &:\Leftrightarrow \exists b \in l, b' \in l' : b \sqsubseteq_E b' && \text{(partial order)} \\
l \triangleq_L l' &:\Leftrightarrow \exists b \in l, b' \in l' : b \triangleq_E b' && \text{(relational equality)}
\end{aligned}$$

Eventhough we could add a top and bottom element to complete the lattice of limits, it is not necessary since we will remove limits whenever a bound turns into  $\perp_E$  or  $\top_E$ .

### 7.1.4 Predicate Abstract Domain

We will use the simplest possible abstraction, where the concrete indices of the list covered by an abstract segment is abstracted by a single element, called predicate, from a domain  $\mathbf{P}$ . We assume the lattice  $L_P = \langle P, \sqsubseteq_P, \sqcup_P, \sqcap_P, \perp_P, \top_P \rangle$  is given. One example for  $\mathbf{P}$  is our usage property domain  $\mathbf{U}$  from Section 5, which we will use as the predicate domain instantiation in Section 7.3. It is possible to use a more complex list element abstraction that relates to the list element index or even the segment index.

## 7.2 Generic List Segmentation Abstract Domain

We now combine the domains introduced to build the list segmentation abstract domain with the signature  $\mathbf{Seg}(\mathbf{L}(\mathbf{E}_V), \mathbf{P}, \mathbf{V})$ . A segmentation  $g$  has the form

$$\underbrace{\{b_1^0 \dots b_{m^0}^0\}}_{L_0} p_0[?]_0 \underbrace{\{b_1^1 \dots b_{m^1}^1\}}_{L_1} p_1[?]_1 \quad \dots \quad p_{n-1}[?]_{n-1} \underbrace{\{b_1^n \dots b_{m^n}^n\}}_{L_n},$$

where  $b_j^i \in E$ ,  $i \in \{0, \dots, n\}$ ,  $j \in \{0, \dots, m^j\}$  and  $p_k \in P$ ,  $[?]_k \in \{?, -\}$ ,  $k \in \{0, \dots, n-1\}$ . We use  $\text{len}(g) = n$  to denote the length of the segmentation. Note that, the number of limits is  $\text{len}(g) + 1$ . We use the notation  $g.p_i$  to retrieve the predicate at index  $i$ . We introduce a formalism for one or multiple updates of predicates at indices  $\mathcal{I}$  to the new predicate  $P$ , conditioned on  $c$  (may depend on the old predicate at index  $i$ ):

$$g \left[ p_i \xleftarrow{c(p_i)} P \mid i \in \mathcal{I} \right] \stackrel{\text{def}}{=} L_0 p'_0[?]_0 L_1 p'_1[?]_1 \quad \dots \quad p'_{n-1}[?]_{n-1}$$

where  $p'_i = \begin{cases} P & \text{if } i \in \mathcal{I} \wedge c(p_i) \\ p_i & \text{otherwise} \end{cases}$

This is similar to the formalism for stores in Section 4.1.

The interpretation of a segment is that the predicate  $p_i$  covers the segment range in the concrete. The upper bound is always exclusive while the lower bound is inclusive except when the segment is empty, i.e., when the limits coincide. The question mark  $[?]_k$  after each predicate is either present (?) or absent (-) and denotes if the segment can possibly be empty or we know that at least one bound in the left limit is strictly smaller than a bound in the right limit:

$$[?]_k = \begin{cases} - & \text{if } \exists b \in L_k, b' \in L_{k+1} : b \sqsubseteq_E b' \\ ? & \text{otherwise} \end{cases}$$

In case two limits actually coincide in the concrete, the segment has zero length and the abstract predicate of that segment has no effect on the concrete. With the order  $- \preceq ?$  and the least



upper bound  $\vee$  and greatest lower bound  $\wedge$  defined as usual we get the lattice  $L_\tau = \langle \{?, -\}, \leq, \vee, \wedge, -, ? \rangle$ .

The consecutive, non-overlapping segments cover the whole list at any time, which means that the extremal limits are supposed to be at the limits of the concrete list. They are never removed and the invariant holds, that  $0 \in L_0$  and  $list\_len \in L_n$ . This implies that there is always at least one limit in the segmentation (since the two extremal limits can be equal for an empty list).

### 7.2.1 Segmentation Unification

To be able to define a partial order and join, meet, widening, narrowing operators on segmentations we need a way to *unify* two segmentations  $g, g'$ . After unifying, they must have the same number and equal limits at each index, i.e.,  $\forall i \in \{0, \dots, n\} : L_i = L'_i$ , where at this point  $n = \text{len}(L_i) = \text{len}(L'_i)$ . The operators can then be defined by lifting the operators of the predicate abstract domain segment-wise.

We adapt the unification algorithm from [CCL11] in Algorithm 1 to use relational constraints from the variable environment abstract domain. The algorithm, denoted by  $\text{UNIFY}_{(\mathbb{I}_l, \mathbb{I}_r)}$ , is parameterized by the left and right neutral element  $\mathbb{I}_l$  and  $\mathbb{I}_r$ . The parameterization is necessary because we use the neutral elements as placeholders during unification for newly inserted segments that get later compared/joined/etc.

We give an overview of the algorithm. Note that,  $\text{UNIFY}_{(\mathbb{I}_l, \mathbb{I}_r)}(g, g')$  is not symmetrical with respect to its segmentation parameters  $g, g'$ , since the respective predicates are preserved (as good as possible) while unifying the *limits*. The algorithm is presented as a recursive function, which can be thought as handling *tails* of the two input segmentations from a specific starting index onwards<sup>14</sup>. Let us call the first limit of a tail the *front limit*. The algorithm maintains the invariant<sup>15</sup>, that the left parts up to the front limits of both segmentations are already unified.

The algorithm inspects, in each recursion, the front limits of the tails. The outermost case distinction separates cases whether the front limits are syntactically comparable, i.e., if the bounds of one front limit are contained in the other front limit. If there are comparable, splitting the limits into parts allows ending up with syntactical equivalent front limits and move on. For the bounds not agreeing in the front limits, the inner case distinction of the algorithm decides about keeping them for later unification or drop them definitely.

If the front limits have no common bound, we use the relational comparison based on the variable environment abstract domain to know if we can still order the limits and keep the larger one for later unification. If this check does not help either, we have to drop both front limits.

<sup>14</sup>This indices do not appear explicitly in the algorithm but implicitly given by the first limit of the tails,  $L_i$  and  $L_j$ .

<sup>15</sup>The invariant is stated as a precondition in the algorithm because of the recursive definition.

---

**Algorithm 1** Recursive Unification of  $g = L_0 p_0[?]_0 L_1 \dots$  with length  $n$  and  $g' = L'_0 p'_0[?]'_0 L'_1 \dots$  with length  $n'$  using relational constraints from the variable environment abstract domain.

---

```

UNIFY( $\mathbb{I}_l, \mathbb{I}_r$ )( $L_i p_i[?]_i L_{i+1} \dots$ ,  $L'_j p'_j[?]'_j L'_{j+1} \dots$ )
  Precondition: The left parts  $\dots p_{i-2}[?]_{i-2} L_{i-1}$  and  $\dots p'_{i-2}[?]_{i-2} L'_{i-1}$  are already unified

  // Recursion ending criteria
  if  $i = n + 1 \wedge j = n' + 1$  then
    return

  if  $L_i = L'_j$  then
    UNIFY( $L_{i+1} \dots$ ,  $L'_{j+1} \dots$ )
  else if  $L_i \supset L'_j$  then
    Let  $B$  be set of bounds such that  $L_i \supset B$  and  $L_i - B = L'_j$ .
    Let  $\bar{B} \subseteq B$  be set of bounds appearing later in limits of second segmentation  $L'_{j+1} \dots$ 
    if  $\bar{B} = \emptyset$  then
      UNIFY( $L'_j p_i[?]_i L_{i+1} \dots$ ,  $L'_j p'_j[?]'_j L'_{j+1} \dots$ )
    else
      UNIFY( $L'_j \mathbb{I}_l ? B p_i[?]_i L_{i+1} \dots$ ,  $L'_j p'_j[?]'_j L'_{j+1} \dots$ )
  else if  $L_i \subset L'_j$  then
    // Symmetrical case not shown for brevity
  else if  $L_i \cap L'_j \neq \emptyset$  then
    Let  $A = L_i \cap L'_j$ ,  $B = L_i - A$ ,  $B' = L'_j - A$ .
    Let  $\bar{B} \subseteq B$  be set of bounds appearing later in limits of second segmentation  $L'_{j+1} \dots$ 
    Let  $\bar{B}' \subseteq B'$  be set of bounds appearing later in limits of first segmentation  $L_{i+1} \dots$ 
    if  $\bar{B} = \emptyset = \bar{B}'$  then
      UNIFY( $A p_i[?]_i L_{i+1} \dots$ ,  $A p'_j[?]'_j L'_{j+1} \dots$ )
    else if  $\bar{B} = \emptyset$  then
      UNIFY( $A p_i[?]_i L_{i+1} \dots$ ,  $A \mathbb{I}_r ? \bar{B}' p'_j[?]'_j L'_{j+1} \dots$ )
    else if  $\emptyset = \bar{B}'$  then // Symmetrical case
      UNIFY( $A \mathbb{I}_l ? \bar{B} p_i[?]_i L_{i+1} \dots$ ,  $A p'_j[?]'_j L'_{j+1} \dots$ )
    else //  $\bar{B} \neq \emptyset \wedge \bar{B}' \neq \emptyset$ 
      UNIFY( $A \mathbb{I}_l ? \bar{B} p_i[?]_i L_{i+1} \dots$ ,  $A \mathbb{I}_r ? \bar{B}' p'_j[?]'_j L'_{j+1} \dots$ )
  else //  $L_i \cap L'_j = \emptyset$ , which implies we are not at start, i.e  $i \neq 0 \wedge j \neq 0$ 
    if  $L_i \hat{=} L'_j$  then
      UNIFY( $(L_i \cup L'_j) p_i[?]_i L_{i+1} \dots$ ,  $(L_i \cup L'_j) p'_j[?]'_j L'_{j+1} \dots$ )
    else if  $L_i \sqsubseteq L'_j$  then
      UNIFY( $L_i p_i[?]_i L_{i+1} \dots$ ,  $L_i \mathbb{I}_r L'_j p'_j[?]'_j L'_{j+1} \dots$ )
    else if  $L_i \supseteq L'_j$  then // Symmetrical case
      UNIFY( $L'_j \mathbb{I}_r L_i p_i[?]_i L_{i+1} \dots$ ,  $L'_j p'_j[?]'_j L'_{j+1} \dots$ )
    else
      // Remove limit  $L_i, L'_j$ , merge consecutive segments
      // If either  $L_i$  or  $L'_j$  is the segmentations upper limit, do not remove it
      // (This case distinction is not shown for brevity)
      UNIFY( $L_{i-1} (p_{i-1} \sqcup p_i) ([?]_{i-1} \vee [?]_i) L_{i+1} \dots$ ,  $L'_{j-1} (p'_{j-1} \sqcup p'_j) ([?]_{j-1} \vee [?]_j) L'_{j+1} \dots$ )

```

---

**Lattice Operators** As mentioned, the operators are applied segment-wise after unification. For the segmentation join and widening, we use the parameterization  $\text{UNIFY}_{(\perp_P, \perp_P)}$ , for segmentation meet and narrowing  $\text{UNIFY}_{(\top_P, \top_P)}$  and for partial order  $\text{UNIFY}_{(\perp_P, \top_P)}$ :

$$\begin{aligned}
g \sqsubseteq_{\text{Seg}} g' &\stackrel{\text{def}}{=} \forall i \in \{0, \dots, \text{len}(\check{g}) - 1\} : \check{g}.p_i \sqsubseteq_P \check{g}'.p_i && \text{where } \check{g} = \text{UNIFY}_{(\perp_P, \top_P)}(g, g') \\
&&& \check{g}' = \text{UNIFY}_{(\perp_P, \top_P)}(g', g) \\
g \sqcup_{\text{Seg}} g' &\stackrel{\text{def}}{=} \check{g} [p_i \leftarrow \check{g}.p_i \sqcup_P \check{g}'.p_i \mid i \in \{0, \dots, \text{len}(\check{g}) - 1\}] && \text{where } \check{g} = \text{UNIFY}_{(\perp_P, \perp_P)}(g, g') \\
&&& \check{g}' = \text{UNIFY}_{(\perp_P, \perp_P)}(g', g) \\
g \sqcap_{\text{Seg}} g' &\stackrel{\text{def}}{=} \check{g} [p_i \leftarrow \check{g}.p_i \sqcap_P \check{g}'.p_i \mid i \in \{0, \dots, \text{len}(\check{g}) - 1\}] && \text{where } \check{g} = \text{UNIFY}_{(\top_P, \top_P)}(g, g') \\
&&& \check{g}' = \text{UNIFY}_{(\top_P, \top_P)}(g', g)
\end{aligned}$$

We add a bottom element  $\perp_{\text{Seg}}$  representing an unreachable segmentation, i.e., when any of the bounds or the predicates<sup>16</sup> is unreachable. We do not need a special top element since we can set

$$\top_{\text{Seg}} \stackrel{\text{def}}{=} \{0\} \top? \{list\_len\}.$$

Finally we get lattice  $L_{\text{Seg}} = \langle \text{Seg}, \sqsubseteq_{\text{Seg}}, \sqcup_{\text{Seg}}, \sqcap_{\text{Seg}}, \perp_{\text{Seg}}, \top_{\text{Seg}} \rangle$ .

**Widening** The widening of segmentations is twofold. We can widen the predicates segment-wise like the other operators. This is not sufficient because it can happen that we infinitely add new limits to a segmentation. We propose a widening operator  $\nabla_{\text{Seg}}$ , that when applied to two segmentations  $g \nabla_{\text{Seg}} g'$ , drops all limits from  $g'$  that can not be found in  $g$  based on a syntactic comparison of equality.

## 7.2.2 Segmentation Operators

Let  $\text{CLEAN}(g)$  be an operation on segmentations that merges all pairs of consecutive limits in  $g$  if they can be determined to be equal (using relational equality  $\hat{=}_L$ ). The predicate of the segment between the merged limits is discarded. Secondly,  $\text{CLEAN}(g)$  removes any question mark if a segment cannot be empty, i.e.,  $\forall i \in \{0, \dots, n\}$ , it sets  $[?]_i = \_$  if  $L_i \sqsubset_L L_{i+1}$ <sup>17</sup>. The  $\text{CLEAN}(g)$  operator is introduced here for the formal presentation, in an implementation it is better performance-wise to clean up adjacent limits whenever a limit is added or changed during an operation than to do a full clean at the end of an operation.

We need a way to update a segmentation not only at specific predicate indices, but also in a range given by expressions. To find the worst-case range of affected predicates in the abstract, we introduce the *least upper limit* and *greatest lower limit* for an index given as an expression  $e$  in normal form. Let again  $g = L_i p_i [?]_i L_{i+1} \dots$  be a segmentation with length  $n$ .

$$\begin{aligned}
\text{luli}[g](e) &\stackrel{\text{def}}{=} \underset{i}{\text{argmin}} (\{i \in \{0, \dots, n\} \mid e < L_i\} \cup \{n\}) && \text{(index of least upper limit)} \\
\text{glli}[g](e) &\stackrel{\text{def}}{=} \underset{i}{\text{argmax}} (\{i \in \{0, \dots, n\} \mid L_i \leq e\} \cup \{0\}) && \text{(index of greatest lower limit)}
\end{aligned}$$

We assume that the query expression  $e$  is always between the extremal limits, i.e.,  $L_0 \leq e < L_n$  (no under-/overflow of list indices). This assumption manifests in the added indices  $n$  and  $0$ ,

<sup>16</sup> Note that, not all predicate abstract domains include an element for representing unreachability, and even though often used as,  $\perp$  does not represent an unreachable state in general.

<sup>17</sup> We did not define the  $\sqsubset_L$ , but for some combinations of expression forms and variable environments it can be defined.

which guarantee that `argmin`, `argmax` and in turn all of the following operators have a well-defined result.

**Set Predicate** When we want to set a predicate in the segmentation, we may need to adjust the limits to enable the most precise reflection of that predicate change. We do not only want to support updates at a constant index  $c$ , but also within an interval  $[l, u]$  or at an expression  $e$  in normal form. To suit all these requirements, we define a more general function that takes two expressions  $e_l, e_u$  to bound the range where we want to set a predicate  $p$ :

$$\text{SET}_{\text{Seg}}[g](e_l, e_u; p) \stackrel{\text{def}}{=} \text{CLEAN} \left( \cdots L_{\text{glli}[g](e_l)} p_{\text{glli}[g](e_l)} \{e_l\} p \{e_b + 1\} p_{\text{luli}[g](e_u)-1} ? L_{\text{luli}[g](e_u)} \cdots \right)$$

This operator creates a new segment with two new limits, the lower limit  $\{e_l\}$  (inclusive) and the upper limit  $\{e_b + 1\}$  (exclusive), and sets the new predicate  $p$  set between those. This segment is inserted into the segmentation between the tightest bounds such that the new segment's limits conform to the order of existent limits. Uncomparable limits between those tightest bounds are removed in favor of the new segment. By the final application of the `CLEAN` operator the newly inserted limits might be merged with existent limits.

It is easy to overload this function to support the different ways of defining the range mentioned above, like  $\text{SET}_{\text{Seg}}[g](l, u; p)$ ,  $\text{SET}_{\text{Seg}}[g](c; p)$  and  $\text{SET}_{\text{Seg}}[g](e; p)$ .

**Get Predicate** To query the segmentation for a predicate at a specified index, given in one of the formats listed in last paragraph, we again offer a general function to get the predicate within a query range given by  $e_l, e_u$ :

$$\text{GET}_{\text{Seg}}[g](e_l, e_u) \stackrel{\text{def}}{=} \bigsqcup_P \{p_{\text{glli}[g](e_l)}, \dots, p_{\text{luli}[g](e_u)-1}\}$$

This operator finds the tightest limits bounding the query range and returns the least upper bound of the predicate between those. We can again overload this function to support the different ways of defining the range, like  $\text{GET}_{\text{Seg}}[g](l, u)$ ,  $\text{GET}_{\text{Seg}}[g](c)$  and  $\text{GET}_{\text{Seg}}[g](e)$ .

### 7.2.3 Abstract Transformers for Backward Analysis

In [CCL11] the authors present the abstract transformers for a forward analysis. For a backward analysis we need to replace the forward transformer for *variable assignment* by a backward transformer. The steps to be performed when transforming the state upon an assignment can be divided into two parts:

1. Replace all symbolic variables used in the bounds of the segmentation.
2. Invoke the transformer of the abstract predicate domain, e.g., to reflect the flow of information (depending on the goal of the analysis). This usually involves transforming some of the predicates in the segmentation, but can also require changing the segmentation itself, e.g., adding or removing limits.

**Assignment to a scalar variable (variable substitution)** In an assignment  $x = e$ , where  $x$  is a scalar variable and  $e$  is any expression, we replace all occurrences of  $x$  in the segmentation bounds by the normal form of  $e$ . If  $e$  cannot be put in normal form we have to use a fallback. The idea is to replace each bound  $b = (x + c)$  containing  $x$  by a segment with predicate  $\top_P$  covering the worst-case range  $[b_l, b_u]$  that bound may take in the concrete. To achieve this, we use the variable environment abstract domain from the program point *below* the assignment to

get an interval abstraction  $[e_l, e_u]$  of  $e$ , calculate  $[b_l, b_u] = [e_l, e_u] + c$  and set the segmentation to  $\top_P$  in that range  $[b_l, b_u]$ .

$$\text{REPLACE}_{\text{Seg}}[x = e]g \stackrel{\text{def}}{=} \bigsqcup_{\text{Seg}} \left\{ \text{SET}_{\text{Seg}}[g](b_l, b_u, \top_P) \mid b \in \text{bounds of } g \text{ containing } x \right\}$$

**Example 3.** Assume we have integers variables  $x, y, a$  in a program and the variable environment provides the constraints  $0 \leq a \leq 1$  and  $0 \leq y$ . Assume we want to transform

$$g = [\{0\}p_0? \underbrace{\{x + 1\}}_{L_1} p_1? \{y\} p_2 \{y + 1\} p_3 \{y + 2\} p_4 \{y + 3\} p_5 \{y + 4\} p_6? \{\text{len}(g)\}]$$

upon the assignment  $[x = 2 * a]$  with a right side that cannot be put in normal form. We evaluate  $[e_l, e_u] = [0, 2]$  and thus we get the worst-case range  $[b_l, b_u] = [1, 3]$  for limit  $L_1$ . Extending  $L_1$  to the worst-case range and setting the predicate to top inside that range results in:

$$\text{REPLACE}_{\text{Seg}}[x = 2 * a]g = [\{0\}p_0\{1\}\top_P\{4\}p_5?\{y + 4\}p_6?\{\text{len}(g)\}]$$

Implementation-wise, instead of joining all the individual segmentations reflecting a change in one bound each, we can use an iterative application of  $\text{SET}_{\text{Seg}}$ . Note that, if  $e$  contains any  $\text{input}()$  expressions, this leads to an interval abstraction of the right side of  $[b_l, b_u] = [-\infty, \infty]$ . In that case, all inner limits have to be dropped. However, the way  $\text{SET}_{\text{Seg}}$  is defined, it ensures the extremal limits are never removed.

**Assignment to a list** In a list assignment  $l = l'$ , where  $l$  is a list variable and  $l'$  is either another list variable or a list display expression, we do not have to replace anything inside a segmentation since segmentations do not contain list variables. So we only perform step 2.

In an assignment  $l[t] = e$  to a list index, where  $l$  is a list variable and  $t$  is a scalar expression and  $e$  any expression, we again only have to perform step 2.

### 7.3 List Segmentation for Usage Property

We instantiate the generic list segmentation abstract domain for the usage analysis to get  $\mathbf{USeg} = \mathbf{Seg}(\mathbf{L}(\mathbf{E}_{\text{Oct}}), \mathbf{U}, \mathbf{Oct})$ . We use the octagon domain as the variable environment abstract domain and the expressions of the form as in Section 7.1.2 using octagonal constraints for the rich comparison operators.

#### 7.3.1 Integer/List stores

We track the scalar (integer) variables of the program with members of the usage domain  $\mathbf{U}$  as in the integer-only programs from Section 5.2, and we track list variables with segmentations. In  $\mathbf{USeg}$  we use abstract stores  $q \in Q$ ,  $q = q_{\text{int}} \cup q_{\text{list}}$ , where  $q_{\text{int}} : \text{Var}_{\text{int}} \rightarrow \mathcal{U}$  and  $q_{\text{list}} : \text{Var}_{\text{list}} \rightarrow \mathbf{USeg}$ .

**Information Flow Operators** We lift the function  $\text{REPLACE}$  to stores:

$$\text{REPLACE}_Q[x = e]q \stackrel{\text{def}}{=} q[l \leftarrow \text{REPLACE}_{\text{Seg}}[x = e]q(l) \mid l \in \text{Var}_{\text{list}}]$$

For other types of assignment this operator is the identity.

We extend usage stores to be able to filter expressions containing both scalar variable accesses and list index expressions. We define the condition when the operands in a filtered expression are considered used, i.e., when a store is *dirty*, as its own operator:

$$\begin{aligned} \text{DIRTY}(q) &\stackrel{\text{def}}{=} \exists y \in \mathbb{V}\text{ar}_{\text{int}} : q(y) \in \{\downarrow, \dagger\} \\ &\quad \vee \exists l \in \mathbb{V}\text{ar}_{\text{list}} : \exists i \in \{0, \dots, \text{len}(q(l)) - 1\} : q(l).p_i \in \{\downarrow, \dagger\} \end{aligned}$$

**Implicit Flow** We account for implicit flow by checking if a store is dirty and updating the store accordingly:

$$\begin{aligned} \text{FILTER}_Q[[e]]q &\stackrel{\text{def}}{=} q \left[ x \xleftarrow{\text{DIRTY}(q)} \downarrow \mid x \in \text{opd}(e) \cap \mathbb{V}\text{ar}_{\text{int}} \right] \\ &\quad \dot{\downarrow} q \left[ l \xleftarrow{\text{DIRTY}(q)} \text{SET}_{\text{Seg}}[q(l)](t; \downarrow) \mid l[t] \in \text{opd}(e) \cap \mathbb{I}\text{ndex} \right] \end{aligned}$$

**Explicit Flow** We account for explicit flow by checking if left side of assignment is currently used:

$$\begin{aligned} \text{USE}_Q[[x = e]]q &\stackrel{\text{def}}{=} q \left[ y \xleftarrow{q(x) \in \{\downarrow, \downarrow\}} \downarrow \mid y \in \text{opd}(e) \cap \mathbb{V}\text{ar}_{\text{int}} \right] \\ &\quad \dot{\downarrow} q \left[ l \xleftarrow{q(x) \in \{\downarrow, \downarrow\}} \text{SET}_{\text{Seg}}[q(l)](t; \downarrow) \mid l[t] \in \text{opd}(e) \cap \mathbb{I}\text{ndex} \right] \end{aligned}$$

The transformer  $\text{USE}_Q[[l[t] = e]]q$  is similar except that the condition when to update the store is  $\text{GET}_{\text{Seg}}[q(l)](t) \in \{\downarrow, \downarrow\}$ . When assigning to a list variable, we propagate the used property naturally to the right side of the assignment, is it another list variable or a list display expression:

$$\begin{aligned} \text{USE}_Q[[l = l']]q &\stackrel{\text{def}}{=} q \left[ l' \leftarrow q(l) \left[ p_i \xleftarrow{l \neq l' \wedge p_i \in \{\downarrow, \downarrow\}} \downarrow \mid i \in \{0, \dots, \text{len}(l)\} \right] \right] \\ \text{USE}_Q[[l = [t_0, \dots, t_{m-1}]]]q &\stackrel{\text{def}}{=} q \left[ x \xleftarrow{\text{GET}_{\text{Seg}}[q(l)](k) \in \{\downarrow, \downarrow\}} \downarrow \mid x \in \text{opd}(t_k) \cap \mathbb{V}\text{ar}_{\text{int}}, k \in \{0, \dots, m-1\} \right] \end{aligned}$$

**Intermitted Flow** The intermitted flow for an assignment to a scalar variable is as before:

$$\text{KILL}_Q[[x = e]]q \stackrel{\text{def}}{=} q \left[ x \xleftarrow{q(x) \in \{\downarrow, \downarrow\} \wedge x \notin \text{opd}(e)} \dagger \right]$$

We add the definitions for the assignment to list variables:

$$\begin{aligned} \text{KILL}_Q[[l = l']]q &\stackrel{\text{def}}{=} q \left[ l \leftarrow q(l) \left[ p_i \xleftarrow{l \neq l' \wedge p_i \in \{\downarrow, \downarrow\}} \dagger \mid i \in \{0, \dots, \text{len}(l)\} \right] \right] \\ \text{KILL}_Q[[l = [t_0, \dots, t_{m-1}]]]q &\stackrel{\text{def}}{=} q \left[ l \leftarrow q(l) \left[ p_i \xleftarrow{p_i \in \{\downarrow, \downarrow\}} \dagger \mid i \in \{0, \dots, \text{len}(l)\} \right] \right] \end{aligned}$$

To determine the intermitted flow for an assignment to a list index  $l[t] = e$ , we have to do a little more work. We have to find the subset of indices of  $t$  that do not appear in the expression

$e$  on the right side:

$$\text{KILL}_Q \llbracket l[t] = e \rrbracket q \stackrel{\text{def}}{=} q[l \leftarrow$$

$$q(l) \left[ \begin{array}{l} p_i \xleftarrow{p_i \in \{\downarrow, \downarrow\}} \text{ } \\ \text{ } \end{array} \right] \left[ \begin{array}{l} i \in \underbrace{\{\text{glli}[q(l)](t), \dots, \text{luli}[q(l)](t) - 1\}}_{\text{predicate indices overwritten left}} \\ - \underbrace{\bigcup \left\{ \{\text{glli}[q(l)](r), \dots, \text{luli}[q(l)](r) - 1\} \mid l[r] \in \text{opd}(e) \cap \text{Var}_{\text{list}} \right\}}_{\text{predicate indices used right}} \end{array} \right]$$

$$\text{ ]}$$

Finally we can combine this operators to adapt the partial abstract semantics given in (3) for int/list stores and the backwards substitution required by segmentations:

$$\llbracket \psi = \omega \rrbracket q \stackrel{\text{def}}{=} (\text{REPLACE}_Q \llbracket \psi = \omega \rrbracket \circ \text{KILL}_Q \llbracket \psi = \omega \rrbracket \circ \text{USE}_Q \llbracket \psi = \omega \rrbracket) q$$

$$\llbracket \text{print}(e) \rrbracket q \stackrel{\text{def}}{=} q[x \leftarrow \downarrow \mid x \in \text{opd}(e) \cap \text{Var}_{\text{int}}]$$

$$\quad \sqcup q[l \leftarrow \text{SET}_{\text{Seg}}[q(l)](t; \downarrow) \mid l[t] \in \text{opd}(e) \cap \text{Index}]$$

$$\quad \sqcup q[l \leftarrow \text{T}_{\text{Seg}} \mid l \in \text{opd}(e) \cap \text{Var}_{\text{list}}]$$

### 7.3.2 Abstract Transformers

A notable difference to the transformers for the analysis for int-only programs is that upon an assignment, we have to transform all stack frames. In the top frame we have to call the complete transformer  $\llbracket \psi = \omega \rrbracket$  as before (which includes the variable replacement and update of usage). In all other frames except the the top frame we have to only replace the assigned variable with  $\text{REPLACE}_Q$ .

$$\llbracket \text{pass} \rrbracket s \stackrel{\text{def}}{=} s$$

$$\llbracket \psi = \omega \rrbracket s \stackrel{\text{def}}{=} \langle \text{REPLACE}_Q \llbracket \psi = \omega \rrbracket s_1, \dots, \text{REPLACE}_Q \llbracket \psi = \omega \rrbracket s_{k-1}, \llbracket \psi = \omega \rrbracket s_k \rangle$$

$$\llbracket \text{print}(e) \rrbracket s \stackrel{\text{def}}{=} \langle s_1, \dots, s_{k-1}, \llbracket \text{print}(e) \rrbracket s_k \rangle$$

$$\llbracket \text{if } b : \text{stmt}_1 \text{ else } \text{stmt}_2 \rrbracket s \stackrel{\text{def}}{=} \text{POP}((\text{FILTER}_S \llbracket b \rrbracket \circ \llbracket \text{stmt}_1 \rrbracket)(\text{PUSH}(s)))$$

$$\quad \sqcup \text{POP}((\text{FILTER}_S \llbracket \text{not } b \rrbracket \circ \llbracket \text{stmt}_2 \rrbracket)(\text{PUSH}(s)))$$

$$\llbracket \text{while } b : \text{stmt} \rrbracket s \stackrel{\text{def}}{=} \text{lfp}_s [\lambda t. \text{FILTER}_S \llbracket \text{not } b \rrbracket t \sqcup \text{POP}((\text{FILTER}_S \llbracket b \rrbracket \circ \llbracket \text{stmt} \rrbracket)(\text{PUSH}(t)))]$$

### 7.3.3 Concretization

The generic concretization given in [CCL11] is fully applicable to our instantiation, since we did not change parts of the segmentation relevant for the concretization.

### 7.3.4 Example

We inspect the two bugs in Code 1 from the introduction separately. Code 4a shows a variant with only Bug A in place. Since Bug A makes the whole while-loop effectless on the output, the analysis correctly detects on line 4:

$$count \rightarrow \perp, i \rightarrow \perp, g \rightarrow \perp, list\_grades \rightarrow [\{0\} \perp \{len(list\_grades)\}]$$

Code 4b shows a variant with only Bug B in place. The while-loop now has an effect, but the first item in the list is excluded from the counting. Thus, the result on line 4 indicates that the list is only used from index 1 onwards:

$$count \rightarrow \perp, i \rightarrow \perp, g \rightarrow \perp, \\ list\_grades \rightarrow [\{0\} \perp \{i\} \perp \{i+1\} \perp \{i+1\} \perp \{i+1\} \perp \{len(list\_grades)\}]$$

```

1 list_grades = listinput()
   count = 0
5 i = 0
  while i < len(list_grades):
    g = list_grades[i]
    if g >= 4:
      g += 1 # Bug A
10
   i += 1
  print(count)

```

(a) with only Bug A

```

1 list_grades = listinput()
   count = 0
5 i = 1 # Bug B
  while i < len(list_grades):
    g = list_grades[i]
    if g >= 4:
      count += 1
10
   i += 1
  print(count)

```

(b) with only Bug B

Code 4: Programs calculating the number of passing grades, each with a different bug.



## 8 Implementation

All domains described in this report have been implemented and tested in the abstract interpretation framework *Lyra*<sup>18</sup>, which was built alongside this project as a generic framework to explore abstract interpretation analysis.

### 8.1 Framework for Abstract Interpretation

The framework is written in Python v3.5. It follows a generic approach to support different input languages by translating the input programs into an intermediate control flow graph (CFG) representation. The CFG contains links to the program points the statements derive their origin from. The interpreter<sup>19</sup> transforms the abstract states while traversing the CFG and collects the analysis results. How the states are transformed can be specified by customizable semantics. The states itself are built of a hierarchy of abstract domains or partial domains. A state offers only a narrow interface to the interpreter and does not have any knowledge about the CFG structure but knows how to transform itself upon assignments, conditions etc.

**Lattice and BoundedLattice** At the core of each domain is at least one lattice structure. In case a domain is used directly by an interpreter, it must built on a lattice that has a dedicated bottom element. This is used by the interpreter to represent unreachability.

**State** A state is itself a lattice, combining one or several lattices to an abstract domain, handling the interaction between them and offering all necessary interface methods for the interpreter, i.e., the abstract transformers.

**Generic Domains** The abstract domain functors from Section 4 have their direct equivalent in the framework, called *Generic Domains*. They are designed as superclasses or mixins to generalize frequent structures of lattices like the variable store.

### 8.2 Automated Testing

While not every analysis supports all constructs of a modern programming language, it is still helpful to limit the input programs to a real subset of, e.g., Python, to make the sample programs actually runnable.

Apart from usual unit tests, we added a full-chain testing framework. This operates on files, each representing a test, and does all the steps necessary to test the analysis on the code in the file:

1. Parse the file into the abstract syntax tree (AST),
2. translate the AST to the CFG,
3. run one or multiple connected analysis on it,
4. check the actual analysis results against expected results provided as special comments in the input files,

---

<sup>18</sup> <https://github.com/caterinaurban/Lyra>

<sup>19</sup> The interpreter of the analysis must not be confused with the Python interpreter.

5. plot (intermediate) results for manual inspection.

### 8.3 Examples

Code 5 shows an erroneous example program. The programmers intention was to halve a given integer  $x$   $n$ -times via integer (floor) division. The bug is marked on line 6. The analysis correctly detects on line 4 that  $i$  and  $n$  (and therefore the complete body of the loop) are effectless. Figure 2 shows the CFG for this program with all intermediate results. If we correct the bug, the analysis result on line 4 becomes  $n \rightarrow \downarrow$  (and  $i \rightarrow \sharp$ ) as expected.

```

1  n = int(input())
   x = int(input())
   i = 0
   # RESULT: i→⊥, n→⊥, x→⧫
5  while i < n:
       x = i / 2 # Bug: should be x = x / 2
       i = i + 1
   else:
       x = -1
10 print(x)

```

Code 5: Program demonstrating expected result specification.

A second example is shown in Code 6. For input grades from 1 (worst) to 6 (best), the program calculates the points a student gained with core and minor subjects throughout a semester. Points are awarded for every grade above 4 (line 15 and 30), and deducted for grades below 4. Grades of core subjects below 4 are subtraced twice (line 17).<sup>20</sup> There are 3 bugs in the code. Bug A causes *total\_point* not be updated with the points accumulated with core subjects. Note that all variables are *truly-live* according to the definition in the introduction. Our analysis is able to detect that the variable *points* is unused at line 15 and line 17, thus the variable *diff* and *g* are also unused in the first while loop. Finally the analysis correctly detects that none of the values in the list of core grades are used at line 3:

$$list\_core\_grades \rightarrow [\{0\} \perp \{len(list\_core\_grades)\}]$$

Bug B and bug C causes the first and the last element of the list of minor grades to be unused. While our algorithm is able to detect Bug B as demonstrated in the example in Section 7.3.4, Bug C is not found because the widening operator is applied in the loop header node after some iterations. This results in the segmentation abstraction for the list of minor grades at line 3:

$$list\_minor\_grades \rightarrow [\{0\} \perp \{i\} \downarrow \{i + 1\} \downarrow \{i + 2\} \downarrow \{i + 3\} \downarrow \{len(list\_minor\_grades)\}]$$

This result could be improved by a more sophisticated widening strategy. There is research that tries to improve widening generically under reasonable performance losses [GR06].

<sup>20</sup>The double deducting of grades below 4 is not made up, this was the rule at the authors secondary school.

```

1 list_core_grades = list(map(int, input().split()))
  list_minor_grades = list(map(int, input().split()))

  # Total accumulated points
5 points = 0
  total_points = 0

  # CALCULATE CORE SUBJECTS POINTS
  i = 0
10 while i < len(list_core_grades):
    g = list_core_grades[i]

    diff = g - 4
    if diff >= 0: # passing grade
15     points += diff
    else: # diff < 0
        points += 2 * diff

    i += 1
20
  # Bug A: missing 'total_points += points'
  points = 0

  # CALCULATE MINOR SUBJECTS POINTS
25 i = 1 # Bug B: should be 'i = 0'
  while i < len(list_core_grades) - 1: # Bug C: should be <=
    g = list_minor_grades[i]

    diff = g - 4
30     points += diff

    i += 1

  total_points += points
35
  if total_points >= 0:
    result = 1 # passed
  else:
    result = -1 # failed
40
  print(result)

```

Code 6: Program collecting points of core and minor subjects and determining if a student passed the semester.

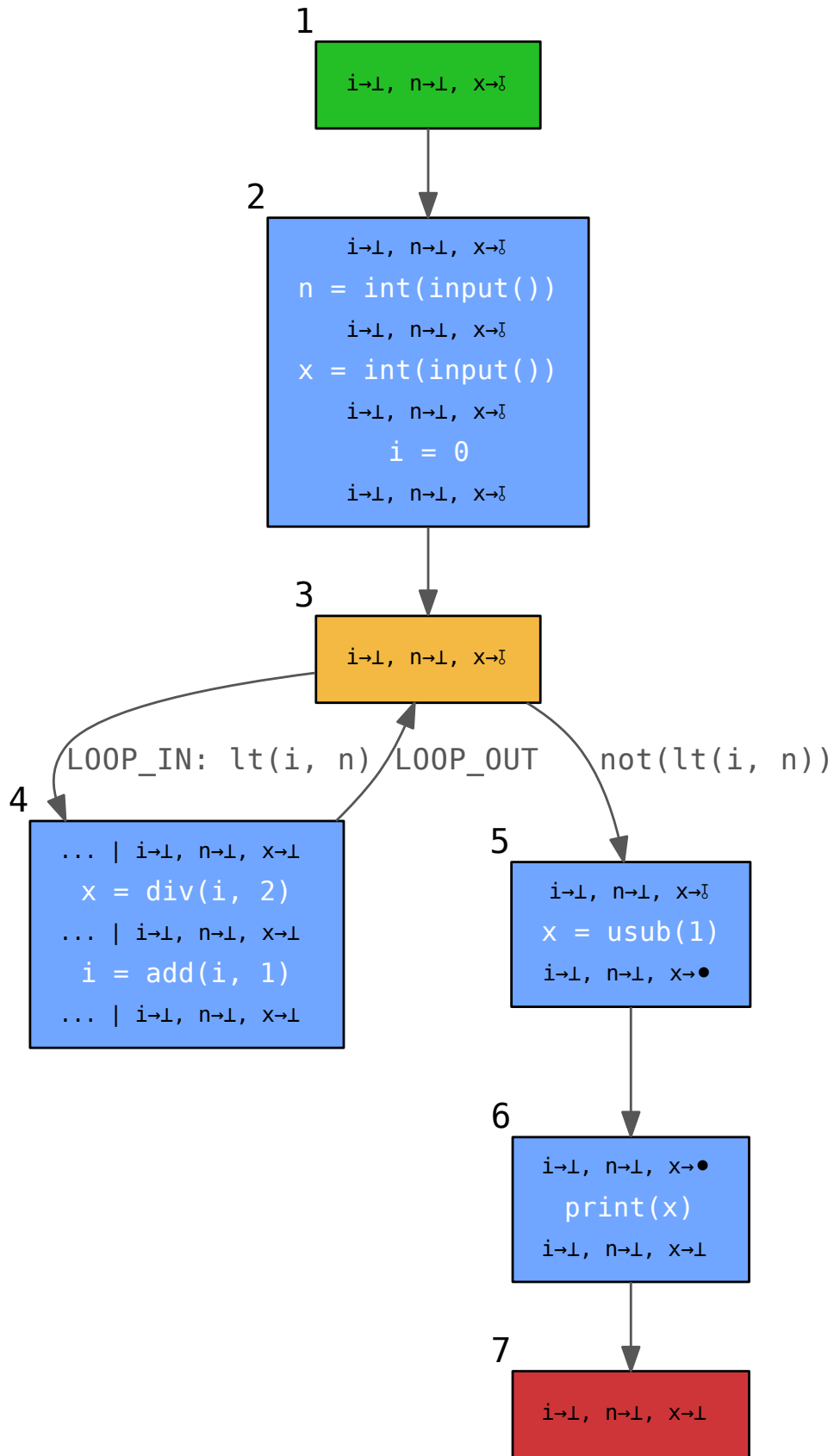


Figure 2: The CFG with usage analysis results for Code 5.

## 9 Conclusion

In this thesis, we started off from well studied problems involving properties of variables of a program, such as the property of *live* and *truly-live* variables and corresponding analysis. We motivated the need for a stronger property, which describes if an input or a variable is used directly, transitively or implicitly by an output statement. We developed a notion and formal definition for the new usage property.

In the framework of Abstract Interpretation, we constructed a simple but powerful abstract domain that approximates the usage property for scalar variables with a backwards analysis. We showed that the concept can be extended to lists, using two approaches: A simple lifting of the usage property to the start of a list already allows detecting bugs involving the content of lists, but the precision is not satisfying. As a next step, we built on known concepts to track list contents by segmentations and instantiated these for the usage property. In this process, we optimized segmentations for our property and improved some of the required operators like the unification of two segmentations.

The prototype developed alongside this thesis implements all the theoretically derived analysis for a simplified Python language and shows their potential in real-world applications. Because the analysis finds a property at the level of variables, we could easily extend it to local variables in functions or methods and parameters of those. The analysis could then be used as an ad-hoc analyzer in programming IDEs.

### Future Work

Possible improvements are twofold: First, we could refine the usage property itself to preserve more detail about how an input or variable is used for the analysis report, e.g., if it is used explicitly or implicitly (in a condition). We could even go as far as finding *how often* or *how much* an input is used and detect irregularities when inputs are expected to contribute equally to the result. Additionally, we could improve the precision of the presented analysis. In the list content usage analysis, the widening operator can be improved to drop limits of a segmentation more carefully. This would allow finding bugs where we miss items at the end of a list.

Secondly, the analysis support can be extended to other language constructs, such as sets, nested lists, procedures and classes. The latter two would ask for an interprocedural analysis. The prototype implementation could be made more user-friendly by reporting back analysis results only for inputs. Because sometimes a program intentionally skips inputs unused, we could either introduce a way for the user to specify which inputs should be checked for the usage property, or we add a pre-analysis that tries to detect this by heuristics.

## References

- [BDR04] Gilles Barthe, Pedro R D’Argenio, and Tamara Rezk. Secure information flow by self-composition. In *Computer Security Foundations Workshop, 2004. Proceedings. 17th IEEE*, pages 100–114. IEEE, 2004.
- [BHZ08] Roberto Bagnara, Patricia Hill, and Enea Zaffanella. An improved tight closure algorithm for integer octagonal constraints. In *Verification, Model Checking, and Abstract Interpretation*, pages 8–21. Springer, 2008.
- [CC77] Patrick Cousot and Radhia Cousot. Abstract interpretation: a unified lattice model for static analysis of programs by construction or approximation of fixpoints. In *Proceedings of the 4th ACM SIGACT-SIGPLAN symposium on Principles of programming languages*, pages 238–252. ACM, 1977.
- [CCL11] Patrick Cousot, Radhia Cousot, and Francesco Logozzo. A parametric segmentation functor for fully automatic and scalable array content analysis. In *ACM SIGPLAN Notices*, volume 46, pages 105–118. ACM, 2011.
- [GM81] Robert Giegerich and Ulrich Möncke. Invariance of approximative semantics with respect to program transformations. In *GI11. Jahrestagung*, pages 1–10. Springer, 1981.
- [GR06] Denis Gopan and Thomas W Reps. Lookahead widening. In *CAV*, volume 4144, pages 452–466. Springer, 2006.
- [Min06] Antoine Miné. The octagon abstract domain. *Higher-order and symbolic computation*, 19(1):31–100, 2006.
- [RHS95] Thomas Reps, Susan Horwitz, and Mooly Sagiv. Precise interprocedural dataflow analysis via graph reachability. In *Proceedings of the 22nd ACM SIGPLAN-SIGACT symposium on Principles of programming languages*, pages 49–61. ACM, 1995.
- [Urb15] Caterina Urban. *Static analysis by abstract interpretation of functional temporal properties of programs*. PhD thesis, École Normale Supérieure, 2015.



## Declaration of originality

The signed declaration of originality is a component of every semester paper, Bachelor's thesis, Master's thesis and any other degree paper undertaken during the course of studies, including the respective electronic versions.

Lecturers may also require a declaration of originality for other written papers compiled for their courses.

I hereby confirm that I am the sole author of the written work here enclosed and that I have compiled it in my own words. Parts excepted are corrections of form and content by the supervisor.

**Title of work** (in block letters):

Static Program Analysis of Data Usage Properties

**Authored by** (in block letters):

*For papers written by groups the names of all authors are required.*

**Name(s):**

Wehrli

**First name(s):**

Simon

With my signature I confirm that

- I have committed none of the forms of plagiarism described in the 'Citation etiquette' information sheet.
- I have documented all methods, data and processes truthfully.
- I have not manipulated any data.
- I have mentioned all persons who were significant facilitators of the work.

I am aware that the work may be screened electronically for plagiarism.

**Place, date**

Zurich, 27.8.2017

**Signature(s)**

S. Wehrli

*For papers written by groups the names of all authors are required. Their signatures collectively guarantee the entire content of the written paper.*