Static Inference of Regular Grammars for Ad Hoc Parsers

ANONYMOUS AUTHOR(S)

 Parsing, the process of structuring a linear representation according to a given grammar, is a fundamental activity in software engineering. While formal language theory has provided theoretical foundations for parsing, the most common kind of parsers used in practice are written ad hoc. They use common string operations for parsing, without explicitly defining an input grammar. These ad hoc parsers are often intertwined with application logic and can result in subtle semantic bugs. Grammars, which are complete formal descriptions of input languages, can enhance program comprehension, facilitate testing and debugging, and provide formal guarantees for parsing code. But writing grammars-e.g., in the form of regular expressions-can be tedious and error-prone. Inspired by the success of type inference in programming languages, we propose a general approach for static inference of regular input string grammars from unannotated ad hoc parser source code. We approach this problem as an intersection of refinement typing and abstract interpretation. We use refinement type inference to synthesize logical and string constraints that represent regular parsing operations, which we then interpret with an abstract semantics into regular expressions. Our contributions include a core calculus λ_{Σ} for representing ad hoc parsers, a formulation of (regular) grammar inference as refinement inference, an abstract interpretation framework for solving refinement variables, and a set of abstract domains for efficiently representing the kinds of numeric and string values encountered during regular ad hoc parsing. We implement our approach in the PANINI system and evaluate its efficacy on a benchmark of 204 Python ad hoc parsers.

1 Introduction

Parsing is one of the fundamental activities in software engineering. It is an activity so common that pretty much every program performs some kind of parsing at one point or another. Yet in every-day software engineering, only a small minority of programs, mainly compilers and some protocol implementations, make use of formal grammars to document their input languages or use formalized parsing techniques such as combinator frameworks [Leijen and Meijer 2001] or parser generators [Johnson and Sethi 1990; Parr and Quong 1995; Warth and Piumarta 2007]. The vast majority of parsing code in software today is *ad hoc*.

Ad hoc parsers are pieces of code that use combinations of common string operations like slice, index, or trim to effectively perform parsing. A programmer manipulating strings in an ad hoc fashion would probably not even think about the fact that they are actually writing a parser. These string-manipulating programs can be found in functions handling command-line arguments, reading configuration files, or as part of any number of minor programming tasks involving strings, often deeply entangled with application logic—a phenomenon known as *shotgun parsing* [Momot et al. 2016]. They have also been shown to produce subtle and difficult to identify semantic bugs [Eghbali and Pradel 2020; Kapugama et al. 2022].

A *grammar* is a complete formal description of all values an input string may assume. It can elucidate the corresponding parsing code, revealing otherwise hidden features and potentially subtle bugs or security issues. By focusing on *data* rather than code, grammars provide a high-level perspective, allowing programmers to grasp an input language directly, without being distracted by the mechanics of the parsing process and the intricacies of imperative string manipulation. Augmenting regular documentation with formal grammars can increase program comprehension by providing alternative representations for a programming task [Fitter and Green 1979; Gilmore and Green 1984]. Because a grammar is also a *generating device*, it is possible to construct any sentence of its language in a finite number of steps—manually or in an automated fashion. Being able to reliably generate concrete examples of possible inputs is invaluable during testing and

Unpublished working draft. Not for distribution.

1:2

Full Source Code String Function Specifications Abstract Domains program slicing SSA/ANF transform Add Hoc Parser Source Code Intermediate Representation

Fig. 1. The complete PANINI system.

debugging. But despite providing all these benefits, hardly anyone ever bothers to write down a grammar, even for more complex ad hoc parsers. Grammars share the same fate as most other forms of specification: they are tedious to write, hard to get right, and seem hardly worth the trouble—especially for such small pieces of code like ad hoc parsers.

The only type of grammar that people actually routinely write down are *regular expressions* [Thompson 1968], probably the biggest and most widely known success story of applied formal language theory. However, in practical use they are often embedded within bigger pieces of ad hoc parsing code and thus usually only describe part of the actual input grammar of an ad hoc parser.

But there is another form of specification that we can draw inspiration from: *types*. Formal grammars are similar to types, in that an ad hoc parser without a grammar is very much like a function without a type signature—it might still work, but you will not have any guarantees about it before actually running the program. Types have one significant advantage over grammars, however: most type systems offer a form of *type inference*, allowing programmers to omit type annotations because they can be automatically recovered from the surrounding context [MacQueen et al. 2020, § 4]. If we could infer grammars like we can infer types, we would reap all the rewards of having a complete specification of our program's input language.

Our Contribution. In this paper, we present a general approach for static inference of regular
 string grammars from unannotated ad hoc parser source code. We pose the problem of (regular)
 grammar inference as a sub-goal of refinement type inference. During type inference, we synthesize
 logical constraints that declaratively represent the parsing operations performed on the input string.
 We then interpret this complex first-order formula using an abstract semantics, in order to find a
 minimal sub-constraint to use as an input string refinement, rendered as a regular expression.
 In brief, the contributions of our work are:

- A core calculus λ_{Σ} for representing ad hoc parsers.
- A formulation of (regular) grammar inference as refinement inference.
- An abstract interpretation framework for solving string refinement variables.
- A set of abstract domains related to regular ad hoc parsing.
- An implementation of our approach in the PANINI system.

2 Overview

Figure 1 presents a schematic overview of PANINI,¹ our end-to-end grammar inference system. At the center of our approach is λ_{Σ} , a language-agnostic intermediate representation for ad hoc parsers. It is powerful enough to represent all relevant parsing operations and simple enough to

¹Named in honor of the Sanskrit grammarian Pāṇini [Bhate and Kak 1991], as well as the delicious Italian sandwiches.

97 98

86

87

88

89

90

91 92

93

94

95





Proc. ACM Program. Lang., Vol. 1, No. OOPSLA, Article 1. Publication date: January 2018025-03-25 18:26. Page 2 of 1-30.

Python	λ_{Σ} specification	
assert b	assert : $\{b : \mathbb{B} \mid b\} \to \mathbb{1}$	assertion
a == b	$eq: (a:\mathbb{Z}) \to (b:\mathbb{Z}) \to \{c:\mathbb{B} \mid c \Leftrightarrow a = b\}$	integer equality
len (s)	length : $(s : \mathbb{S}) \rightarrow \{n : \mathbb{N} \mid n = s \}$	string length
s[i]	charAt: $(s: \mathbb{S}) \rightarrow \{i: \mathbb{N} \mid i < s \} \rightarrow \{c: \mathbb{Ch} \mid c = s[i]\}$	character at ind
s == t	$eqChar: (s: \mathbb{C}h) \to (t: \mathbb{C}h) \to \{b: \mathbb{B} \mid b \Leftrightarrow s = t\}$	character equal

Table 1. Python operations as axiomatic λ_{Σ} specifications.

enable straight-forward refinement type inference. The refinement type system of λ_{Σ} allows us to synthesize constraints over a parser's input string—i.e., it allows us to infer a parser's grammar.

The PANINI system can be separated into a front- and a back-end. In the front-end, ad hoc parsers written in a general-purpose programming language, e.g., Python, are translated into λ_{Σ} programs. In the back-end, those λ_{Σ} programs are statically analyzed and their input string constraints extracted. *This paper is about the back-end*. In short order, we will describe the λ_{Σ} calculus, its refinement type system, our grammar inference algorithm, and the underlying abstract domains. To situate this work, we first want to briefly sketch the front-end process.

2.1 The Front End: From Source to λ_{Σ}

118 After locating an ad hoc parser slice, it is first translated into static single assignment (SSA) 119 form [Braun et al. 2013] and then, via an SSA-to-ANF transformation [Chakravarty et al. 2004], 120 into a PANINI program. This requires a library of string function specifications that map the source 121 language's string operations to equivalent λ_{Σ} functions. Note that it is not necessary to have actual 122 λ_{Σ} implementations of these operations. We only need axiomatic specifications—in the form of type 123 signatures-to capture those properties of the original functions that are necessary to synthesize 124 string grammars. Table 1 shows some examples of such axioms, based on the semantics of certain 125 Python functions. We will use these axioms in the examples throughout this paper.

The front-end transformations (parser slicing, source-to- λ_{Σ} translation) and accompanying axiomatic string function specifications need to be defined and implemented (and proven correct) only once per source programming language. For the remainder of this paper, we will assume a Python-to- λ_{Σ} transformation and a library of Python string function specifications.

131 2.2 The Back End: From λ_{Σ} to Grammar

¹³² λ_Σ is a simple λ-calculus with a refinement type system in the style of *Liquid Types* [Rondon et al. ¹³³ 2008; Vazou et al. 2014]. Refinement types allow us to extend base types with logical constraints. ¹³⁴ This is useful to precisely describe subsets of values, as well as track complex relationships between ¹³⁵ values, all on the type level. For example, the type of natural numbers can be defined as a subset of ¹³⁶ the integers, $\mathbb{N} = \{v : \mathbb{Z} \mid v \ge 0\}$, and we can give a precise definition of the length function on ¹³⁷ strings using a dependent function type and the string length operator |□| of the refinement logic,

length :
$$(s:\mathbb{S}) \to \{n:\mathbb{N} \mid n = |s|\}$$

Checking a refinement type reduces to proving a so-called *verification condition* (VC), a first-order constraint in the refinement logic generated by the type system. The validity of the VC entails the correctness of the program's given and inferred types [Nelson 1980]. For example, in order to check whether $\{x : \mathbb{Z} \mid x > 42\}$ (the type of all numbers greater than forty-two) is a subtype of \mathbb{N} , the VC constraint $\forall x. x > 42 \Rightarrow x \ge 0$ has to be verified. For verification to remain practical, the refinement logic is typically chosen to allow *satisfiability modulo theories* (SMT) [Barrett et al. 2021], which means VCs can be discharged using an off-the-shelf constraint solver such as Z3 [De Moura

138

139

107 108

109

110

111

112

113

114

115

Anon.

1:4

assert s[0] == "a" $\lambda(s:\mathbb{S}).$ let $x = charAt \circ 0$ in let p = eqChar x 'a' in assert p

 $\forall s. \kappa(s) \Rightarrow$ $0 < |s| \land \forall x. x = s[0] \Rightarrow$ $\forall p. (p \Leftrightarrow x = `a`) \Rightarrow$ р

Fig. 2. A simple Python expression (left) and the equivalent λ_{Σ} program (middle) with an incomplete verification condition (right) for its inferred type $\{s : \mathbb{S} \mid \kappa(s)\} \rightarrow \mathbb{1}$.

and Bjørner 2008]. Our system uses quantifier-free linear arithmetic with uninterpreted functions (OF UFLIA) [Barrett et al. 2016] for its refinement predicates, extended with a theory of operations over strings [Berzish et al. 2017].

Refinement Inference. To infer a refinement type for any given term, one must find a predicate that describes (at most) all possible values the term could have during any (successful) run of the program. To facilitate this, refinement type systems typically first infer the basic shapes of all types in the program, using standard type inference à la Hindley-Damas-Milner [Damas and 166 Milner 1982; Hindley 1969], with placeholder variables standing in for as-yet-unknown refinement 167 predicates. These placeholder variables—variously called " κ variables" [Cosman and Jhala 2017], 168 "Horn variables" [Jhala and Vazou 2020], or "liquid type variables" [Rondon et al. 2008]—are also 169 present in the VC at this point and prevent it from being discharged (since they are unknown). The 170 type system then tries to find the strongest satisfying assignments for all refinement variables in 171 the VC constraints, in order to both validate the VC and complete the inferred type. 172

Example 2.1. The simple λ_{Σ} program **if** true **then 1 else** 2 can be inferred to have an incomplete type of shape $\{v : \mathbb{Z} \mid \kappa(v)\}$, where κ is an unknown refinement variable, together with the VC

174 175 176

177

178

179

173

 $(\mathsf{true} \Rightarrow \forall v. v = 1 \Rightarrow \kappa(v)) \land (\mathsf{false} \Rightarrow \forall v. v = 2 \Rightarrow \kappa(v)).$

In order to complete the type, we now have to find an assignment for κ . With such a simple constraint, the refinement solver can easily infer the correct assignment $\kappa(v) \mapsto v = 1$, which validates the VC and produces the final type { $v : \mathbb{Z} \mid v = 1$ }.

180 **Grammar Inference.** If a κ variable stands for an input string refinement, we call this a grammar 181 variable, because its solution must be some finite description of all strings that are accepted by the 182 program, i.e., a grammar. To be practical, we would like this grammar to be as complete as possible. 183

Example 2.2. To illustrate, consider the Python expression in Figure 2 (on the left). Assuming the 184 function specifications from Table 1, we can transform this expression to an equivalent λ_{Σ} program 185 (in the middle) with an inferred top-level refinement type of $\{s : \mathbb{S} \mid \kappa(s)\} \to \mathbb{1}$ and a VC (on the 186 right) that closely matches the program. In order to complete both the VC and the top-level type, we 187 have to find an appropriate assignment for the grammar variable κ . The assignment must take the 188 form of a single-argument function constraining the string s. It is clear that choosing $\kappa(s) \mapsto$ true, 189 i.e., allowing any string for s, is not a valid solution because it does not satisfy the VC. On the 190 other hand, choosing $\kappa(s) \mapsto$ false trivially validates the VC, but it implies that the function could 191 never actually be called, as no string satisfies the predicate false. One possible assignment could be 192 $\kappa(s) \mapsto s =$ "a", which only allows exactly the string "a" as a value for *s*. While this validates the 193 VC and produces a correct type in the sense that it ensures the program will never go wrong, it 194 is much too strict: we are disallowing an infinite number of other strings that would just as well 195

fulfill these criteria (e.g., "aa", "ab", and so on). The correct assignment is $\kappa(s) \mapsto s[0] = a$, which ensures that the first character of the string is "a" but leaves the rest of the string unconstrained. This is equivalent to the regular language $a\Sigma^*$, where Σ is any letter from the input alphabet. Note

The key insight that allows us to find suitable assignments for grammar variables in a general 202 manner is that the top-level VC for a parser will always be of the form $\forall s. \kappa(s) \Rightarrow \varphi$, where s is 203 the input string and φ is a constraint that precisely captures all parsing operations the program 204 performs on s. By simply taking $\kappa(s) \mapsto \varphi$, the VC becomes a trivially valid tautology and the 205 string refinement captures exactly those inputs that the parser accepts. But clearly this solution 206 is practically useless: we want a succinct predicate in a grammar-like form, but the top-level VC 207 consequent φ is a complex term in first-order logic that is basically identical to the program code 208 itself; it does not lead to any further insight about the parser's actual input language. 209

that the solution is a minimized version of the top-level consequent in the VC.

However, we can use φ as a starting point and reduce it into a simpler predicate by performing an *abstract interpretation* of φ . Our approach utilizes an abstract semantics of first-order formulas over string constraints in order to soundly eliminate quantifiers and approximate the parser's input language. The precision of this approximation depends on the language complexity of the parser, the ability of the refinement inference system to synthesize program invariants, and the expressiveness of the underlying abstract value domains.

Example 2.3. To give an idea of how grammar inference works, we present a brief example.

The program shown in Figure 3a is a simple parser written in Python that checks if the first character of the input string is an 'a' and, if so, asserts that the length of the string must be one and thus the string must be exactly "a"; otherwise, if the first character is not an 'a', then the program asserts that the second character must be a 'b', with no further restrictions on the input string. Note that the Python string indexing operations s[0] and s[1] are themselves types of assertions, since they will cause the program to crash if the index is out of bounds. This behavior is captured in the λ_{Σ} equivalent of the source program via its function axioms (Table 1).

Figure 3b shows the parser's top-level typing derivation (§ 3). The refinement inference judgement, applied to the parser function and its axioms, results in an incomplete refinement type containing an unsolved κ variable, and a verification condition of the form $\forall s. \kappa(s) \Rightarrow \varphi$, which tells us that this κ variable is indeed a grammar variable. We can see that the VC's consequent φ captures all of the parsers explicit and implicit constraints over its program variables.

Finally, Figure 3c shows a (possible) step-by-step reduction of φ into a simple quantifier-free predicate using abstract interpretation (§ 4). First, the innermost quantifiers $\forall v_1, \forall n, \forall v_2$, and $\forall y$ are eliminated and their quantified variables replaced by semantically equivalent expressions. Then we eliminate $\forall p_1$, followed by $\forall x$. At this point, there are no more quantified variables and we can fully abstract each occurrence of the only free variable *s*. Finally, we normalize the quantifier-free predicate into a single abstract string relation, equivalent to a regular expression.

3 Refinement Inference

We now give a formalization of λ_{Σ} , our core abstraction for representing ad hoc parsers. The language and its type system were heavily inspired by the SPRITE tutorial language by Jhala and Vazou [2020], and incorporate ideas from various other refinement type systems [Cosman and Jhala 2017; Dunfield and Krishnaswami 2021; Montenegro et al. 2020]. Our main contribution in this area is an extended refinement variable solving procedure that synthesizes precise grammars for input string constraints (§§ 3.3 and 4).

197

198

199

200 201

216

217

218

219

220

221

222

223

224

225

226

227

228

229

230

231

232

233

234

235 236 237

def parser(s): parser = $\lambda(s : \mathbb{S})$. **if** s[0] == "a": let $x = charAt \ s \ 0$ in assert len(s) == 1 let $p_1 = eqChar x 'a' in$ else: if p_1 then **assert** s[1] == "b" let *n* = length *s* in let $p_2 = eq n 1$ in assert p_2 >>> parser("") IndexError else >>> parser("a") \checkmark let y = charAt s 1 in IndexError let $p_3 = eqChar y$ 'b' in >>> parser("b") AssertionError >>> parser("ab") assert p_3 >>> parser("bb") \checkmark (a) A Python program (left) and its λ_{Σ} equivalent (right). incomplete refinement verification condition axioms $\vdash \text{ parser } \nearrow \{s : \mathbb{S} \mid \kappa(s)\} \to \mathbb{1} \ \exists \ \forall s. \kappa(s) \Rightarrow \omega$ Г $\varphi \doteq (\forall v_1. v_1 = 0 \Rightarrow v_1 \ge 0 \land v_1 < |s|) \land$ $(\forall x. x = s[0] \Rightarrow (\forall p_1. p_1 = \text{true} \Leftrightarrow x = \text{`a'} \Rightarrow$ $(p_1 = \text{true} \Rightarrow (\forall n. n \ge 0 \land n = |s| \Rightarrow n = 1)) \land$ $(p_1 = \text{false} \Rightarrow (\forall v_2, v_2 = 1 \Rightarrow v_2 \ge 0 \land v_2 < |s|) \land$ $(\forall y. y = s[1] \Rightarrow y = b'))))$ (b) An incomplete refinement typing of the λ_{Σ} program above. $\stackrel{1}{\rightsquigarrow} \langle |s| > 0 \land (\forall x. \ x = s[0]) \Rightarrow (\forall p_1. \ p_1 = \text{true} \Leftrightarrow x = \text{`a'} \Rightarrow$ $(p_1 = \text{true} \Rightarrow |s| = 1) \land (p_1 = \text{false} \Rightarrow |s| > 1 \land s[1] = b))$ $\stackrel{2}{\sim}$ $|s| > 0 \land (\forall x. x = s[0] \Rightarrow (x = a^{\prime} \land |s| = 1) \lor (x \neq a^{\prime} \land |s| > 1 \land s[1] = b^{\prime}))$ $|s| > 0 \land ((s[0] = a \land |s| = 1) \lor (s[0] \neq a \land |s| > 1 \land s[1] = b))$ $s \in \Sigma^+ \land ((s \in a\Sigma^* \land s \in \Sigma)) \lor (s \in (\Sigma \setminus a)\Sigma^* \land s \in \Sigma^2\Sigma^* \land s \in \Sigma b\Sigma^*))$ $\stackrel{5}{\sim}$ $s \in (a + (\Sigma \setminus a)b\Sigma^*)$ (c) Possible steps in the abstract interpretation of φ : (1) eliminate $\forall v_1, \forall n, \forall v_2, \forall y$; (2) eliminate $\forall p_1$; (3) eliminate $\forall x$; (4) abstract s; (5) normalize. Highlights indicate changes relative to previous step.

Fig. 3. An example of grammar inference.

1:6

3.1 Syntax 295

301

305

307

323

296 λ_{Σ} is a small λ -calculus in A-normal form (ANF) [Bowman 2022; Flanagan et al. 1993]. It exists 297 solely for type synthesis and its programs are neither meant to be executed nor written by hand. 298 Its syntax, collected in Figure 4, is thus minimal and has few affordances.

299 Values are either primitive constants or variables. Terms are comprised of values and the usual 300 constructs: function applications, function abstractions, bindings, recursive bindings, and branches. Applications are in ANF as a natural result of the SSA translation performed on the original source 302 code (see \S 2.1), but we also generally enforce this in the syntax to simplify the typing rules (\S 3.2). 303 λ_{Σ} terms have no type annotations, except on λ -binders and recursive **rec**-binders, whose (base) 304 types we assume are inferred in the pre-processing phase or provided by the programmer.

The primitive **Base Types** are 1 (unit), \mathbb{B} (Boolean), \mathbb{Z} (integer), $\mathbb{C}\mathbb{h}$ (character), and \mathbb{S} (string). 306 **Types** are formed by decorating base types with refinement predicates, or by constructing (dependent) function types, whose output types can refer to input types.

308 Predicates are terms in a Boolean logic, with the usual Boolean connectives, plus (in)equality 309 relations and arithmetic comparisons between predicate expressions, membership queries for 310 regular expression matching, applications of unknown κ variables, and existential quantifiers, 311 which arise during the FUSION phase of κ solving (§ 3.3). Both κ applications and existentials are not 312 part of the user-visible surface syntax. Expressions within predicates are built from lifted values 313 and functions, in particular linear integer arithmetic and operations over strings, e.g., length |□|, 314 character-at-index $\Box[\Box]$, substring $\Box[\Box..\Box]$, etc. To simplify the presentation, predicate expressions 315 are not further syntactically stratified, but are assumed to always occur well-typed (which is assured 316 by the implementation).

317 VC generation (§ 3.2) results in Constraints that are Horn clauses [Bjørner et al. 2015] in 318 Negation Normal Form (NNF), basically tree-like conjunctions of refinement predicates, where 319 each root-to-leaf path is a Constrained Horn Clause (CHC). This representation of VCs is due to 320 Cosman and Jhala [2017], who cleverly employ the constraints' nested scoping structure to make 321 κ solving tractable. 322

3.2 Type System

324 The main purpose of the type system of λ_{Σ} is to generate constraints, in particular input string 325 constraints for parser functions. Thus, the type system is focused on inference/synthesis, rather 326 than type checking. Terms need only be minimally annotated, at λ -abstractions and recursive 327 bindings. The types of applied functions need to be known, however, and available in the typing 328 context (see the discussion of axioms, § 2.1 and Table 1). Our system borrows heavily from Liquid 329 Haskell [Vazou et al. 2014] and the expositions given by Cosman and Jhala [2017] and Jhala and 330 Vazou [2020]. We combine typing rules and VC generation into one syntax-driven declarative 331 system of inference rules, given in Figure 5 and discussed below. 332

 $t_1 \leq t_2 \exists c$ **Subtyping.** A type t_1 is a subtype of t_2 (meaning, the values denoted by t_1 are 333 subsumed by t_2), if the entailment constraint c is satisfied. In the SUB/BASE case, this means the 334 335 refinement predicate of t_1 must imply the predicate of t_2 , for all possible values the types can have. 336 In the SUB/FUN case, where the contra-variant input constraint is joined with the co-variant output 337 constraint, we add an additional implication to the output constraint, strengthening it with the 338 supertype's input predicate. This is done using a generalized implication operation $(x :: t) \Rightarrow c$, which simply ensures that only base types are bound to quantifiers in the refinement logic. 339

 $\Gamma \vdash t \triangleright \hat{t}$ Template Generation. To enable complete type synthesis for all intermediate terms, it 341 is sometimes necessary to turn a type t into a template \hat{t} , where the refinement predicate is denoted 342

343

344 345 variables Values υ ::= x, y, z, \ldots 346 ... varies ... constants 347 348 Terms value е ••= 7) 349 application e v 350 $\lambda(x:b).e$ abstraction 351 let $x = e_1$ in e_2 binding 352 **rec** $x : t = e_1 e_2$ recursion 353 if v then e_1 else e_2 branch 354 355 Base Types 1 | B | Z | Ch | S b ::= 356 357 Types t ::= $\{x : b \mid p\}$ refined base 358 $(x:t_1) \rightarrow t_2$ dependent function 359 360 true | false **Boolean** constants **Predicates** p ::= 361 connectives $p_1 \wedge p_2 \mid p_1 \vee p_2$ 362 $p_1 \Rightarrow p_2 \mid p_1 \Leftrightarrow p_2$ implications 363 negation ¬p 364 (in)equality $w_1 = w_2$ $w_1 \neq w_2$ 365 $w_1 < w_2 \mid w_1 \leq w_2$ arithmetic comparison 366 $w \in RE$ regular language membership 367 κ application $\kappa(\overline{v})$ 368 $\exists (x:b).t$ existential quantification 369 370 Expressions value **d=** w 7) 371 function $f(\overline{w})$ 372 373 Constraints С ::= predicate p 374 $C_1 \wedge C_2$ conjunction 375 $\forall (x:b). \ p \Rightarrow c$ universal implication 376 377 Fig. 4. Syntax of λ_{Σ} terms, types, and refinements. 378

by a placeholder variable whose resolution is deferred (see §§ 2.2, 3.3, and 4). The rule KAP/BASE introduces a fresh κ variable, representing an *n*-ary relation between the type itself and all variables in the current environment Γ . Usually this environment is empty, but if *t* is a function, KAP/FUN recursively generates input and output templates, extending the environment along the way.

 $[Γ \vdash e \land t \neq c]$ *Type/Constraint Synthesis.* Given a typing context Γ mapping values to types, and a term *e*, we can synthesize a type *t* whose correctness is implied by the constraint *c*. The rule Syn/Con synthesizes built-in primitive types, denoted by prim(*c*) in the obvious way, e.g., prim(0) $\stackrel{\text{def}}{=} \{v : \mathbb{Z} \mid v = 0\}$ and so on. Syn/VAR retrieves the type of a variable from the current context, using *selfification* to produce the most precise possible type [Ou et al. 2004]. The self function lifts the variable into the refinement, allowing each occurrence of the variable in different

379 380 381

382

383

384 385

386

387

388

389

390

Proc. ACM Program. Lang., Vol. 1, No. OOPSLA, Article 1. Publication date: January 2018025-03-25 18:26. Page 8 of 1-30.

 $|t_1 \leq t_2 \exists c$ Subtyping $\frac{1}{\{v_1:b \mid p_1\} \leqslant \{v_2:b \mid p_2\} \exists \forall (v_1:b). p_1 \Rightarrow p_2[v_2 \coloneqq v_1]}$ SUB/BASE $\frac{s_2 \leqslant s_1 \exists c_i \qquad t_1[x_1 \coloneqq x_2] \leqslant t_2 \exists c_o}{(x_1:s_1) \to t_1 \leqslant (x_2:s_2) \to t_2 \exists c_i \land ((x_2:s_2) \Rightarrow c_o)}$ SUB/FUN $(x::t) \Rightarrow c \stackrel{\text{def}}{=} \begin{cases} \forall (x:b). \ p[v:=x] \Rightarrow c & \text{if } t \equiv \{v:b \mid p\}, \\ c & \text{otherwise.} \end{cases}$ $\boxed{\Gamma \vdash t \triangleright \hat{t}} \quad \text{Template Generation}$ $t \triangleright \hat{t} \stackrel{\text{def}}{=} \varnothing \vdash t \triangleright \hat{t}$ $\frac{\kappa \text{ is a fresh variable of sort } b \times \overline{t}}{\overline{x:t} \vdash \{v:b \mid p\} \triangleright \{v:b \mid \kappa(v,\overline{x})\}} \text{ Kap/Base } \frac{\Gamma \vdash t_1 \triangleright \hat{t}_1 \qquad \Gamma[x \mapsto t_1] \vdash t_2 \triangleright \hat{t}_2}{\Gamma \vdash (x:t_1) \to t_2 \triangleright (x:\hat{t}_1) \to \hat{t}_2} \text{ Kap/Fun}$ $\Gamma \vdash e \nearrow t \exists c$ **Type/Constraint Synthesis** $\frac{\Gamma(x) = t}{\Gamma \vdash x \not\subset \text{self}(x, t) \neq \text{true}} \xrightarrow{\text{Syn/Var}} \frac{\text{prim}(c) = t}{\Gamma \vdash c \not\subset t \neq \text{true}} \xrightarrow{\text{Syn/Con}}$ $\frac{\Gamma \vdash e \nearrow (y:t_1) \to t_2 \exists c_e \qquad \Gamma \vdash v \nearrow t_v \qquad t_v \leqslant t_1 \exists c_v}{\Gamma \vdash e v \nearrow t_2 [y:=v] \exists c_v \land c_v} \text{Syn/App}$ $\begin{array}{c|c} \tilde{t}_1 \succ \hat{t}_1 & \Gamma[x \mapsto \hat{t}_1] \vdash e \nearrow t_2 \exists c_2 \\ \hline \Gamma \vdash \lambda(x : \tilde{t}_1), e \nearrow (x : \hat{t}_1) \rightarrow t_2 \exists (x :: \hat{t}_1) \Rightarrow c_2 \end{array} \\ \begin{array}{c} \text{Syn/Lam} \end{array}$ $\frac{\Gamma[x \mapsto t_1] \vdash e_2 \nearrow t_2 \exists c_2 \qquad t_2 \triangleright \hat{t}_2 \qquad t_2 \leqslant \hat{t}_2 \exists \hat{c}_2}{\Gamma \vdash \mathbf{let} \ x = e_1 \ \mathbf{in} \ e_2 \nearrow \hat{t}_2 \exists c_1 \land ((x :: t_1) \Rightarrow c_2 \land \hat{c}_2)} \operatorname{Syn/Let}$ $\frac{\begin{array}{cccc} & \mathbf{r}_{1} \succ \mathbf{r}_{1} \\ & \Gamma[x \mapsto t_{1}] \vdash e_{1} \nearrow t_{1} \exists c_{1} \\ & \Gamma[x \mapsto t_{1}] \vdash e_{2} \nearrow t_{2} \exists c_{2} \\ \hline \Gamma \vdash \mathbf{rec} \ x : \tilde{t}_{1} = e_{1} \ e_{2} \nearrow \hat{t}_{2} \exists ((x :: \hat{t}_{1}) \Rightarrow c_{1} \land \hat{c}_{1}) \land ((x :: t_{1}) \Rightarrow c_{2} \land \hat{c}_{2}) \end{array}}{\left(\mathbf{Syn/Rec} \right)} \operatorname{Syn/Rec}$ $\frac{\Gamma \vdash x \nearrow \mathbb{B} \qquad \begin{array}{ccc} \Gamma \vdash e_1 \nearrow t_1 \dashv c_1 \qquad t_1 \rhd \hat{t} \qquad t_1 \leqslant \hat{t} \dashv \hat{c}_1 \\ \Gamma \vdash e_2 \nearrow t_2 \dashv c_2 \qquad t_2 \leqslant \hat{t} \dashv \hat{c}_2 \end{array}}{\Gamma \vdash \mathbf{if} \ x \ \mathbf{then} \ e_1 \ \mathbf{else} \ e_2 \nearrow \hat{t} \dashv (x = \mathsf{true} \Rightarrow c_1 \land \hat{c}_1) \land (x = \mathsf{false} \Rightarrow c_2 \land \hat{c}_2)} \\ \end{array} \\ \xrightarrow{\mathsf{Syn/IF}}$

Fig. 5. Typing rules for λ_{Σ} .

1:10

branches of the program to be precisely typed.

self(x, t)
$$\stackrel{\text{def}}{=} \begin{cases} \{v : b \mid p \land v = x\} & \text{if } t \equiv \{v : b \mid p\}, \\ t & \text{otherwise.} \end{cases}$$

SYN/APP synthesizes the result type of a function application, where the given input can be a subtype of the function's declared input type. In the result type, the declared input variable is replaced by the given input, using standard capture-avoiding substitution. Note that because our terms are in ANF, no arbitrary expressions are introduced into the type during this substitution. The synthesized VC is simply a conjunction of the function's VC and the constraint created by the subtyping judgement. Syn/Lam produces a function type whose input refinement is a fresh κ variable, based on the annotation on the λ -binder. Syn/Let first synthesizes the type t_1 for the bound term, then the type t_2 for the expression's body under an environment where x is bound to t_1 . To ensure that x does not escape its scope, the whole **let**-expression is given the templated type \hat{t}_2 , a supertype of t_2 with a κ variable in place of its refinement. Syn/Rec is similar to Syn/Let, with the addition that we first assume the bound term's type as a placeholder \hat{t}_1 based on the annotation \tilde{t}_1 on the binder, before synthesizing it as t_1 , allowing for recursion in the bound term. SYN/IF synthesizes the type of the whole conditional as a supertype of both branches. In the VC, we imply the **then**-branch's constraints if the condition is true, and the **else**-branch's constraints if the condition is false. The κ variable in the templated supertype \hat{t} allows this path-sensitive information to travel back upwards.

3.3 Variable Solving

Before the synthesized VCs can be sent off to an SMT solver to be proven valid, we have to replace all κ variables with concrete refinement predicates. Some of these variables represent input string refinements and this is the point where we need to find the grammars describing those strings.

Algorithm 1 gives a high-level overview of the VC solving procedure. We start with an incomplete VC constraint *c* (i.e., a constraint containing κ variables) and an initial κ assignment σ (usually empty). We use FUSION [Cosman and Jhala 2017] and predicate abstraction [Rondon et al. 2008] to deal with non-grammar κ variables. For grammar variables, we use our abstract interpretation approach detailed in § 4 to find an abstract solution for κ , i.e., a grammar. If the complete VC, with all κ variables replaced by their solutions, is satisfiable (modulo theories) then the inferred type is valid and σ contains concrete assignments for all of the type's refinement variables. Otherwise, no valid solution could be found, either because there does not exist one (e.g., the program has a type error) or due to insufficient invariants or limits of the abstract domain.

4 Grammar Inference

Given a program *P* with a partially inferred refinement type $\{s : S \mid \kappa(s)\} \rightarrow \mathbb{1}$ and an incomplete verification condition of the form $\forall (s : S). \kappa(s) \Rightarrow \varphi$, our goal is to find an assignment for κ that both validates the VC and meaningfully refines the type of *s*. Trivially, the assignment $\kappa(s) \mapsto$ false always validates the VC, but is hardly a meaningful refinement. Assuming that *P* is a parser, we ideally want an assignment for κ that constrains *s* in a way that

- (1) disallows any string rejected by *P* (soundness) and
- (2) allows all strings accepted by *P* (completeness).

In other words, we want a refinement for *s* that is a *grammar* for *P*, i.e., a finite description of the (possibly infinite) set of strings contained in the language L(P).

Luckily, the VC's consequent φ already appears to be what we are looking for: it is a first-order formula over a free string variable *s* that exactly captures the semantics of the parsing operations

Algorithm 1 Solve an incomplete VC and find all κ assignments 491 492 1: procedure SOLVE(c, σ) 493 $c \leftarrow \sigma(c)$ 2: 494 $c \leftarrow \text{Fusion}(c)$ \triangleright eliminate local acyclic κ variables [Cosman and Jhala 2017] 3: 495 $\sigma \leftarrow \sigma \cup \text{Liouid}(c)$ \triangleright solve residual non-grammar κ variables [Rondon et al. 2008] 4: 496 $c \leftarrow \sigma(c)$ 5: 497 **for all** constraints in *c* of the form $\forall (s : \mathbb{S})$. $\kappa(s) \Rightarrow \varphi$ **do** 6: 498 $\hat{\sigma} \leftarrow \llbracket \varphi \rrbracket^{\sharp}$ ▶ infer grammar using abstract interpretation [§ 4] 7: 499 $\sigma \leftarrow \sigma \left[\kappa(s) \mapsto s \in \hat{\sigma}(s) \right]$ 8: 500 $c \leftarrow \sigma(c)$ 9: 501 10: end for 502 if SATISFIABLE(c) then return VALID else return INVALID 11: 503 12: end procedure 504 505

performed by *P* on *s*. Assuming that φ itself does not contain any other unsolved κ variables (these having been eliminated by prior inference and solving steps, e.g., FUSION and predicate abstraction, see § 3.3), as well as the correctness of all involved axiomatic string function specifications (Table 1), then, under some model \mathfrak{M} of the refinement logic (e.g., linear integer arithmetic and basic string operations), by construction,

 $\mathfrak{M}, [s \mapsto t] \models \varphi \quad \Longleftrightarrow \quad P \text{ successfully parses } t \quad \Longleftrightarrow \quad t \in L(P).$

The parser *P* represented by φ accepts some particular string *t* if and only if φ is true under an assignment of *s* to *t*. The set of all strings that satisfy φ is exactly the language L(P) accepted by *P*. Thus, φ is technically a complete grammar for *P*.

⁵¹⁷ Practically, however, φ makes for a rather poor grammar. It is a first-order formula close in size ⁵¹⁸ and structure to the parsing program *P* itself. While $\kappa(s) \mapsto \varphi$ is an ideal assignment in the sense ⁵¹⁹ that it is both sound and complete, it results in a rather impractical refinement that would puzzle a ⁵²⁰ human programmer. The presence of quantifiers within φ complicates any further analysis.

⁵²¹ In order to obtain a better grammar for L(P), we will use φ as a starting point to derive a quantifier-⁵²² free predicate that is much simpler than φ but semantically equivalent. They key components of ⁵²³ our approach are:

- (1) an abstract interpretation of certain first-order string formulas as grammars (§ 4.1),
- (2) an extended constraint syntax that mixes symbolic and abstract representations (§ 4.2),
- (3) an abstract semantics of relations between predicate expressions (§ 4.3),
- (4) a procedure for eliminating quantified variables by abstract substitution (§ 4.4),
- (5) abstract value representations of all base types (§ 5).
- 4.1 Abstract Interpretation

Background. Abstract interpretation [Cousot and Cousot 1977, 1979] is a well-established framework for formalizing static program analyses. It involves the sound approximation of all possible states of a program, usually trading precision for efficiency. The concrete semantics of a program \mathcal{P} are defined by a semantic function $[\Box] : \mathcal{P} \to \wp(C)$, which produces a power set of concrete values *C*. The concrete domain $\wp(C)$ can be approximated by an abstract domain \mathcal{A} , with an abstraction function $\alpha : \wp(C) \to \mathcal{A}$ and a concretization function $\gamma : \mathcal{A} \to \wp(C)$ mapping elements between the domains. Usually, \mathcal{A} is a complete lattice $\langle \mathcal{A}, \sqsubseteq, \sqcap, \sqcap, \sqcup, \bot, \top \rangle$ and the two domains form

539

506

507

508

509

510

511

512

513

524

525

526

527

528

529 530

Table 2. Summary of abstract domains in PANINI

base type				abstract domain	
unit	$\wp(1)$	=	î	the two-element lattice	§ 5.1
Booleans	$\wp(\mathbb{B})$	=	Â	Boolean subset lattice	§ 5.2
integers	$\wp(\mathbb{Z})$	⊇	Ź	open-ended interval lists	§ 5.3
characters	$\wp(\mathbb{Ch})$	=	Ĉ	Unicode character sets	§ 5. 4
strings	$\wp(\mathbb{S})$	⊇	Ŝ	regular expressions over \hat{C}	§ 5.5

a Galois connection $\langle \wp(C), \subseteq \rangle \xrightarrow[\alpha]{\alpha} \langle \mathcal{A}, \sqsubseteq \rangle$, which intuitively means that relationships between elements of $\wp(C)$ also hold between the corresponding abstracted elements of \mathcal{A} . We can then define an abstract semantics $\llbracket \square \rrbracket^{\sharp}: \mathcal{P} \to \mathcal{A}$ to abstractly interpret programs and directly produce abstract values. The abstract interpretation is *sound* iff, for all programs \mathcal{P} , $\alpha(\llbracket \mathcal{P} \rrbracket) \sqsubseteq \llbracket \mathcal{P} \rrbracket^{\sharp}$, and it is also *complete* iff $\alpha(\llbracket \mathcal{P} \rrbracket) = \llbracket \mathcal{P} \rrbracket^{\sharp}$. Completeness in this context means that the abstract semantics incurs no loss of precision relative to the underlying abstract domain, i.e., the abstract semantics can take full advantage of the whole domain. Finally, an abstract interpretation is *exact* iff $\llbracket \mathcal{P} \rrbracket = \gamma(\llbracket \mathcal{P} \rrbracket^{\sharp})$. meaning that the abstraction loses no information and the abstract semantics exactly captures the concrete semantics of the program [Cousot 1997; Giacobazzi and Quintarelli 2001].

Parser Semantics. In our case, the kind of program we want to abstractly interpret is a parser P represented by a first-order formula φ . We take the concrete semantics $[\![\varphi]\!]$ to be an assignment σ of free variables to values, and in particular denote the parser's input string by the free variable *s*, such that

 $\mathfrak{M}, \sigma \models \varphi \iff \sigma \in \llbracket \varphi \rrbracket \iff \sigma(s) \in L(P).$

Trying to compute $\llbracket \varphi \rrbracket$ directly will generally result in an infinite set of values. Hence our desire for an abstract semantics $\llbracket \varphi \rrbracket^{\sharp}$ that produces a finite approximation of this set such that

$$\mathfrak{M}, \hat{\sigma} \models \varphi \quad \longleftrightarrow \quad \hat{\sigma} = \llbracket \varphi \rrbracket^{\sharp} \quad \Longrightarrow \quad \hat{\sigma}(s) \subseteq L(P)$$

where $\hat{\sigma}$ is an assignment of free variables to abstract values. In particular, $\hat{\sigma}(s)$ is now a grammar describing (a subset of) the strings in L(P).

The completeness of the approximation $\hat{\sigma}$ depends on the underlying abstract domains. We give a brief summary of the domains currently used by PANINI in Table 2 and provide complete definitions in § 5. Note that our abstract string domain \hat{S} represents sets of strings with regular expressions. This means that our approach must under-approximate any L(P) above regular in the Chomsky hierarchy [Chomsky and Schützenberger 1963]. For super-regular languages, we can only infer a partial grammar representing a regular subset, if one exists. However, if L(P) is (at most) regular, then we can infer a complete grammar for P.

Using our abstract interpretation, we can construct a finite but quantifier-free solution for the refinement variable,

 $\kappa(s) \mapsto s \in \hat{\sigma}(s), \quad \text{where } \hat{\sigma} = \llbracket \varphi \rrbracket^{\sharp}.$

The definition of the abstract semantics function $\llbracket \Box \rrbracket^{\sharp}$ is given in Figure 6. It depends on a novel variable-focused relational semantics $\llbracket \Box \rrbracket^{\uparrow}_{\Box}$ and a quantifier elimination procedure qelim(\Box). We describe these in the remainder of this section, after establishing some syntactic extensions to λ_{Σ} constraints.

589	
590	$\llbracket a \rrbracket^{\sharp} \div \lbrace r \mapsto \llbracket a \rrbracket^{\uparrow} \mid r \in \operatorname{vare}(a) \rbrace$
591	$\llbracket \Psi \rrbracket = \{x + \gamma \llbracket \Psi \rrbracket x \mid x \in \operatorname{var} \mathfrak{s}(\Psi) \}$
592	
593	$\llbracket P_1 \land P_2 \rrbracket \uparrow_x \doteq \llbracket P_1 \rrbracket \uparrow_x \sqcap \llbracket P_2 \rrbracket \uparrow_x \text{if } P_1, P_2 \text{ quantifier-free}$
594	$[P_1 \lor P_2] \uparrow_x \doteq [P_1] \uparrow_x \sqcup [P_2] \uparrow_x$ if P_1, P_2 quantifier-free
595	$\begin{bmatrix} -D \end{bmatrix}^{\uparrow} \doteq = \begin{bmatrix} D \end{bmatrix}^{\uparrow} \qquad \qquad$
596	$\ \neg r \ _{x} = \neg \ r \ _{x}$ If r quantifier-free
597	$\llbracket \varphi \rrbracket \uparrow_x \doteq \llbracket \operatorname{qelim}(\varphi) \rrbracket \uparrow_x$
598	(ω) as defined by domain:
599	$\llbracket \rho \rrbracket \uparrow_x \doteq \begin{cases} \omega, & \text{ab defined by domain,} \\ (\omega, \sigma) & \text{at herming} \end{cases}$
600	$(\langle x: \rho \rangle, \text{ otherwise.}$
601	
602	Fig. 6. Abstract semantics of $\lambda_{\rm T}$ constraints
603	Fig. 0. Abstract semantics of π_{Σ} constraints.
604	
605	
606	Constraints φ ::= true false Boolean constants
607	$ \varphi_1 \wedge \varphi_2 \varphi_1 \vee \varphi_2$ connectives
608	$ \varphi_1 \Rightarrow \varphi_2 \varphi_1 \Leftrightarrow \varphi_2$ implications
609	$\neg \varphi$ negation
610	$\exists x. \varphi$ existential quantification
611	$\forall x. \varphi$ universal quantification
612	ρ relations
613	
614	Relations ρ ::= $\omega_1 \bowtie \omega_2$ where $\bowtie \in \{=, \neq, <, \leq, \in, \notin, \emptyset, \ \}$
615	
616	Expressions $\omega ::= x, y, z, \dots$ variables
617	c concrete value
618	\hat{c} abstract value
619	$\langle x: \rho \rangle$ abstract relation
620	$f(\overline{\omega})$ function
621	
622	

Fig. 7. Extended syntax of λ_{Σ} constraints.

4.2 Extended Constraint Syntax

We extend the formal syntax of predicates and constraints from Figure 4 to include abstract values and relations. The extended syntax is given in Figure 7.

Constraints φ , in addition to Boolean constants and the usual logical connectives, include both universal and existential quantification, and generalized relations between predicate expressions. There are no κ applications at this point. The base types of all bound variables are known, but we elide them from the presentation to reduce clutter.

Expressions ω , in addition to variables and concrete values, now include abstract values \hat{c} and abstract relations $\langle x: \rho \rangle$. Abstract values \hat{c} are taken from our abstract domains (§ 5) and are representations of potentially infinite sets of concrete values. For example, the abstract value $[1, \infty]$ represents an infinite set of integers $\{1, 2, 3, ...\}$ and the abstract string Σ^* a represents the set of all

strings that end with the character 'a'. For functions, we assume the usual notational conveniences, e.g., we write a + b for +(a, b) and |s| for str.len(s) and so on. If an abstract value appears in a function, it lifts the whole expression into the abstract realm, e.g., $x + [1, \infty]$ represents the set of expressions {x + 1, x + 2, x + 3, ...}. Abstract relations $\langle x : \rho \rangle$ similarly represent sets of values, specifically those values that are described by the given relation. For example, the abstract relation $\langle x : x > 5 \rangle$ represents "the set of all values x for which x > 5 is true" and $\langle i : s[i] = c \rangle$ represents the set of all indexes *i* at which the string in variable *s* contains the character in variable *c*.

Relations ρ are comprised of the usual (in)equality and arithmetic comparisons, as well as set (non)inclusion (\in , \notin) and (non-)empty intersection (\emptyset , \parallel). The latter are simply abbreviations for common set operations, with

648 649 650

655

656

657

- $A \notin B \equiv A \cap B \neq \emptyset$, meaning "A and B have at least one element in common,"
- $A \parallel B \equiv A \cap B = \emptyset$, meaning "A and B have no elements in common."

Since abstract values are essentially sets, relating them works the same way. Note the distinction between $x \in [0, \infty]$ ("x is a member of the set of natural numbers"), $x = [0, \infty]$ ("x is the set of natural numbers"), and $x \notin [0, \infty]$ ("x has at least one element in common with the set of natural numbers").

Normalization. Both expressions and relations can be normalized by partial evaluation. If abstract values are involved, the semantics of the particular abstract domains apply. If possible, relations are fully abstracted using their relational semantics (§ 4.3).

$$1 + x - 2 \quad \rightsquigarrow \quad x - 1$$

$$|\text{``abc''}| \quad \rightsquigarrow \quad 3$$

$$x - 1 \in [1, 5] \quad \rightsquigarrow \quad x \in [1, 5] + 1 \quad \rightsquigarrow \quad x \in [2, 6]$$

$$x < 5 \quad \rightsquigarrow \quad x \in [-\infty, 4]$$

$$|s| + 1 > 2 \quad \rightsquigarrow \quad |s| > 1 \qquad \rightsquigarrow \quad s \in \Sigma^{+}$$

$$[1, 2] \emptyset \langle i: s[i] = `a` \rangle \quad \rightsquigarrow \quad s[[1, 2]] \ni `a` \quad \rightsquigarrow \quad s \in \Sigma\Sigma?a\Sigma^{*}$$

In the remainder, we assume that expressions and relations are always fully normalized.

4.3 Relational Semantics

The concrete semantics $[\![\rho]\!]$ of a single relation are all those assignments of the relation's free variables that make the relation true, e.g.,

671 672 673

674 675

676

677

678

679

680

681

682

683

684 685 686

666

667

668

669 670

$$\llbracket x > 3 \rrbracket \doteq \{x \mapsto 4, x \mapsto 5, \dots\},$$

$$\llbracket [|s| > 3 \rrbracket \doteq \{\dots, s \mapsto \text{``abaa''}, s \mapsto \text{``abab''}, \dots\}$$

$$\llbracket |s| > x \rrbracket \doteq \left\{\dots, \begin{pmatrix} s \mapsto \text{``ab''} \\ x \mapsto 0 \end{pmatrix}, \begin{pmatrix} s \mapsto \text{``ab''} \\ x \mapsto 1 \end{pmatrix}, \dots \right\}.$$

Unfortunately, constructing an abstract semantics even for such simple relations is decidedly nontrivial. It requires complex relational domains to capture just a limited set of constraints between multiple variables [Cousot and Halbwachs 1978; Logozzo and Fähndrich 2010; Simon et al. 2003].

Fortunately, we do not actually need a full abstract semantics that condenses relations into singular abstract values. For our purposes, it is sufficient to consider relational semantics on a per-variable basis. To this end, we introduce a variable-focused abstract semantics function $[\![\rho]\!]\uparrow_x$ that for a given variable *x* occurring free in ρ produces an abstract expression whose concrete values are exactly those that could be substituted for *x* to make ρ true, i.e.,

$$\mathfrak{M}, [x \mapsto c] \models \rho \quad \iff \quad c \in \llbracket \rho \rrbracket \uparrow_x.$$

Abstract relations $\langle x: \rho \rangle$ provide a convenient "default" implementation,

$$\llbracket \rho \rrbracket \uparrow_x \doteq \langle x \colon \rho \rangle,$$

since by definition they exactly capture the abstract semantics of the relation ρ for the given variable *x*. But we can often do better than this. Depending on the underlying abstract domains and domain-specific knowledge about the involved operations, more specific definitions of $[\Box] \uparrow_\Box$ are possible (§ 5). For example, for the domains of regular expressions and integers, we can define precise abstractions for simple string length relations,

$$\llbracket |s| = n \rrbracket \uparrow_s \doteq \Sigma^n, \qquad \llbracket |s| \ge n \rrbracket \uparrow_s \doteq \Sigma^n \Sigma^*, \qquad \llbracket |s| \le n \rrbracket \uparrow_s \doteq \Sigma^0 + \Sigma^1 + \dots + \Sigma^n$$

where n is a meta-variable standing for some concrete integer value. With these and other domainspecific semantics, we can construct variable-focused abstractions for the earlier examples. Note how in all cases the abstraction effectively eliminates the chosen variable:

$$\llbracket x > 3 \rrbracket \uparrow_x \doteq \llbracket 4, \infty \rrbracket \qquad \llbracket |s| > 3 \rrbracket \uparrow_s \doteq \Sigma^4 \Sigma^* \qquad \llbracket |s| > x \rrbracket \uparrow_s \doteq \langle s \colon |s| > x \rangle$$
$$\llbracket |s| > x \rrbracket \uparrow_x \doteq |s| - \llbracket 1, \infty \rrbracket$$

4.4 Quantifier Elimination

Figure 8 presents our procedure for eliminating quantifiers by abstract substitution. The top-level function qelim traverses a given constraint to eliminate all of its quantifiers. During this traversal, we apply standard logical transformations to simplify the problem.

$$\begin{array}{l} \forall x.\varphi \iff \neg \exists x. \neg \varphi & (\text{DeMorgan's law}) \\ \exists x. \lor \land \rho \iff \lor \exists x. \land \rho & (\text{distributivity of } \exists \text{ over } \lor) \end{array}$$

714 The actual variable elimination happens in gelim1, where we only need to consider a single 715 conjunctive set of relations R. To eliminate a particular variable x from the relations in this set, we 716 first compute the *x*-relative relational semantics $[\rho] \uparrow_x$ for all relations ρ in the subset of relations in 717 R that contain x as a free variable. This results in a set E of (abstract) expressions, all representing 718 some aspect of x in R. The expressions in E do not contain the variable x but they might contain 719 other free variables. We now generate all unordered pairwise combinations of expressions in E 720 and create a new equality relation between each pair of expressions (using a relational operator 721 appropriate for the pair's types), resulting in a set of relations \hat{R} that contain no reference to x yet 722 preserve the semantics of the original conjunction of relations. Finally, we return \hat{R} together with 723 those relations in *R* that did not contain *x* in the first place. 724

5 Abstract Domains

We now define abstract domains for each of the base types in our system. Each domain efficiently captures (possibly infinite) sets of concrete values of the corresponding type, lifts certain operations of the type to the abstract domain, and defines some abstract relational semantics $[\Box] \uparrow_{\Box}$ for expression involving those operations.

To simplify the definitions and avoid boilerplate repetition, we generally assume that expressions have already been arranged in a uniform manner and are fully normalized (§ 4.2). We also forego subscripting domain-specific operations and elements if there is no ambiguity, e.g., we write \top instead of differentiating \top_{1} , $\top_{\mathbb{B}}$, etc.

735

725

726

691

692

693

694

699

700

707

708

709

710 711 712

736 $\operatorname{qelim}(\exists x. \varphi) \doteq \bigvee \left\{ \operatorname{qelim}(x, R) \mid R \in \operatorname{dnf}(\operatorname{qelim}(\varphi)) \right\}$ 737 738 $\operatorname{qelim}(\forall x. \varphi) \doteq \operatorname{qelim}(\neg \exists x. \neg \varphi)$ 739 $\operatorname{qelim}(\varphi_1 \land \varphi_2) \doteq \operatorname{qelim}(\varphi_1) \land \operatorname{qelim}(\varphi_2)$ 740 $\operatorname{qelim}(\varphi_1 \lor \varphi_2) \doteq \operatorname{qelim}(\varphi_1) \lor \operatorname{qelim}(\varphi_2)$ $\operatorname{qelim}(\neg \varphi) \doteq \neg \operatorname{qelim}(\varphi)$ $\operatorname{qelim}(\rho) \doteq \rho$ $qelim(true) \doteq true$ $qelim(false) \doteq false$ $R) \doteq \hat{R} \cup \{ \rho \in R \mid x \notin vars(\rho) \}$ ere $\hat{R} = \left\{ \omega_1 \bowtie \omega_2 \middle| (\omega_1, \omega_2) \in {E \choose 2} \right\} \bowtie \in \{=, \in, \ni, \emptyset \}$ $\operatorname{qelim1}(x, R) \doteq \hat{R} \cup \{\rho \in R \mid x \notin \operatorname{vars}(\rho)\}$ where $E = \{ \llbracket \rho \rrbracket \uparrow_{x} \mid \rho \in R, x \in \operatorname{vars}(\rho) \}$

Fig. 8. Quantifier elimination

Unit 5.1

The abstract unit type $\hat{1}$ simply adds a bottom element, forming the two-element lattice, with

 $\alpha = \gamma = \mathrm{id}, \qquad [\![x = \mathrm{unit}]\!]\uparrow_x \doteq \mathrm{unit}, \qquad [\![x \neq \mathrm{unit}]\!]\uparrow_x \doteq \bot.$

5.2 **Booleans**

The Boolean subset lattice $\langle \varphi(\mathbb{B}), \subseteq \rangle$ is a complete abstraction of the Boolean values, adding \emptyset and {true, false} as bottom and top elements, respectively, and forming a complete complemented lattice via subset inclusion and set complement. For consistency with the other definitions, we call this abstraction $\hat{\mathbf{B}}$ and will use the notation $\langle \hat{\mathbf{B}}, \sqsubseteq, \bot, \top, \sqcup, \neg, \neg \rangle$ for the lattice elements and operations. The abstract semantics are defined by

$$\alpha = \gamma = \mathrm{id}, \qquad [\![x = b]\!]\uparrow_x \doteq \{b\}, \qquad [\![x \neq b]\!]\uparrow_x \doteq \{\neg b\}.$$

5.3 Integers

Any concrete set of contiguous integers $\{x \in \mathbb{Z} \mid a \le x \le b\}$ can be represented efficiently as an interval [a, b], even allowing for a or b to be $\pm \infty$. The domain of integer intervals

$$\mathbf{IZ} \stackrel{\text{der}}{=} \left\{ [a, b] \mid a, b \in \mathbb{Z} \cup \{-\infty, \infty\} \text{ and } a \leq b \right\}$$

forms a pseudo-semi-lattice $(IZ, \subseteq, \top, \sqcup, \sqcap)$ with a bounded meet but no \perp element:

	1011110 a poetado benni idence $(12, 2, 3, 3, 5)$ with a bounded incert such
778	•
779	$[a,b] \sqsubseteq [c,d] \iff c \le a \land b \le d$
780	$\top = [-\infty, \infty]$
781	$\begin{bmatrix} a & b \end{bmatrix} + \begin{bmatrix} c & d \end{bmatrix} = \begin{bmatrix} \min(a & c) & \max(b & d) \end{bmatrix}$
782	$\begin{bmatrix} u, v \end{bmatrix} \doteq \begin{bmatrix} v, u \end{bmatrix} = \begin{bmatrix} \min(u, v), \min(v, u) \end{bmatrix}$
783	$[a,b] \sqcap [c,d] = [\max(a,c),\min(b,d)]$
784	

763

764

765

766

767

768

773

774

We can also define some standard operations and convenient relations between intervals:²

$$[a,b] + [c,d] = [a+c,b+d]$$

$$[a,b] \text{ precedes } [c,d] \iff b < (c-1)$$

$$[a,b] - [c,d] = [a-d,b-c]$$

$$[a,b] \text{ is before } [c,d] \iff b < c$$

$$[a,b] \text{ contains } [c,d] \iff a \le c \land d \le b$$

$$[a,b] \text{ overlaps } [c,d] \iff a \le c \land c \le b \land b < d$$

While IZ can abstractly represent infinite contiguous sets, such as those defined by a single inequality relation like $\{x \in \mathbb{Z} \mid x > 5\}$, it can not represent even finite non-contiguous sets like $\{1, 5, 7\}$ or inequalities like $\{x \in \mathbb{Z} \mid x \neq 2\}$. Thus the connection between \mathbb{Z} and IZ is only partial:

$$\alpha : \varphi(\mathbb{Z}) \hookrightarrow \mathbf{IZ} \qquad \qquad \gamma : \mathbf{IZ} \to \varphi(\mathbb{Z})$$
$$\alpha(\{x \in \mathbb{Z} \mid a \le x \le b\}) = [a, b] \qquad \qquad \gamma([a, b]) = \{x \in \mathbb{Z} \mid a \le x \le b\}$$

To completely represent non-contiguous sets of integers, we can use ordered lists of nonoverlapping intervals; e.g., $\{1, 2, 3, 7, 8, ...\}$ can be represented as $[1, 3, 7, \infty]$. As long as the number of gaps between intervals is bounded, $\varphi(\mathbb{Z})$ can be efficiently abstracted by the domain

 $\hat{\mathbf{Z}} \stackrel{\text{def}}{=} \{ x_1 x_2 \cdots x_n \mid x_i \in \mathbf{IZ} \text{ and } x_i \text{ is before } x_{i+1} \text{ and } i \leq n \},\$

which forms a complete complemented lattice $\langle \hat{\mathbf{Z}}, \sqsubseteq, \top, \top, \sqcup, \neg, \neg \rangle$, with $\bot = \emptyset$ and $\top = [-\infty, \infty]$ and the operations defined below. We use Haskell list notation (x : xs) to peel off (or add on) the first interval x in a list, with xs denoting the remaining intervals.

$$(x:xs) \sqsubseteq (y:ys) \iff (x \sqsubseteq_{\mathbf{IZ}} y \land xs \sqsubseteq (y:ys)) \lor (y \text{ is before } x \land (x:xs) \sqsubseteq ys)$$

$$(x : xs) \sqcup (y : ys) \longleftrightarrow (x \sqcup y : ys)) \text{ if } x \text{ precedes } y \\ y : ((x : xs) \sqcup ys) = \begin{cases} x : (xs \sqcup (y : ys)) & \text{if } x \text{ precedes } y \\ y : ((x : xs) \sqcup ys) & \text{if } y \text{ precedes } x \\ ((x \sqcup_{\mathrm{IZ}} y) \sqcup xs) \sqcup ys) & \text{otherwise} \end{cases}$$

$$(x : xs) \sqcap (y : ys) = \begin{cases} x : (xs \sqcup (y : ys)) & \text{if } x \text{ precedes } x \\ ((x \sqcup_{\mathrm{IZ}} y) \sqcup xs) \sqcup ys) & \text{otherwise} \end{cases}$$

$$(x : xs) \sqcap (y : ys) = \begin{cases} x : (xs \sqcup (y : ys)) & \text{if } x \text{ is before } y \\ (x : xs) \sqcap ys) & \text{if } x \text{ is before } x \\ (x : xs) \sqcap ys) & \text{if } x \text{ contains } y \text{ or } y \text{ overlaps } x \\ (x \sqcap_{\mathrm{IZ}} y) : ((x : xs) \sqcap ys) & \text{if } x \text{ contains } x \text{ or } x \text{ overlaps } x \\ (x \sqcap_{\mathrm{IZ}} y) : (xs \sqcap (y : ys)) & \text{if } y \text{ contains } x \text{ or } x \text{ overlaps } y \\ (x \sqcap_{\mathrm{IZ}} y) : (xs \sqcap ys) & \text{otherwise} \end{cases}$$

$$\neg [-\infty, b_1 \mid a_2, b_2 \mid \dots \mid a_n, \infty] = [b_1 + 1, a_2 - 1 \mid b_2 + 1, a_3 - 1 \mid \dots \mid b_{n-1} + 1, a_n - 1] \\ \neg [-\infty, b_1 \mid a_2, b_2 \mid \dots \mid a_n, b_n] = [b_1 + 1, a_2 - 1 \mid \dots \mid b_{n-1} + 1, a_n - 1] \\ \neg [a_1, b_1 \mid a_2, b_2 \mid \dots \mid a_n, \infty] = [-\infty, a_1 - 1 \mid b_1 + 1, a_2 - 1 \mid \dots \mid b_{n-1} + 1, a_n - 1]$$

$$\neg [a_1, b_1 | a_2, b_2 | \dots | a_n, b_n] = [-\infty, a_1 - 1 | b_1 + 1, a_2 - 1 | \dots | b_{n-1} + 1, a_n - 1 | b_n + 1, \infty]$$

The standard arithmetic operations are lifted into $\hat{\mathbf{Z}}$ via pointwise mapping of the equivalent operations on IZ. We assume a standard semantics for abstract integer expressions, allowing us to reduce and rewrite arithmetic expressions into a canonical form, e.g., $[3, 5] + 1 \rightarrow [4, 6]$.

We can construct the abstraction function $\alpha : \varphi(\mathbb{Z}) \to \hat{Z}$ for any finite subset of integers $X \in \wp(\mathbb{Z})$ by first sorting all elements of X in ascending order and then identifying all non-overlapping intervals of consecutive integers. This strategy does not work if X is infinite, and in

²These interval relations are reminiscent of the temporal interval algebra of Allen [1983]. Our definitions of "precedes" and "overlaps" are identical to Allen's, whereas "is before" corresponds to Allen's "precedes or meets," and our "contains" is equivalent to Allen's "contains or equals or is started by or is finished by."

any case is rather inefficient. Likewise, the concretization function $\gamma : \hat{\mathbf{Z}} \to \wp(\mathbb{Z})$, defined as

834 835

836 837

838

839

840 841 842

843

844 845

846

847

848

849

850

851

852

853

854

855

856

857

858

859

860 861 862

863

864

865

866

867

868

$$\gamma(x_1x_2\cdots x_n) = \bigcup_{i\leq n} \gamma_{\mathrm{IZ}}(x_i),$$

is not very practical. Fortunately, for our use case we only need to abstract and concretize sets of integers defined via relational predicates, which is easily tractable, even with infinite bounds. The corresponding semantic functions are defined below.

$$\begin{split} \llbracket x &= a \rrbracket \uparrow_x \doteq \llbracket a, a \rrbracket & \llbracket x \geq a \rrbracket \uparrow_x \doteq \llbracket a, \infty \rrbracket & \llbracket x \leq a \rrbracket \uparrow_x \doteq \llbracket -\infty, a \rrbracket \\ \llbracket x \neq a \rrbracket \uparrow_x \doteq \llbracket -\infty, a - 1 \mid a + 1, \infty \rrbracket & \llbracket x > a \rrbracket \uparrow_x \doteq \llbracket a + 1, \infty \rrbracket & \llbracket x < a \rrbracket \uparrow_x \doteq \llbracket -\infty, a - 1 \rrbracket \end{aligned}$$

5.4 Characters

Abstract characters, i.e., sets of possible characters, are a common component of string grammars whitespace, for example, is usually defined as a set of certain invisible characters. While the size of any string alphabet is always bounded and thus the maximum number of possibilities for a single character is finite, these bounds can be quite large—the Unicode standard currently defines 149 878 characters [Unicode Consortium 2023]. Additionally, it is often desirable to define elements of a string by what characters are *not* allowed to be there. Thus the need for an efficient abstract representation of large sets of (im)possible characters.

Formally, we can define the domain of abstract characters \hat{C} via the alphabet subset lattice $\langle \wp(\Sigma), \subseteq \rangle$, with the usual operations and elements $\langle \hat{C}, \sqsubseteq, \bot, \top, \sqcup, \sqcap, \neg \rangle$ (cf. abstract Booleans \hat{B}). We use the familiar notation Σ interchangeably with \top , indicating the set of all characters of the alphabet. We use set difference to indicate exclusion, e.g., $\Sigma \setminus \{a, b\}$ means the set of all characters excluding 'a' and 'b'. For singleton sets like $\{a\}$ we will usually drop the braces and just write a. Note that $\bot = \emptyset$ is *not* equivalent to the empty string; rather, it is the empty set of characters, representing a space that is impossible to fill or a character that is impossible to produce. The abstract semantics of \hat{C} are defined by

$$\alpha = \gamma = \mathrm{id}, \qquad [\![x = c]\!]\uparrow_x \doteq \{c\}, \qquad [\![x \neq c]\!]\uparrow_x \doteq \Sigma \setminus \{c\}.$$

5.5 Strings

To abstractly represent infinite sets of strings, we define a domain \hat{S} of regular expressions over abstract characters \hat{C} , consisting of the empty language \emptyset , the empty string ε , abstract character literals $\hat{c} \in \hat{C}$ such that $\hat{c} = \{c_1, c_2, ..., c_n\} \equiv c_1 + c_2 + \cdots + c_n$, concatenation $\Box \cdot \Box$ (usually elided), alternation $\Box + \Box$, the Kleene star \Box^* , and optionals $\Box^? \equiv \varepsilon + \Box$. We also allow the abbreviations $\Box^+ = \Box \Box^*$ and $\Box^n = \Box \Box \cdots \Box$ where \Box is repeated *n* times, with n < 1 resulting in \emptyset and $\Box^0 = \varepsilon$.

The domain \hat{S} forms a complete complemented lattice $\langle \hat{S}, \sqsubseteq, \bot, \top, \sqcup, \sqcap, \neg \rangle$ with $\bot = \emptyset, \top = \Sigma^*$, and the standard operations on regular sets [Hopcroft and Ullman 1979]. The language $L(\hat{s})$ describes all strings generated/recognized by a regular expression $\hat{s} \in \hat{S}$, thus $\gamma(\hat{s}) = L(\hat{s})$. We can of course only abstract sets of strings that form a regular language, i.e., $\alpha(S) = \hat{s} \iff S = L(\hat{s})$. While α is not generally realizable [Gold 1967], we can define an abstract semantics over string relations, shown in Figure 9, that allows us to abstract all strings expressed via relations between common string operations.

6 Implementation and Evaluation

We have implemented our approach in the PANINI prototype system, available at https://anonymous. 40pen.science/r/panini. It includes a basic Python frontend with a small default set of λ_{Σ} axioms mimicking the semantics of common Python string functions. Our implementation is written in Haskell and uses the Z3 theorem solver [De Moura and Bjørner 2008] and is modular with respect to

882

876

Fig. 9. Abstract relational semantics for strings. Abstract values are denoted $\hat{\Box}$. We assume single characters have been lifted to singleton character sets, e.g., $s[i] = c \rightsquigarrow s[i] \in \hat{c}$ where $\hat{c} = \{c\}$.

the abstract domains used during grammar inference. PANINI can be used as a library, a standalone batch-mode command-line application, or interactively via a read-eval-print loop.

6.1 Efficient Regular Expressions

An important aspect for the practical viability of our approach is an efficient implementation of the underlying abstract domains, in particular abstract strings \hat{S} and abstract characters \hat{C} . The machine representation of \hat{C} is based on complemented PATRICIA tries [Kmett 2012; Morrison 1968; Okasaki and Gill 1998], which enable compact representation of negated sets while allowing all

1:20

942

943

common set operations. For \hat{S} , we implemented an efficient regular expression type whose literals 932 are abstract characters from \hat{C} and whose operations, like regular intersection and complement, 933 934 are performed purely algebraically, without going through finite automata. Our approach uses Brzozowski derivatives [Antimirov 1996; Brzozowski 1964] and is based on the equation-solving 935 technique by Acay [2018], using Arden's lemma [Arden 1961] and Gaussian elimination to solve 936 a system of regular equations. It is similar to the approach by Liang et al. [2015], but makes use 937 of the local mintermization trick employed by Keil and Thiemann [2014] to effectively compute 938 939 precise derivative over large alphabets. To keep regular expressions as concise and human-readable as possible, we aggressively simplify intermediate expressions using the transformation rules and 940 heuristics described by Kahrs and Runciman [2022]. 941

6.2 Experiments

To provide insights into the efficacy and applicability of our approach, we used the PANINI prototype to infer regular grammars for a number of ad hoc parser implementations.

946 Dataset. We created a benchmark dataset comprising 947 204 regular ad hoc parsers written in Python and then 948 translated to C. The dataset is diverse across two di-949 mensions: complexity of accepted grammar and struc-950 ture of parser code. We generated the dataset by writ-951 ing straightforward parsers for increasingly complex 952 combinations of regular language operations and then 953 added variations of each parser, e.g., using loops to 954 iterate over the characters in a string, using high-level 955 Python string functions such as index, parsing the 956 input from back-to-front, and so on. Figures 3 and 10 957 to 13 show subjects from the benchmark dataset. 958

To facilitate a structured analysis, we classified the 959 parsers within the dataset into three overarching cate-960 gories based on their structural features: Straight-Line 961 Programs are ad hoc parsers characterized by linear 962 execution flow without branching constructs such as 963 conditionals or loops (e.g., Figure 10); Programs with 964 Conditionals are ad hoc parsers that incorporate con-965 ditional statements to make decisions based on input 966

```
def getAddrSpec(email):

b1 = email.index('<',\emptyset)+1

b2 = email.index('>',b1)

return email[b1:b2]

getAddrSpec : {email : S | *} \rightarrow S

getAddrSpec = \lambda(email : S).

let v_0 = indexFrom email "<" 0 in

let v_1 = ge v_0 0 in

let v_1 = ge v_0 0 in

let b_1 = add v_0 1 in

let b_2 = indexFrom email ">" b_1 in

let v_3 = ge b_2 0 in

let _ = assert v_3 in

slice email b_1 b_2
```

Fig. 10. A straight-line ad hoc parser.

characteristics (e.g., Figure 3); and *Programs with Loops* are ad hoc parsers containing iterative constructs alongside conditional statements for string manipulation tasks (e.g., Figure 11).

Methodology. For each ad hoc parser in the dataset, we used PANINI to transpile the original Python 969 source to λ_{Σ} and automatically infer a grammar from the parser's λ_{Σ} representation. We compared 970 each inferred grammar G_i against a manually derived ground truth grammar \hat{G}_i . We differentiate 971 between exact matches, where $L(G_i) = L(\hat{G}_i)$; successful under-approximations $L(G_i) \subset L(\hat{G}_i)$, 972 which correctly identify a subset of allowed strings; trivial under-approximations to the empty 973 language $L(G_i) = \emptyset$; and cases where PANINI reports an error due to limitations of the prototype 974 implementation or the underlying abstract domains. We additionally computed the precision and 975 *recall* of each inferred grammar. Precision measures the percentage of inputs accepted by G_i that 976 are also accepted by \hat{G}_i , whereas recall measures the pecentage of inputs accepted by \hat{G}_i that are 977 also accepted by G_i . We generate up to 1000 random sample inputs from the respective grammars 978 to compute each measure. Low precision indicates over-approximation, i.e., the inferred grammar 979

Table 3. Results of evaluating PANINI on a varied dataset of ad hoc parsers.									
			Infer	red	Lan	guage			
Parser Example #		#	=	\subset	Ø	Error	SR	Р	R
getAddrSpec	[^<>]*<[^>]*>.*	89	84	0	0	5	.94	1.00	1.00
Figure 3	a [^a]b.*	51	50	0	0	1	.98	1.00	1.00
lsb check	0*1	64	40	7	7	10	.84	1.00	.76

Table

SR = success rate, P = precision, R = recall

981 982 983

984

985

986

987

988

989 990

991

992

993

994

995

996

997

998

999

1000

1001

1002

1003

1004

1005

1006

1007

1008

1009

1010

1011

1012

1013

1014

1015

1016

1017

1018

1019

1020

1021

1022

1023

1024

1025

1026

1027

1028 1029 Type

Straight

+ Cond.

+ Loops

allows more inputs than would be safely accepted by the parser program, while low recall indicates under-approximation, i.e., the inferred grammar is unnecessarily stricter than the actual program.

174 7 7

204

We aggregated precision and recall across each of the three parser categories and additionally report the success rate, the percentage of benchmark subjects for which PANINI is able to infer a grammar. Finally, we also report average wall-clock execution time. All benchmarks were run on a MacBook Pro with an Apple M4 processor and 24 GB RAM, using Z3 4.8.10 for SMT solving.

```
def lsb_check(s):
   i = 0
   while i < len(s)-1:</pre>
       assert s[i] == '0'
       i += 1
   assert s[i] == '1'
lsb_check : {s : \mathbb{S} \mid \star} \rightarrow \mathbb{1}
lsb_check = \lambda(s : S).
   \operatorname{rec} L_2: \{i: \mathbb{Z} \mid \star\} \to \mathbb{1} = \lambda i.
      let v_0 = length s in
      let v_1 = \operatorname{sub} v_0 \ 1 in
      let v_2 = \operatorname{lt} i v_1 in
      if v_2 then
         let v_3 = \text{slice1} s i in
          let v_4 = match v_3 "0" in
          let _ = assert v_4 in
          let i = add i 1 in
          L_2 i
      else
          let v_5 = \text{slice1 } s i in
          let v_6 = match v_5 "1" in
          assert v_6
   in L_2 0
```

Fig. 11. An ad hoc parser with loops.

Results. Table 3 presents our experimental results. Out of the 204 parser programs in our dataset, PANINI is able to successfully infer a grammar for 188, a success rate of 92 %. As expected, the precision of inferred grammars is 100 % across all types of parsers: PANINI never over-approximates. The average overall recall is 93 %. PANINI achieves exact inference for all straight-line and purely conditional programs and for 40 programs containing loops; it safely under-approximates the grammar of 7 loop programs; and it fails to find a meaningful refinement beyond the empty language for another 7 loop programs.

.92

1.00

.93

16

The performance of grammar inference is mostly in the sub-second range, with some loop programs as outliers. Most of the inference time is spent during classical refinement inference, where SMT solving is still the biggest bottleneck.

6.3 Comparison with Other Approaches

Table 4 presents a comparison of PANINI and other grammar inference systems. We compare operating requirements (whether the parser's source code needs to be available; whether the parser needs to be executed in some way, perhaps instrumented or modified; and whether the approach requires pre-existing input samples), the type of output (a regular expression, a context-free grammar, an automaton), and the results of running the approach on the PANINI benchmark dataset (see § 6.2). The systems and algorithms selected for comparison represent the state-ofthe-art in grammar inference:

Exbar [Lang 1999] is the fastest known algorithm for finding minimal DFAs from labeled samples only. This type of grammar inference-known as *passive automata learning*-is notable in that it does not require the existence of a parser program at all, neither to run nor inspect.

Time (s)

 0.16 ± 0.27

 0.09 ± 0.05

 2.32 ± 4.77

 0.82 ± 2.85

1030 1031 PANINI benchmark Requirements 1032 Р Approach Source Execution Samples Output SR R Time (s) 1033 PANINI Python or λ_{Σ} Regex .92 1.00 .93 0.82 ± 2.85 1034 CFG .95 STALAGMITE С symbolic .40 .56 75.07 ±249.25 1035 .97 Python or C traced CFG Mimid positive .67 .81 27.39 ± 42.78 1036 TREEVADA oracle positive CFG .98 .99 .66 72.63 ±100.84 1037 _ TTT .99 _ oracle positive DFA .13 1.00 2.74 ± 15.17 1038 Exbar pos. + neg. DFA _ _ _ 1039 1040

Table 4. Comparison of state-of-the-art grammar inference approaches.

Other prominent algorithms in this category include RPNI [Oncina and García 1992] and the 1041 approximative ED-BEAM [Lang 1999]. Unfortunately, the requirement of a well-labeled set of representative input samples is unrealistic in many settings, including ours.

TTT [Isberner 2015; Isberner et al. 2014] is a leading algorithm in active automata learning, 1044 specifically within the minimally adequate teacher (MAT) framework. Established by Angluin 1045 [1987] with the introduction of the seminal L^* algorithm, MAT formulates grammar learning as 1046 an interactive process, in which a teacher-an oracle or black-box parser program-can answer 1047 two types of question: whether a certain input is a *member* of the target language, and whether a 1048 hypothesized language is *equivalent* to the target language, providing a counterexample if it is not. 1049 In practice, the requirement of equivalence queries is quite onerous, which is why they are usually 1050 approximated by conformance testing [Aichernig et al. 2024] using known positive input samples. 1051

When running TTT on the PANINI benchmark, we generated up to 10 000 positive input samples 1052 for each parser, simulating a best-case scenario. The results (13 % average precision, 100 % recall) 1053 indicate that TTT tends to significantly over-approximate the true input language for these kinds 1054 of ad hoc parsers, even with a very large number of input samples. 1055

TREEVADA [Arefin et al. 2024] is the state-of-the-art in black-box inference of context-free 1056 grammars. Other approaches in this vein are ARVADA [Kulkarni et al. 2021] and the pioneering 1057 GLADE [Bastani et al. 2017]. These systems do not require equivalence queries, which makes them 1058 much more practical than traditional MAT algorithms, but they do all require a well-covering set 1059 of positive input samples in order to produce accurate grammars [Bendrissou et al. 2022]. 1060

When given up to 20 positive input samples per parser in the PANINI benchmark, TREEVADA 1061 achieves near-perfect 99 % precision (on average) but under-approximates the true grammars with 1062 only 66 % average recall. We found that doubling the amount of input samples improves recall by 1063 less than 10 % while more than doubling the runtime. 1064

Mimid [Gopinath et al. 2020a] generalizes positive sample inputs into a context-free grammar 1065 by analyzing execution traces of an instrumented version of the parser, which needs to be available 1066 in source form. Mimid significantly improves on the previous white-box approach AUTOGRAM 1067 [Höschele and Zeller 2017], but it still requires a good set of pre-existing input samples to produce 1068 an accurate grammar. A general advantage of white-box approaches is that the inferred grammars 1069 tend to be very readable, because they can incorporate identifiers from the source code. 1070

On the PANINI benchmark, with up to 1000 positive sample inputs for each parser, Mimid produces 1071 grammars of reasonably high precision (97 % on average) but tends to under-approximate (81 % 1072 average recall). It also has a low success rate of only 67 %, meaning Mimid failed to infer any 1073 grammar for a third of the parsers in the dataset. 1074

STALAGMITE [Bettscheider and Zeller 2024] is a recent white-box approach that obviates the 1075 need for input samples by transforming the source program into a version amenable for symbolic 1076 execution. In addition to inserting tracing calls, this includes limiting recursion depth and the 1077

1042 1043

number of loop iterations, which enables a symbolic test generator like KLEE [Cadar et al. 2008a]
to automatically generate input samples that cover all execution paths. After running the modified

parser on a large enough number of samples, the collected symbolic input traces are woven togetherto produce a context-free grammar.

STALAGMITE performs rather poorly on the PANINI benchmark, with both low precision (40 % on
 average) and low recall (56 % on average). While it does not require any pre-existing input samples,
 the symbolic execution of the parser programs is quite resource intensive.

Note that **PANINI** is the only approach that requires neither positive input samples nor any interaction with the parser. It is therefore uniquely suited to deal with ad hoc parsers, for which input samples are generally not available and which occur as code fragments that cannot always be expected to run as-is. In its current form, PANINI is limited to inference of at-most regular languages, but we conjecture that these make up the highest share of ad hoc parsers in the wild [Schröder et al. 2023]. We argue that all other existing approaches have requirements that make them impractical to use in this setting, including prohibitive resource usage.

1094 6.4 Real-World Case Studies

1093

To investigate PANINI's suitability for real-world grammar inference, we look at two regularlanguage ad hoc parsers from the *Mimid* benchmark suite [Gopinath et al. 2020b].

1097 cgidecode.py. This is a Python program to decode 1098 CGI-encoded strings, an encoding where spaces are re-1099 placed by plus signs and other invalid characters are re-1100 placed by a hexadecimal encoding prefixed with a per-1101 cent sign. If an improperly encoded string is encoun-1102 tered, the program raises an error. The regular expres-1103 sion ([^%] |%[0-9A-Fa-f][0-9A-Fa-f])* encompasses 1104 the input language of this ad hoc parser.³ As Table 5 shows, 1105 all approaches except for STALAGMITE do fairly well on

Table 5. cgidecode.py benchmark					
	Р	R	Time (s)		
Panini	1.00	1.00	1.49		
STALAGMITE	.12	1.00	26.74		
Mimid	.90	.96	139.07		
TreeVada	.96	.60	288.40		
TTT	1.00	1.00	3.60		

precision. TREEVADA exhibits the lowest recall, likely due to prioritizing larger-scale inference
 of context-free properties, to the detriment of regular language accuracy. Only PANINI and TTT
 achieve both perfect precision and recall, and only PANINI does so without having access to any
 known positive inputs.

	Р	R	Time (s)
Panini	-	-	
STALAGMITE	-	-	
Mimid	1.00	0.19	157.0
TreeVada	.00	0.96	625.3
TTT	.00	1.00	1.7

urlparse.py. This is the URL parser part of the Python *urllib* library. Even though this parser is quite large and complex (>1000 LoC), it is written very defensively. The only time it actually rejects an input is when the network location part of an otherwise valid URL contains an insufficiently bracketed IPv6 address. The parser's input language is equivalent to the regular expression (([a-zA-Z0-9.+-]+:)?//([^/?#]*(\[[^/?#]*[]]]][]][^/?#]*\[).*|[^][]*([/?#].*)?))|(.*([^:]]

|[^a-zA-Z0-9.+-].*)//.*), which looks somewhat baroque but is surprisingly permissive.⁴ This
 proves to be a difficult case for grammar inference: As shown in Table 6, TREEVADA and TTT
 over-approximate the correct grammar to the extent that they achieve no precision; *Mimid* achieves

1127

1

1

1

1

1

1119

³The golden grammar included by the *Mimid* artifact is insufficient: it does not include all possible two-digit hexadecimals and is limited to a subset of the ASCII alphabet.

⁴This is again significantly different from the golden grammar included by the *Mimid* artifact, which is much more restrictive, yet at the same time does not prevent the parser's actual error case.

1128 perfect precision at the cost of severely under-approximating the true grammar; and neither 1129 STALAGMITE nor PANINI are able to infer any grammar at all. STALAGMITE requires C source 1130 code, which does not exist for this program, while PANINI's Python frontend is not yet capable of 1131 automatically translating this syntactically complex parser into λ_{Σ} .

1133 6.5 Current Limitations and Future Work

Our PANINI prototype is a proof-of-concept that shows the viability of our approach; it is not yet a
practical end-user system. Based on the results of our evaluation, including the two case studies,
we see the following main areas of improvement as part of future work:

¹¹³⁷ **Parser Extraction.** PANINI is built around λ_{Σ} , a calculus for representing ad hoc parsers. While ¹¹³⁸ the present work focuses on synthesizing grammars from λ_{Σ} programs, a significant part of our ¹¹³⁹ envisioned grammar inference process is the extraction of λ_{Σ} from general-purpose programming ¹¹⁴⁰ languages (cf. § 2.1). The PANINI prototype does include some machinery to automatically extract ¹¹⁴¹ and transpile ad hoc parser slices from Python to λ_{Σ} , but this is still very limited. We are currently ¹¹⁴² experimenting with different parser slicing techniques and continue to improve our Python front-¹¹⁴³ end. We are also hoping to add prototypical support for other languages in the near future.

Language Features. The λ_{Σ} language is by design minimal, in order to provide a common intermediate representation for a wide range of ad hoc parser implementations and to simplify many aspects of refinement and grammar inference. However, its lack of language features makes it difficult to easily capture many real-world parser programs, such as those involving generic data types. We are currently working on extending the λ_{Σ} language to add support for polymorphism and user-defined data types.

- **Invariant Inference.** In our approach, the solving of non-grammar κ variables is delegated to the classical refinement inference machinery (§ 3.3). In particular, this includes the generation of loop/recursion invariants, which our prototype infers using a textbook implementation of purely conjunctive predicate abstraction, a technique that is inherently limited in the shape of invariants that can be produced and is highly dependent on a good set of candidate predicates; it is also a runtime bottleneck. We plan to improve our qualifier extraction heuristics and integrate external state-of-the-art invariant generators.
- 1158 Abstract Domains. Abstract interpretation is bounded by the
- limits of the underlying abstract domains (§ 5). For example,
 our integer domain cannot efficiently represent congruence
 classes, e.g., infinite sets of even or odd numbers. This can
 lead to under-approximations, as in Figure 12, where PANINI
 can only infer the subset (ab)? of the true grammar (ab)*.
 By design, our system is modular in the choice of underlying
 abstract domains, and we are working on extending them.
- Relational Semantics. If PANINI lacks the semantics to 1166 rewrite, normalize, or abstract a particular relation (see § 4.3), 1167 it will eventually become stuck. We can increase the capabil-1168 ities of the system by extending the set of semantic rules-i.e., 1169 adding more equations to Figure 9. But take the parser in 1170 Figure 13, which leads to PANINI trying to compute the ab-1171 straction $[s[0] = s[1]]\uparrow_s$. In isolation, this constraint is not 1172 even expressible in a regular string domain-even though the 1173 final grammar is regular. To handle such tricky cases, a more 1174 complex non-local rewriting strategy might be needed. 1175

```
def f250(s):
    i = 0
    while i < len(s):
        assert s[i] == "a"
        assert s[i+1] == "b"
        i += 2
```



```
def f521(s):
    assert s[0] == "a"
    assert s[1] == s[0]
```

Fig. 13. A parser for aa.*

1132

Static Inference of Regular Grammars for Ad Hoc Parsers

1177 Context-Free Grammars. If Panini encounters a context-free parser, it will likely yield Ø or get
1178 stuck, since the underlying string domain is built on regular expressions and cannot represent
1179 context-free constructs. Non-trivial recursion exhibited by a parser might impede invariant infer1180 ence, however we do not see any limitations inherent to our technique that would in principle
1181 prevent us from eventually inferring (deterministic) context-free grammars.

¹¹⁸³ 7 Related Work

1182

1184

String Constraint Solving. Precise formal reasoning over strings can be accomplished using string 1185 constraint solving (SCS), a declarative paradigm of modeling relations between string variables and 1186 solving attendant combinatorial problems [Amadini 2021]. It is usually assumed that collecting 1187 string constraints requires some kind of (dynamic) symbolic execution [Kausler and Sherman 2014], 1188 1189 and practical SCS applications are generally concerned with the inverse of our problem: modeling the possible strings a function can return or express [Bultan et al. 2018], instead of the strings a 1190 function can accept. In our approach, we use purely static means (viz. refinement type inference) 1191 to essentially collect input string constraints (see § 3), which we then simplify/solve in ways not 1192 dissimilar but nonetheless different from traditional SCS techniques (see § 4). There have been 1193 1194 many recent advances in SCS for SMT [Abdulla et al. 2015; Kan et al. 2022; Kiezun et al. 2009; Trinh et al. 2014, 2020; Zheng et al. 2013]. We particularly make use of string theories embedded in the 1195 Z3 constraint solver [Berzish et al. 2017] in our implementation (§6). 1196

1197 Related work in (dynamic) symbolic execution make use of constraint solvers over string domains at their core to reason about how strings are manipulated in programs, with applications ranging 1198 1199 from generating test inputs [Bjørner et al. 2009; Cadar et al. 2008b; Li et al. 2011] and detecting vulnerabilities (e.g., cross-site scripting, SQL injection) [Holik et al. 2017; Loring et al. 2017; Saxena 1200 et al. 2010]. These approaches differ from our work on two specific aspects. First, they rely on 1201 traces from dynamic executions to infer more precise constraints, while we are able to reason about 1202 1203 string constraints statically. Second, in the case of detecting vulnerabilities, they reason about the 1204 resulting output induced through string operations, and do not infer a grammar over the input 1205 language, which is our main goal.

Abstract Domains. Abstract string domains approximate strings to track information precisely
 enough to analyze particular behaviors of interest while only preserving relevant information.
 Most of the existing work in string domains differs in what kind of behavior is of interest and how
 the approximation is achieved efficiently.

Costantini et al. [2015] introduces a suite of different abstract semantics for concatenation, charac-1211 ter inclusion, and substring extraction (particularly pre- and suffixes). In their work, they explicitly 1212 discuss the trade-off between precision and efficiency. Amadini et al. [2020] review the dashed 1213 string abstraction, an approach that considers strings as blocks of characters and the constraints 1214 on these blocks, which has shown good performance on benchmarks involving constraints on 1215 string length, equality, concatenation, and regular expression membership. M-String [Cortesi and 1216 Olliaro 2018] considers a parametric abstract domain for strings in the C programming language 1217 by leveraging abstract domains for the content of a string and for expressions to infer when a 1218 string index position corresponds to an expression of interest. While most of the existing work has 1219 focused on approximating a single variable, very recent work by Arceri et al. [2022] focuses on 1220 relational string domains that try to capture the relation between string variables and expressions 1221 for which we cannot compute static values, such as user input. 1222

There have been a multitude of abstract domains that aim at specific target languages, such as JavaScript [Amadini et al. 2017; Jensen et al. 2009; Kashyap et al. 2014; Park et al. 2016].

1223

¹²²⁴ 1225

Anon.

1:26

1230

Precondition Inference. Despite a wide variety of approaches for computing preconditions [Barnett
and Leino 2005; Cousot et al. 2013; Dillig et al. 2013; Padhi et al. 2016; Seghir and Kroening 2013],
we are not aware of any that focus specifically on string operations, or that would allow us to
reconstruct an input grammar in a way suitable for our envisioned applications.

1231 Data-Availability Statement

We have implemented our approach in the PANINI prototype system, whose source repository will become publicly available upon acceptance and is currently viewable for reviewers in anonymized form at https://anonymous.4open.science/r/panini. This repository also includes our full evaluation dataset. We will additionally provide a self-contained artifact (most likely in the form of a Docker container) to easily reproduce all claims made in this paper.

1238 References

- Parosh Aziz Abdulla, Mohamed Faouzi Atig, Yu-Fang Chen, Lukáš Holík, Ahmed Rezine, Philipp Rümmer, and Jari Stenman.
 2015. Norn: An SMT solver for string constraints. In *Computer Aided Verification: 27th International Conference, CAV* 2015. San Francisco, CA, USA, July 18-24, 2015, Proceedings, Part I. Springer, 462–469.
- Josh Acay. 2018. A Regular Expression Library for Haskell. (2018). https://github.com/cacay/regexp Unpublished manuscript,
 dated May 22, 2018. LaTeX files and Haskell source code.
- Bernhard K. Aichernig, Martin Tappler, and Felix Wallner. 2024. Benchmarking Combinations of Learning and Testing Algorithms for Automata Learning. *Form. Asp. Comput.* 36, 1, Article 3 (mar 2024), 37 pages. https://doi.org/10.1145/ 3605360
- James F. Allen. 1983. Maintaining Knowledge about Temporal Intervals. *Commun. ACM* 26, 11 (nov 1983), 832–843.
 https://doi.org/10.1145/182.358434
- 1248 Roberto Amadini. 2021. A Survey on String Constraint Solving. arXiv:2002.02376 [cs.AI]
- Roberto Amadini, Graeme Gange, and Peter J. Stuckey. 2020. Dashed strings for string constraint solving. Artificial Intelligence 289 (2020), 103368. https://doi.org/10.1016/j.artint.2020.103368
- Roberto Amadini, Alexander Jordan, Graeme Gange, François Gauthier, Peter Schachte, Harald SØndergaard, Peter J. Stuckey,
 and Chenyi Zhang. 2017. Combining String Abstract Domains for JavaScript Analysis: An Evaluation. In Proceedings,
 Part I, of the 23rd International Conference on Tools and Algorithms for the Construction and Analysis of Systems Volume
 10205. Springer-Verlag, Berlin, Heidelberg, 41–57. https://doi.org/10.1007/978-3-662-54577-5_3
- Dana Angluin. 1987. Learning regular sets from queries and counterexamples. Information and Computation 75, 2 (1987), 87–106. https://doi.org/10.1016/0890-5401(87)90052-6
- Valentin Antimirov. 1996. Partial derivatives of regular expressions and finite automaton constructions. *Theoretical Computer Science* 155, 2 (March 1996), 291–319. https://doi.org/10.1016/0304-3975(95)00182-4
- 1257 Vincenzo Arceri, Martina Olliaro, Agostino Cortesi, and Pietro Ferrara. 2022. Relational String Abstract Domains. In
 1258 Verification, Model Checking, and Abstract Interpretation: 23rd International Conference, VMCAI 2022, Philadelphia, PA,
 1259 USA, January 16–18, 2022, Proceedings (Philadelphia, PA, USA). Springer-Verlag, Berlin, Heidelberg, 20–42. https://doi.org/10.1007/978-3-030-94583-1_2
- Dean N. Arden. 1961. Delayed-logic and finite-state machines. In 2nd Annual Symposium on Switching Circuit Theory and Logical Design (SWCT 1961). 133–151. https://doi.org/10.1109/FOCS.1961.13
- Mohammad Rifat Arefin, Suraj Shetiya, Zili Wang, and Christoph Csallner. 2024. Fast Deterministic Black-box Context-free
 Grammar Inference. In *Proceedings of the IEEE/ACM 46th International Conference on Software Engineering* (Lisbon,
 Portugal) (*ICSE '24*). Association for Computing Machinery, New York, NY, USA, Article 117, 12 pages. https://doi.org/
 10.1145/3597503.3639214
- Mike Barnett and K. Rustan M. Leino. 2005. Weakest-Precondition of Unstructured Programs. In *Proceedings of the 6th ACM SIGPLAN-SIGSOFT Workshop on Program Analysis for Software Tools and Engineering* (Lisbon, Portugal) (*PASTE '05*).
 ACM, New York, NY, USA, 82–87. https://doi.org/10.1145/1108792.1108813
- 1268
 Clark Barrett, Pascal Fontaine, and Cesare Tinelli. 2016. The Satisfiability Modulo Theories Library (SMT-LIB). http://smt-lib.org

 1269
 lib.org
- Clark Barrett, Roberto Sebastiani, Sanjit A. Seshia, and Cesare Tinelli. 2021. Satisfiability Modulo Theories. In *Handbook of Satisfiability* (2nd ed.). IOS Press, Chapter 33, 1267–1329.
- Osbert Bastani, Rahul Sharma, Alex Aiken, and Percy Liang. 2017. Synthesizing Program Input Grammars. In Proceedings of the 38th ACM SIGPLAN Conference on Programming Language Design and Implementation (Barcelona, Spain) (PLDI 2017).
 ACM, New York, NY, USA, 95–110. https://doi.org/10.1145/3062341.3062349
- 1274

Proc. ACM Program. Lang., Vol. 1, No. OOPSLA, Article 1. Publication date: January 202025-03-25 18:26. Page 26 of 1-30.

- Bachir Bendrissou, Rahul Gopinath, and Andreas Zeller. 2022. "Synthesizing input grammars": a replication study. In
 Proceedings of the 43rd ACM SIGPLAN International Conference on Programming Language Design and Implementation
 (San Diego, CA, USA) (PLDI 2022). Association for Computing Machinery, New York, NY, USA, 260–268. https://doi.org/10.1145/3519939.3523716
- Murphy Berzish, Vijay Ganesh, and Yunhui Zheng. 2017. Z3str3: A string solver with theory-aware heuristics. In 2017
 Formal Methods in Computer Aided Design (Vienna, Austria) (*FMCAD 2017*). IEEE, 55–59. https://doi.org/10.23919/
 FMCAD.2017.8102241
- Leon Bettscheider and Andreas Zeller. 2024. Look Ma, No Input Samples! Mining Input Grammars from Code with Symbolic
 Parsing. In Companion Proceedings of the 32nd ACM International Conference on the Foundations of Software Engineering
 (Porto de Galinhas, Brazil) (FSE 2024). 522–526. https://doi.org/10.1145/3663529.3663790
- Saroja Bhate and Subhash Kak. 1991. Pāṇini's Grammar and Computer Science. Annals of the Bhandarkar Oriental Research
 Institute 72/73, 1/4 (1991), 79–94.
- Nikolaj Bjørner, Arie Gurfinkel, Ken McMillan, and Andrey Rybalchenko. 2015. Horn Clause Solvers for Program Verification.
 In Fields of Logic and Computation II: Essays Dedicated to Yuri Gurevich on the Occasion of His 75th Birthday. Springer, 24–51.
- Nikolaj Bjørner, Nikolai Tillmann, and Andrei Voronkov. 2009. Path feasibility analysis for string-manipulating programs. In Tools and Algorithms for the Construction and Analysis of Systems: 15th International Conference, TACAS 2009, Held as Part of the Joint European Conferences on Theory and Practice of Software, ETAPS 2009, York, UK, March 22-29, 2009.
 Proceedings 15. Springer, 307–321.
- 1291 William J. Bowman. 2022. The A Means A. https://www.williamjbowman.com/blog/2022/06/30/the-a-means-a/
- Matthias Braun, Sebastian Buchwald, Sebastian Hack, Roland Leißa, Christoph Mallon, and Andreas Zwinkau. 2013. Simple and Efficient Construction of Static Single Assignment Form. In *Proceedings of the 22nd International Conference on Compiler Construction* (Rome, Italy) (*CC'13*). Springer-Verlag, Berlin, Heidelberg, 102–122. https://doi.org/10.1007/978-3-642-37051-9_6
- Janusz A. Brzozowski. 1964. Derivatives of Regular Expressions. J. ACM 11, 4 (Oct. 1964), 481–494. https://doi.org/10.1145/
 321239.321249
- Tevfik Bultan, Fang Yu, Muath Alkhalaf, and Abdulbaki Aydin. 2018. String Analysis for Software Verification and Security (1st ed.). Springer Cham. https://doi.org/10.1007/978-3-319-68670-7
- Cristian Cadar, Daniel Dunbar, and Dawson Engler. 2008a. KLEE: unassisted and automatic generation of high-coverage
 tests for complex systems programs. In *Proceedings of the 8th USENIX Conference on Operating Systems Design and Implementation* (San Diego, California) (OSDI'08). USENIX Association, USA, 209–224.
- Cristian Cadar, Vijay Ganesh, Peter M Pawlowski, David L Dill, and Dawson R Engler. 2008b. EXE: Automatically generating
 inputs of death. ACM Transactions on Information and System Security (TISSEC) 12, 2 (2008), 1–38.
- Manuel MT Chakravarty, Gabriele Keller, and Patryk Zadarnowski. 2004. A functional perspective on SSA optimisation algorithms. *Electronic Notes in Theoretical Computer Science* 82, 2 (2004), 347–361. https://doi.org/10.1016/S1571-0661(05)82596-4
- N. Chomsky and M.P. Schützenberger. 1963. The Algebraic Theory of Context-Free Languages. In *Computer Programming* and *Formal Systems*, P. Braffort and D. Hirschberg (Eds.). Studies in Logic and the Foundations of Mathematics, Vol. 35. Elsevier, 118–161. https://doi.org/10.1016/S0049-237X(08)72023-8
- Agostino Cortesi and Martina Olliaro. 2018. M-String Segmentation: A Refined Abstract Domain for String Analysis
 in C Programs. In 2018 International Symposium on Theoretical Aspects of Software Engineering (TASE). 1–8. https://doi.org/10.1109/TASE.2018.00009
- Benjamin Cosman and Ranjit Jhala. 2017. Local Refinement Typing. *PACM on Programming Languages* 1, ICFP, Article 26 (Aug. 2017), 27 pages. https://doi.org/10.1145/3110270
- Giulia Costantini, Pietro Ferrara, and Agostino Cortesi. 2015. A Suite of Abstract Domains for Static Analysis of String Values. *Softw. Pract. Exper.* 45, 2 (feb 2015), 245–287. https://doi.org/10.1002/spe.2218
 Internet and Agostino Cortesi. 2015. A Suite of Abstract Domains for Static Analysis of String Values. *Softw. Pract. Exper.* 45, 2 (feb 2015), 245–287. https://doi.org/10.1002/spe.2218
- Patrick Cousot. 1997. Types as Abstract Interpretations. In *Proceedings of the 24th ACM SIGPLAN-SIGACT Symposium on Principles of Programming Languages* (Paris, France) (*POPL '97*). Association for Computing Machinery, New York, NY,
 USA, 316–331. https://doi.org/10.1145/263699.263744
- Patrick Cousot and Radhia Cousot. 1977. 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* (Los Angeles, California) (*POPL '77*). Association for Computing Machinery, New York, NY, USA, 238–252. https://doi.org/10.1145/512950.512973
- Patrick Cousot and Radhia Cousot. 1979. Systematic Design of Program Analysis Frameworks. In *Proceedings of the 6th ACM SIGACT-SIGPLAN Symposium on Principles of Programming Languages* (San Antonio, Texas) (*POPL '79*). Association
 for Computing Machinery, New York, NY, USA, 269–282. https://doi.org/10.1145/567752.567778
- 1322 1323

1:28

1324

1325

(Rome, Italy) (VMCAI 2013). Springer-Verlag, Berlin, Heidelberg, 128-148. https://doi.org/10.1007/978-3-642-35873-9_10 1326 Patrick Cousot and Nicolas Halbwachs. 1978. Automatic discovery of linear restraints among variables of a program. In 1327 Proceedings of the 5th ACM SIGACT-SIGPLAN symposium on Principles of programming languages - POPL '78 (POPL '78). 1328 ACM Press. https://doi.org/10.1145/512760.512770 Luis Damas and Robin Milner. 1982. Principal Type-Schemes for Functional Programs. In Proceedings of the 9th ACM 1329 SIGPLAN-SIGACT Symposium on Principles of Programming Languages (Albuquerque, New Mexico) (POPL '82). Association 1330 for Computing Machinery, New York, NY, USA, 207-212. https://doi.org/10.1145/582153.582176 1331 Leonardo De Moura and Nikolaj Bjørner. 2008. Z3: An Efficient SMT Solver. In Proceedings of the Theory and Practice of 1332 Software, 14th International Conference on Tools and Algorithms for the Construction and Analysis of Systems (Budapest, 1333 Hungary) (TACAS'08/ETAPS'08). Springer-Verlag, Berlin, Heidelberg, 337-340. https://doi.org/10.1007/978-3-540-78800-1334 3 24 Isil Dillig, Thomas Dillig, Boyang Li, and Ken McMillan. 2013. Inductive Invariant Generation via Abductive Inference. In 1335 Proceedings of the 2013 ACM SIGPLAN International Conference on Object Oriented Programming Systems Languages & 1336

Patrick Cousot, Radhia Cousot, Manuel Fähndrich, and Francesco Logozzo. 2013. Automatic Inference of Necessary Preconditions. In Proceedings of the 14th International Conference on Verification, Model Checking, and Abstract Interpretation

- Proceedings of the 2013 ACM SIGPLAN International Conference on Object Oriented Programming Systems Languages & Applications (Indianapolis, Indiana, USA) (OOPSLA '13). ACM, New York, NY, USA, 443–456. https://doi.org/10.1145/2509136.2509511
- Jana Dunfield and Neel Krishnaswami. 2021. Bidirectional Typing. Comput. Surveys 54, 5, Article 98 (May 2021), 38 pages.
 https://doi.org/10.1145/3450952
- Aryaz Eghbali and Michael Pradel. 2020. No strings attached: An empirical study of string-related software bugs. In
 Proceedings of the 35th IEEE/ACM International Conference on Automated Software Engineering. 956–967.
- M Fitter and TRG Green. 1979. When do diagrams make good computer languages? *International Journal of man-machine studies* 11, 2 (1979), 235–261.
- Cormac Flanagan, Amr Sabry, Bruce F. Duba, and Matthias Felleisen. 1993. The Essence of Compiling with Continuations.
 In Proceedings of the ACM SIGPLAN 1993 Conference on Programming Language Design and Implementation (Albuquerque, New Mexico, USA) (PLDI '93). ACM, New York, NY, USA, 237–247. https://doi.org/10.1145/155090.155113
- Roberto Giacobazzi and Elisa Quintarelli. 2001. Incompleteness, Counterexamples, and Refinements in Abstract Model Checking. In *Proceedings of the 8th International Symposium on Static Analysis (SAS '01)*. Springer-Verlag, Berlin, Heidelberg, 356–373.
- David J. Gilmore and Thomas R. G. Green. 1984. Comprehension and recall of miniature programs. *International Journal of Man-Machine Studies* 21, 1 (1984), 31–48.
- 1350 E Mark Gold. 1967. Language identification in the limit. Information and control 10, 5 (1967), 447–474.
- Rahul Gopinath, Björn Mathis, and Andreas Zeller. 2020a. Mining Input Grammars from Dynamic Control Flow. In *Proceedings* of the 28th ACM Joint Meeting on European Software Engineering Conference and Symposium on the Foundations of Software Engineering (Virtual Event, USA) (ESEC/FSE 2020). ACM, New York, NY, USA, 172–183. https://doi.org/10.1145/3368089. 3409679
- Rahul Gopinath, Björn Mathis, and Andreas Zeller. 2020b. Replication package for Mining Input Grammars from Dynamic
 Control Flow. https://doi.org/10.5281/zenodo.3876969
- R. Hindley. 1969. The Principal Type-Scheme of an Object in Combinatory Logic. *Trans. Amer. Math. Soc.* 146 (December 1969), 29–60.
 1357
 146 (December 1969), 29–60.
- Lukáš Holík, Petr Janků, Anthony W. Lin, Philipp Rümmer, and Tomáš Vojnar. 2017. String Constraints with Concatenation and Transducers Solved Efficiently. *Proc. ACM Program. Lang.* 2, POPL, Article 4 (dec 2017), 32 pages. https://doi.org/10.
 1145/3158092
- 1360 John Hopcroft and Jeffrey Ullman. 1979. Introduction to Automata Theory, Languages, and Computation. Addison-Wesley.
- Matthias Höschele and Andreas Zeller. 2017. Mining Input Grammars with AUTOGRAM. In 2017 IEEE/ACM 39th International Conference on Software Engineering Companion (ICSE-C). 31–34. https://doi.org/10.1109/ICSE-C.2017.14
- Malte Isberner. 2015. Foundations of Active Automate Learning: An Algorithmic Perspective. Ph. D. Dissertation. TU Dortmund.
 Malte Isberner, Falk Howar, and Bernhard Steffen. 2014. The TTT Algorithm: A Redundancy-Free Approach to Active
- Automata Learning. In *Runtime Verification*, Borzoo Bonakdarpour and Scott A. Smolka (Eds.). Springer International
 Publishing, Cham, 307–322.
- Simon Holm Jensen, Anders Møller, and Peter Thiemann. 2009. Type analysis for javascript.. In SAS, Vol. 9. Springer, 238–255.
- Ranjit Jhala and Niki Vazou. 2020. Refinement Types: A Tutorial. (2020). arXiv:2010.07763 [cs.PL]
- 1368 Stephen C Johnson and Ravi Sethi. 1990. Yacc: A Parser Generator. UNIX Vol. II: Research System (1990), 347–374.
- Stefan Kahrs and Colin Runciman. 2022. Simplifying regular expressions further. *Journal of Symbolic Computation* 109 (March 2022), 124–143. https://doi.org/10.1016/j.jsc.2021.08.003
- 1371 1372

- Shuanglong Kan, Anthony Widjaja Lin, Philipp Rümmer, and Micha Schrader. 2022. Certistr: a certified string solver. In
 Proceedings of the 11th ACM SIGPLAN International Conference on Certified Programs and Proofs. 210–224.
- 1375 Charaka Geethal Kapugama, Van-Thuan Pham, Aldeida Aleti, and Marcel Böhme. 2022. Human-in-the-loop oracle learning for semantic bugs in string processing programs. In *Proceedings of the 31st ACM SIGSOFT International Symposium on Software Testing and Analysis.* 215–226.
- Vineeth Kashyap, Kyle Dewey, Ethan A Kuefner, John Wagner, Kevin Gibbons, John Sarracino, Ben Wiedermann, and Ben
 Hardekopf. 2014. JSAI: A static analysis platform for JavaScript. In Proceedings of the 22nd ACM SIGSOFT international
 symposium on Foundations of Software Engineering, 121–132.
- Scott Kausler and Elena Sherman. 2014. Evaluation of String Constraint Solvers in the Context of Symbolic Execution. In Proceedings of the 29th ACM/IEEE International Conference on Automated Software Engineering (Vasteras, Sweden) (ASE '14). ACM, New York, NY, USA, 259–270. https://doi.org/10.1145/2642937.2643003
- Matthias Keil and Peter Thiemann. 2014. Symbolic Solving of Extended Regular Expression Inequalities. (2014).
 arXiv:1410.3227 [cs.FL]
- 1384Adam Kiezun, Vijay Ganesh, Philip J Guo, Pieter Hooimeijer, and Michael D Ernst. 2009. HAMPI: a solver for string1385constraints. In Proceedings of the eighteenth international symposium on Software testing and analysis. 105–116.
- Edward Kmett. 2012. charset: Fast unicode character sets based on complemented PATRICIA tries. https://hackage.haskell. org/package/charset
- Neil Kulkarni, Caroline Lemieux, and Koushik Sen. 2021. Learning Highly Recursive Input Grammars. In 2021 36th IEEE/ACM
 International Conference on Automated Software Engineering (ASE). 456–467. https://doi.org/10.1109/ASE51524.2021.
 9678879
- Kevin J Lang. 1999. Faster algorithms for finding minimal consistent DFAs. Technical Report. NEC Research Institute, 4
 Independence Way, Princeton, NJ.
- Daan Leijen and Erik Meijer. 2001. Parsec: Direct Style Monadic Parser Combinators for the Real World. Technical Report
 UU-CS-2001-35. Department of Information and Computing Sciences, Utrecht University. http://www.cs.uu.nl/research/
 techreps/repo/CS-2001/2001-35.pdf
- Guodong Li, Indradeep Ghosh, and Sreeranga P Rajan. 2011. KLOVER: A symbolic execution and automatic test generation tool for C++ programs. In *Computer Aided Verification: 23rd International Conference, CAV 2011, Snowbird, UT, USA, July* 14-20, 2011. Proceedings 23. Springer, 609–615.
- Tianyi Liang, Nestan Tsiskaridze, Andrew Reynolds, Cesare Tinelli, and Clark Barrett. 2015. A Decision Procedure for
 Regular Membership and Length Constraints over Unbounded Strings. Springer International Publishing, 135–150. https:
 //doi.org/10.1007/978-3-319-24246-0_9
- Blake Loring, Duncan Mitchell, and Johannes Kinder. 2017. ExpoSE: practical symbolic execution of standalone JavaScript. In *Proceedings of the 24th ACM SIGSOFT International SPIN Symposium on Model Checking of Software*. 196–199.
- David MacQueen, Robert Harper, and John Reppy. 2020. The History of Standard ML. *PACM on Programming Languages* 4, HOPL, Article 86 (June 2020), 100 pages. https://doi.org/10.1145/3386336
- 1404Falcon Momot, Sergey Bratus, Sven M Hallberg, and Meredith L Patterson. 2016. The seven turrets of babel: A taxonomy of1405langsec errors and how to expunge them. In 2016 IEEE Cybersecurity Development (SecDev). IEEE, 45–52.
- Manuel Montenegro, Susana Nieva, Ricardo Peña, and Clara Segura. 2020. Extending Liquid Types to Arrays. ACM Transactions on Computational Logic 21, 2, Article 13 (Jan. 2020), 41 pages. https://doi.org/10.1145/3362740
- Donald R. Morrison. 1968. PATRICIA Practical Algorithm To Retrieve Information Coded in Alphanumeric. J. ACM 15, 4
 (oct 1968), 514–534. https://doi.org/10.1145/321479.321481
- Charles Gregory Nelson. 1980. Techniques for Program Verification. Ph. D. Dissertation. Stanford University. A revised
 version was published in June 1981 by Xerox PARC as report number CSL-81-10.
- Chris Okasaki and Andy Gill. 1998. Fast mergeable integer maps. In *ACM SIGPLAN Workshop on ML*. 77–86.
- 1411José Oncina and Pedro García. 1992. Inferring Regular Languages in Polynomial Updated Time. World Scientific, 49–61.1412https://doi.org/10.1142/9789812797902_0004
- Xinming Ou, Gang Tan, Yitzhak Mandelbaum, and David Walker. 2004. Dynamic Typing with Dependent Types. In
 Exploring New Frontiers of Theoretical Informatics, Jean-Jacques Levy, Ernst W. Mayr, and John C. Mitchell (Eds.). 437–450.
 https://doi.org/10.1007/1-4020-8141-3_34
- Saswat Padhi, Rahul Sharma, and Todd Millstein. 2016. Data-Driven Precondition Inference with Learned Features. In Proceedings of the 37th ACM SIGPLAN Conference on Programming Language Design and Implementation (Santa Barbara, CA, USA) (PLDI '16). ACM, New York, NY, USA, 42–56. https://doi.org/10.1145/2908080.2908099
- Changhee Park, Hyeonseung Im, and Sukyoung Ryu. 2016. Precise and scalable static analysis of jQuery using a regular
 expression domain. In *Proceedings of the 12th Symposium on Dynamic Languages*. 25–36.
- 1420 1421

1422	$Terence \ J. \ Parr \ and \ Russell \ W. \ Quong. \ 1995. \ ANTLR: \ A \ predicated-LL(k) \ parser \ generator. \ Software: \ Practice \ and \ Experience$
1423	25, 7 (1995), 789-810. https://doi.org/10.1002/spe.4380250705
1424	Patrick M. Rondon, Ming Kawaguci, and Ranjit Jhala. 2008. Liquid Types. In Proceedings of the 29th ACM SIGPLAN Conference
1425	https://doi.org/10.1145/1375581.1375602
1426	Prateek Saxena, Devdatta Akhawe, Steve Hanna, Feng Mao, Stephen McCamant, and Dawn Song. 2010. A symbolic execution
1427	framework for javascript. In 2010 IEEE Symposium on Security and Privacy. IEEE, 513–528.
1428	Michael Schröder, Marc Goritschnig, and Jürgen Cito. 2023. An Exploratory Study of Ad Hoc Parsers in Python.
1429	arXiv:2304.09733 [cs.SE] https://arxiv.org/abs/2304.09733 Accepted as a registered report for MSR 2023 with Con-
1430	tinuity Acceptance (CA). Mohamed Nassim Seghir and Daniel Kroening, 2013. Countereyample-Guided Precondition Inference. In <i>Proceedings of</i>
1431	the 22nd European Conference on Programming Languages and Systems (Rome, Italy) (ESOP'13), Springer-Verlag, Berlin,
1432	Heidelberg, 451–471. https://doi.org/10.1007/978-3-642-37036-6_25
1433	Axel Simon, Andy King, and Jacob M. Howe. 2003. Two Variables per Linear Inequality as an Abstract Domain. Springer
1434	Berlin Heidelberg, 71-89. https://doi.org/10.1007/3-540-45013-0_7
1435	Ken Thompson. 1968. Programming techniques: Regular expression search algorithm. Commun. ACM 11, 6 (1968), 419–422.
1436	applications. In Proceedings of the 2014 ACM SIGSAC Conference on Computer and Communications Security 1232–1243
1437	Minh-Thai Trinh, Duc-Hiep Chu, and Joxan Jaffar. 2020. Inter-theory dependency analysis for SMT string solvers. <i>Proceedings</i>
1438	of the ACM on Programming Languages 4, OOPSLA (2020), 1–27.
1439	The Unicode Consortium. 2023. The Unicode Standard, Version 15.1.0. The Unicode Consortium, South San Francisco, CA.
1440	https://www.unicode.org/versions/Unicode15.1.0/
1441	Niki Vazou, Eric L. Seidel, Kanjit Jhala, Dimitrios Vytiniotis, and Simon Peyton-Jones. 2014. Refinement Types for Haskell. In Proceedings of the 19th ACM SIGPI AN International Conference on Functional Programming (Cothenburg, Sweden).
1442	(ICFP '14). ACM. New York, NY, USA, 269–282. https://doi.org/10.1145/2628136.2628161
1443	Alessandro Warth and Ian Piumarta. 2007. OMeta: An Object-Oriented Language for Pattern Matching. In Proceedings of the
1444	2007 Symposium on Dynamic Languages (Montreal, Quebec, Canada) (DLS '07). 11-19.
1445	Yunhui Zheng, Xiangyu Zhang, and Vijay Ganesh. 2013. Z3-str: A z3-based string solver for web application analysis. In
1446	Proceedings of the 2013 9th Joint Meeting on Foundations of Software Engineering. 114–124.
1447	
1448	
1449	
1450	
1451	
1452	
1453	
1454	
1455	
1456	
1457	
1458	
1459	
1460	
1461	
1462	
1463	
1464	
1465	
1466	
1467	
1468	
1469	
1107	
1470	