

# Morpheus: Automated Safety Verification of Data-Dependent Parser Combinator Programs

Ashish Mishra  

Department of Computer Science, Purdue University, West Lafayette, IN, USA

Suresh Jagannathan  

Department of Computer Science, Purdue University, West Lafayette, IN, USA

---

## Abstract

Parser combinators are a well-known mechanism used for the compositional construction of parsers, and have shown to be particularly useful in writing parsers for rich grammars with data-dependencies and global state. Verifying applications written using them, however, has proven to be challenging in large part because of the inherently effectful nature of the parsers being composed and the difficulty in reasoning about the arbitrarily rich data-dependent semantic actions that can be associated with parsing actions. In this paper, we address these challenges by defining a parser combinator framework called *Morpheus* equipped with abstractions for defining composable effects tailored for parsing and semantic actions, and a rich specification language used to define safety properties over the constituent parsers comprising a program. Even though its abstractions yield many of the same expressivity benefits as other parser combinator systems, *Morpheus* is carefully engineered to yield a substantially more tractable automated verification pathway. We demonstrate its utility in verifying a number of realistic, challenging parsing applications, including several cases that involve non-trivial data-dependent relations.

**2012 ACM Subject Classification** Software and its engineering → General programming languages

**Keywords and phrases** Parsers, Verification, Domain-specific languages, Functional programming, Refinement types, Type systems

**Digital Object Identifier** 10.4230/LIPIcs.ECOOP.2023.20

**Related Version** *Full Version*: <https://doi.org/10.48550/arXiv.2305.07901>

**Funding** Funding for this work was provided in part by DARPA under the SafeDocs program (grant HR0011-19-C-0073).

**Acknowledgements** The authors thank the reviewers for their insightful and useful comments.

## 1 Introduction

Parsers are transformers that decode serialized, unstructured data into a structured form. Although many parsing problems can be described using simple context-free grammars (CFGs), numerous real-world data formats (e.g., pdf [34], dns [9], zip [35], etc.), as well as many programming language grammars (e.g., Haskell, C, Idris, etc.) require their parser implementations to maintain additional context information during parsing. A particularly important class of context-sensitive parsers are those built from *data-dependent grammars*, such as the ones used in the data formats listed above. Such *data-dependent* parsers allow parsing actions that explicitly depend on earlier parsed data or semantic actions. Often, such parsers additionally use global effectful state to maintain and manipulate context information. To illustrate, consider the implementation of a popular class of *tag-length-data* parsers; these parsers can be used to parse image formats like PNG or PPM images, networking packets formats like TCP, etc., and use a parsed length value to govern the size of the input payload that should be parsed subsequently. The following BNF grammar captures this relation for a simplified PNG image.



© Ashish Mishra and Suresh Jagannathan;

licensed under Creative Commons License CC-BY 4.0

37th European Conference on Object-Oriented Programming (ECOOP 2023).

Editors: Karim Ali and Guido Salvaneschi; Article No. 20; pp. 20:1–20:27

Leibniz International Proceedings in Informatics



LIPICs Schloss Dagstuhl – Leibniz-Zentrum für Informatik, Dagstuhl Publishing, Germany

```

png ::= header . chunk*
chunk ::= length . typespec . content

```

The grammar defines a `header` field followed by zero or more `chunks`, where each `chunk` has a single byte `length` field parsed as an unsigned integer, followed by a single byte `chunk type specifier`. This is followed by zero or more bytes of actual `content`. A useful data-dependent safety property that any parser implementation for this grammar should satisfy is that “*the length of content plus typespec is equal to the value of length*”.

Parser combinator libraries [44, 25, 12, 33] provide an elegant framework in which to write parsers that have such data-dependent features. These libraries simplify the task of writing parsers because they define the grammar of the input language and implement the recognizer for it at the same time. Moreover, since combinator libraries are typically defined in terms of a shallowly-embedded DSL in an expressive host language like Haskell [1, 19] or OCaml [25], parser implementations can seamlessly use a myriad of features available in the host language to express various kinds of data-dependent relations. This makes them capable of parsing both CFGs as well as richer grammars that have non-trivial semantic actions. Consequently, this style of parser construction has been adopted in many domains [2, 1, 33], a fact exemplified by their support in many widely-used languages like Haskell, Scala, OCaml, Java, etc.

Although parser combinators provide a way to easily write data-dependent parsers, verifying the correctness (i.e., ensuring that all data dependencies are enforced) of parser implementations written using them remains a challenging problem. This is in large part due to the inherently effectful nature of the parsers being composed, the pervasive use of rich higher-order abstractions available in the combinators used to build them, and the difficulty of reasoning about complex data-dependent semantic actions triggered by these combinators that can be associated with a parsing action.

This paper directly addresses these challenges. We do so by imposing modest constraints on the host language capabilities available to parser combinator programs; these constraints *enable* mostly automated reasoning and verification, *without* comprising the ability to specify parsers with rich effectful, data-dependent safety properties. We manifest these principles in the design of a deeply-embedded DSL for OCaml called **Morpheus** that we use to express and verify parsers and the combinators that compose them. Our design provides a novel (and, to the best of our knowledge, first) automated verification pathway for this important application class. This paper makes the following contributions:

1. It details the design of an OCaml DSL **Morpheus** that allows compositional construction of *data-dependent* parsers using a rich set of primitive parsing combinators along with an expressive specification language for describing safety properties relevant to parsing applications.
2. It presents an automated refinement type-based verification framework that validates the correctness of **Morpheus** programs with respect to their specifications and which supports fine-grained effect reasoning and inference to help reduce specification annotation burden.
3. It justifies its approach through a detailed evaluation study over a range of complex real-world parser applications that demonstrate the feasibility and effectiveness of the proposed methodology.

The remainder of the paper is organized as follows. The next section presents a detailed motivating example to illustrate the challenges with verifying parser combinator applications and presents a detailed overview of **Morpheus** that builds upon this example. We formalize **Morpheus**’s specification language and type system in Secs. 3 and 4. Details about **Morpheus**’s implementation and benchmarks demonstrate the utility of our framework is given in Sec. 5. Related work and conclusions are given in Secs. 6 and 7, respectively.

<pre> 1 decl ::= typedef . type-expr . id=rawident 2       extern ... 3       ... 4 typename ::= rawident 5 type-exp ::= "int"   "bool" 6 expr ::= ...   id=rawident </pre>	<pre> 1 decl ::= typedef . type-expr . id=rawident [ 2     ¬ id ∈ (!identifiers) 3     {types.add id} 4       ... 5 typename ::= x = rawident [x ∈ (!types)]{ 6     return x} 7 type-exp ::= "int"   "bool" 8 expr ::= ...   id=rawident {identifiers.add id ; 9     return id} </pre>
---	--

■ **Figure 1** Context-free and context-sensitive grammars for C declarations.

## 2 Motivation and Morpheus Overview

To motivate our ideas and give an overview of Morpheus, consider a parser for a simplified C language *declarations, expressions and typedefs* grammar. The grammar must handle context-sensitive disambiguation of *typename*s and *identifiers*<sup>1</sup>. Traditionally, C-parsers achieve this disambiguation via cumbersome *lexer hacks*<sup>2</sup> which use feedback from the symbol table maintained in the parsing into the lexer to distinguish variables from types. Once the disambiguation is outsourced to the lexer-hack, the C-decl grammar can be defined using a context-free-grammar. For instance, the left hand side, Figure 1, presents a simplified context-free grammar production for a C declaration.

Unfortunately, ad-hoc lexer-hacks are both tedious and error prone. Further, this convoluted integration of the lexing and parsing phases makes it challenging to validate the correctness of the parser implementation. A cleaner way to implement such a parser is to disambiguate *typename*s and *identifiers* when parsing by writing an actual context-sensitive parser. One approach would be to define a shared *context* of two non-overlapping lists of *types* and *identifiers* and a stateful-parser using this context. The modified *context-sensitive* grammar is shown in right hand side, Figure 1.

The square brackets show context-sensitive checks e.g.  $[\neg id \in (!identifier)]$  checks that the parsed *rawident* token *id* is not in the list of *identifiers*, while the braces show semantic actions associated with parser reductions, e.g.  $\{typed.add id\}$ , adds the token *id* to *types*, a list of *identifiers* seen thus far in the parse.

Given this grammar, we can use parser combinator libraries [25, 30] in our favorite language to implement a parser for C language declarations. Unfortunately, although cleaner than the using unwieldy lexer hacks, it is still not obvious how we might verify that implementations actually satisfy the desired *disambiguation* property, i.e. *typename*s and *identifiers* do not overlap. In the next section we provide an overview of Morpheus that informally presents our solution to this problem.

### 2.1 Morpheus Surface Language

An important design decision we make is to provide a surface syntax and API very similar to conventional monadic parser combinator libraries like Parsec [25] in Haskell or mParser [30] in OCaml; the core API that Morpheus provides has the signature shown in Figure 3. The library

<sup>1</sup> <https://web.archive.org/web/20070622120718/http://www.cs.utah.edu/research/projects/mso/goofie/grammar5.txt>

<sup>2</sup> <https://www.lysator.liu.se/c/ANSI-C-grammar-1.html>

## 20:4 Morpheus: Automated Safety Verification for Parser Combinators

<pre> 1 let ids = ref [] 2 let types = ref [] 3 type decl = 4   Typedecl of {typeexp:string} 5     ... 6 type expression = 7   Address of expression 8     Cast of string * expression 9     ... 10    Identifier of string 11 expression : PE<sup>stexc</sup> {∀ h,  ldisjoint (sel (h, ids),sel (h, types)) = true} ν : expression result {∀ h, ν, h'.ν = Inl (v1) =&gt;  ldisjoint (sel (h', ids),sel (h', types)) = true} ∧ ν = Inr (Err) =&gt; included(inp,h,h') = true } </pre>	<pre> typedecl : PE<sup>stexc</sup> {∀ h,  ldisjoint (sel (h, ids),sel (h, types)) = true} ν : tdecl result {∀ h, ν, h'.ν = Inl (v1) =&gt;  ldisjoint (sel (h', ids),sel (h', types)) = true} ∧ ν = Inr (Err) =&gt; included(inp,h,h') = true } </pre>	28
<pre> 12 let expression = 13   do<sub>m</sub> char '(' 14     tn ← typename 15     char ')' 16     e ← expression 17     return Cast (tn, e) 18   &lt; &gt; ... 19   &lt; &gt; 20   do<sub>m</sub> 21     id ← identifier 22     let b = List.mem id !types 23     if (!b) then 24       ids := id :: (!ids) 25       return (Identifier id) 26     else 27       fail </pre>	<pre> let typedecl = do<sub>m</sub> td ← keyword "typedef" te ← string "bool" &lt; &gt; string "int" id ← identifier (* incorrect-check: if (not(List.mem id !types)) then*) if (not (List.mem id !ids)) then types := id :: (!types) return Tdecl {typeexp; id} else fail </pre>	29 30 31 32 33 34 35 36 37 38 39 40
<pre> typename : PE<sup>stexc</sup> {∀ h,  ldisjoint (sel (h, ids),sel (h, types)) = true} ν : string result {∀ h, ν, h'.ν = Inl (v) =&gt;  mem (sel (h', types), v) = true ∧ ν = Inr (Err) =&gt; included(inp,h,h') = true} </pre>	<pre> let typename = do<sub>m</sub> x ← identifier if (List.mem x !types) then return x else fail </pre>	41 42 43 44 45 46 47

■ **Figure 2** A simplified C-declaration parser written in Morpheus. Specifications in blue are provided by the programmer; specifications in gray are inferred by Morpheus. Line number 28 represents the complete multiline type specification.

defines a number of primitive combinators: `eps` defines a parser for the empty language, `bot` always fails, and `char c` defines a parser for character `c`. Beyond these, the library also provides a `bind (>>=)` combinator for monadically composing parsers, a `choice (<|>)` combinator to non-deterministically choose among two parsers, and a `fix` combinator to implement recursive parsers. The `return x` is a parser which always succeeds with a value `x`. As we demonstrate, these combinators are sufficient to derive a number of other useful parsing actions such as `many`, `count`, etc. found in these popular combinator libraries. From the parser writer's perspective, Morpheus programs can be expressed using these combinators along with a basic collection of other non-parser expression forms similar to those found in an ML core language, e.g., first-class functions, `let` expressions, references, etc.

```

type 'a t
val eps : unit t
val bot : 'a t
val char : char → char t
val (>>=) : 'a t → (a → 'b t) → 'b t
val <|> : 'a t → 'a t → 'a t
val fix : ('b t → 'b t) → 'b t
val return : 'a → 'a t

```

■ **Figure 3** Signatures of primitive parser combinators supported by Morpheus.

For instance a parser for option  $p$ , which either parses an empty string or anything that  $p$  parses can be written:

```
let option p = (eps >>= λ_. return None) <|> (p >>= λ x. return Some x)
```

We can also define more intricate parsers like *Kleene-star* and *Kleene-plus*:

```
let star p = fix (λ p_star. eps <|> p >>= λ x. p_star >>= λ xs . return (x :: xs) )
let plus p = fix (λ p_star. p <|> p >>= λ x. p_star >>= λ xs . return (x :: xs) )
```

Figure 2 shows a Morpheus implementation that parses a valid C language decl.<sup>3</sup> The parser uses two mutable lists to keep track of types and identifiers. The structure is similar to the original data-dependent grammar, even though the program uses ML-style operators for assignment and dereferencing. For ease of presentation, we have written the program using *do-notation* as syntactic sugar for Morpheus’s monadic bind combinator.

The `typedcl` parser follows the grammar and parses the keyword *typedef* using the keyword parser (not shown).<sup>4</sup> It uses a choice combinator (`<|>`) (line 32), which has a semantics of a non-deterministic choice between two sub-parsers. The interesting case occurs while parsing an identifier (lines 33 - 39), in order to enforce disambiguation between *typenames* and *identifiers*, the parser needs to maintain an invariant that the two lists, `types` for parsed *typenames* and `ids` for parsed *identifiers* are always *disjoint* or *non-overlapping*.

In order to maintain the non-overlapping list invariant, a parsed identifier token (line 33) can be a valid typename only if it is not parsed earlier as an identifier expression. i.e. it is not in the `ids` list. The parser performs this check at (line 35). If this check succeeds, the list of typenames (`types`) is updated and a `decl` is returned, else the parsing fails.

The disambiguation decision is required during the parsing of an *expression*. The expression parser defines multiple choices. The parser for the *casting* expression parses a typename followed by a recursive call to expression. The `typename` parser in turn (line 41) parses an identifier token and checks that the identifier is indeed a `typename` (line 44) and returns it, or fails.

The `ids` list is updated during parsing an identifier expression (line 20), here again to maintain disambiguation, before adding a string to the `ids` list, its non-membership in the current `types` list is checked (line 22).

Although the above parser program is easy to comprehend given how closely it hews to the grammar definition, it is still nonetheless non-trivial to verify that the parser actually satisfies the required disambiguation safety property. For example, an implementation in

<sup>3</sup> For now, ignore the specifications given in `gray` and `blue`.

<sup>4</sup> Morpheus, like other parser combinator libraries provides a library of parsers for parsing keywords, identifiers, natural numbers, strings, etc.

which line 34 is replaced with the commented expression above it would incorrectly check membership on the wrong list. We describe how Morpheus facilitates verification of this program in the following section.

## 2.2 Specifying Data-dependent Parser Properties

Intuitively, verifying the above-given parser for the absence of overlap between the *typenames* and *identifiers* requires establishing the following partial correctness property: if the *types* and *identifiers* lists do not overlap when the *typedecl* parser is invoked, and the parser terminates without an error, then they must not overlap in the output state generated by the parser. Additionally, it is required that the parser consumes some prefix of the input list. Morpheus provides an expressive specification language to specify properties such as these.

Morpheus allows standard ML-style inductive type definitions that can be refined with *qualifiers* similar to other refinement type systems [38, 43, 18]. For instance, we can refine the type of a list of strings to only denote *non-empty* lists as: `type nonempty = {  $\nu$  : [string] | len ( $\nu$ ) > 0 }`. Here,  $\nu$  is a special bound variable representing a list and `(len  $\nu$  > 0)` is a *refinement* where `len` is a *qualifier*, a predicate available to the type system that captures the length property of a list.

### 2.2.1 Specifying effectful safety properties

Standard refinement type systems, however, are ill-suited to specify safety properties for effectful computation of the kind expressible by parser combinators. Our specification language, therefore, also provides a type for effectful computations. We use a specification monad (called a *Parsing Expression*) of the form `PE $\varepsilon$  {  $\phi$  }  $\nu$  :  $\tau$  {  $\phi'$  }` that is parameterized by the *effect* of the computation  $\varepsilon$  (e.g., `state`, `exc`, `nondet`, and their combinations like `stexc` for (both `state` and `exc`), `stnon` (for both `state` and `nondet`), etc.); and Hoare-style pre- and post-conditions [31, 41, 40]. Here,  $\phi$  and  $\phi'$  are first-order logical propositions over qualifiers applied to program variables and variables in the type context. The precondition  $\phi$  is defined over an abstract input heap  $h$  while the postcondition  $\phi'$  is defined over input heap  $h$ , output heap  $h'$ , and the special result variable  $\nu$  that denotes the result of the computation. Using this monad, we can specify a safety property for the *typedecl* subparser as shown at line 28 in Figure 2. The type should be understood as follows: The *effect* label `stexc` defines that the parser may have both `state` effect as it reads and updates the context; and `exc` effect as the parser may fail. The precondition defines a property over a list of identifiers `ids` and a list of typenames `types` in the input heap  $h$  via the use of the built-in qualifier `sel` that defines a select operation on the heap [27]; here,  $\nu$  is bound to the result of the parse. Morpheus also allows user-defined qualifiers, like the qualifier `ldisjoint`. It establishes the *disjointness/non-overlapping* property between two lists. This qualifier is defined using the following definition:

```

qualifier ldisjoint [] l2 → true
    | l1 [] → true
    | (x :: xs) l2 → member (x, l2) = false ∧ ldisjoint (xs, l2)
    | l1 (y :: ys) → member (y, l1) = false ∧ ldisjoint (l1, ys)

```

This definition also uses another qualifier for list membership called `member`. Morpheus automatically translates these user-defined qualifiers to axioms, logical sentences whose validity is assumed by the underlying theorem prover during verification. For instance, given the above qualifier, Morpheus generates axioms like:

```

Axiom1:  $\forall l1, l2 : \alpha \text{ list. } (\text{empty}(l1) \vee \text{empty}(l2)) \Rightarrow \text{ldisjoint}(l1, l2) = \text{true}$ 
Axiom2:  $\forall xs, l2 : \alpha \text{ list, } x : \alpha. \text{ldisjoint}(xs, l2) = \text{true} \wedge \text{member}(x, l2) = \text{false} \Rightarrow \text{ldisjoint}((x::xs), l2) = \text{true}$ 
Axiom3:  $\forall l1, l2 : \alpha \text{ ldisjoint}(l1, l2) \Leftrightarrow \text{ldisjoint}(l2, l1)$ 

```

The specification (at line 28) also uses another qualifier, `included(inp,h,h')`, which captures the monotonic consumption property of the input list `inp`. The qualifier is true when the remainder `inp` after parsing in `h'` is a suffix of the original `inp` list in `h`.

The types for other parsers in the figure can be specified as shown at lines 11, 40, etc.; these types shown in gray are automatically inferred by `Morpheus`'s type inference algorithm. For example, the type for the `typename` parser (line 40) returns an optional string (`result` is a special option type) and records that when parsing is successful, the returned string is added to the `types` list, and when unsuccessful, the input is still monotonically consumed.

## 2.2.2 Verification using Morpheus

Note that the pre-condition in the specification (`ldisjoint(l1, l2) = true`) and the type ascribed to the membership checks in the implementation (line 35) are sufficient to conclude that the addition of a `typename` to the `types` list (line 36) maintains the `ldisjoint` invariant as required by the postcondition.

In contrast, an erroneous implementation that omits the membership check or replaces the check at line 34 with the commented line above it will cause type-checking to fail. The program will be flagged ill-typed by `Morpheus`. For this example, `Morpheus` generated 21 verification conditions (VCs) for the control-path representing a successful parse and generated 5 VCs for the failing branch. We were able to discharge these VCs to the SMT solver Z3 [7], which took 6.78 seconds to verify the former and 1.90 seconds to verify the latter.

## 3 Morpheus Syntax and Semantics

### 3.1 Morpheus Syntax

Figure 4 defines the syntax of  $\lambda_{sp}$ , a core calculus for `Morpheus` programs. The language is a call-by-value polymorphic lambda-calculus with effects, extended with primitive expressions for common parser combinators and a refinement type-based specification language. A  $\lambda_{sp}$  value is either a constant drawn from a set of base types (`int`, `bool`, etc.), as well as a special `Err` value of type `exception`, an abstraction, or a constructor application. Variables bound to updateable locations ( $\ell$ ) are distinguished from variables introduced via function binding ( $x$ ). A  $\lambda_{sp}$  expression  $e$  is either a value, an application of a function or type abstraction, operations to dereference and assign to top-level locations (see below), polymorphic `let` expressions, reference binding expressions, a `match` expression to pattern-match over type constructors, a `return` expression that lifts a value to an effect, and various parser primitive expressions that define parsers for the empty language (`eps`), a character (`char`) parser, and  $\perp$ , a parser that always fails. Additionally, the language provides combinators to monadically compose parsers ( $>>=$ ), to implement parsers defined in terms of a non-deterministic choice of its constituents ( $<|>$ ), and to express parsers that have recursive ( $\mu(x : \tau).p$ ) structure.

We restrict how effects manifest by requiring reference creation to occur only within `let` expressions and not in any other expression context. Moreover, the variables bound to locations so created ( $\ell$ ) can only be dereferenced or assigned to and cannot be supplied



## Expression Language

$c, \text{unit}, \text{Err}$	$\in$	<i>Constants</i>	
$x$	$\in$	<i>Vars</i>	
$\text{inp}, \ell$	$\in$	<i>RefVars</i>	
$v$	$\in$	<i>Value</i>	$::= c \mid \lambda(x : \tau).e \mid \Lambda(\alpha).e \mid D_i \overline{t_k} \overline{v_j}$
$e$	$\in$	<i>Exp</i>	$::= v \mid x \mid p \mid e x \mid e[t] \mid \text{deref } \ell \mid \ell := e$ $\mid \text{let } x = v \text{ in } e \mid \text{let } \ell = \text{ref } e \text{ in } e$ $\mid \text{match } v \text{ with } D_i \overline{\alpha} \overline{x_j} \rightarrow e \mid \text{return } e$
$p$	$\in$	<i>Parsers</i>	$::= \text{eps} \mid \perp \mid \text{char } e \mid (\mu(x : \tau).p)$ $\mid p \gg e \mid p \langle \rangle p$

## Specification Language

$\alpha$	$\in$	<i>TypeVariables</i>	
TN	$\in$	<i>User Defined Types</i>	$::= \alpha \text{ list}, \alpha \text{ tree}, \dots$
$\mathbf{t}$	$\in$	<i>Base Types</i>	$::= \alpha \mid \text{int} \mid \text{bool} \mid \text{unit} \mid \text{heap} \mid \text{TN} \mid \mathbf{t}$ $\text{result} \mid \mathbf{t} \text{ ref} \mid \text{exc}$
$\tau$	$\in$	<i>Type</i>	$::= \{\nu : \mathbf{t} \mid \phi\} \mid (x : \tau) \rightarrow \tau \mid \text{PE}^\varepsilon\{\phi_1\}\nu : \mathbf{t}\{\phi_2\}$
$\varepsilon$	$\in$	<i>Effect Labels</i>	$::= \text{pure} \mid \text{state} \mid \text{exc} \mid \text{nondet} \mid \dots$
$\sigma$	$\in$	<i>Type Scheme</i>	$::= \tau \mid \forall \alpha. \tau$
$Q$	$\in$	<i>Qualifiers</i>	$::= \text{QualifierName}(\overline{x_i})$
$\phi, P$	$\in$	<i>Propositions</i>	$::= \text{true} \mid \text{false} \mid Q \mid Q_1 = Q_2$ $\mid \neg \phi \mid \phi \wedge \phi \mid \phi \vee \phi \mid \phi \Rightarrow \phi \mid \forall(x : \mathbf{t}).\phi$
$\Gamma$	$\in$	<i>Type Context</i>	$::= \emptyset \mid \Gamma, x : \sigma \mid \Gamma, \ell : \tau \text{ ref} \mid \Gamma, \phi$
$\Sigma$	$\in$	<i>Constructors</i>	$::= \emptyset \mid \Sigma, D_i \overline{\alpha_k} \overline{x_j} : \overline{\tau_j} \rightarrow \tau$

■ **Figure 4**  $\lambda_{sp}$  Expressions and Types.

as arguments to abstractions or returned as results since they are not treated as ordinary expressions. This stratification, while arguably restrictive in a general application context, is consistent with how parser applications, such as our introductory example are typically written and, as we demonstrate below, do not hinder our ability to write real-world data-dependent parser implementations. Enforcing these restrictions, however, provides a straightforward mechanism to prevent aliasing of effectful components during evaluation, significantly easing the development of an automated verification pathway in the presence of parser combinator-induced computational effects.

### 3.2 Semantics

Figure 5 presents a big-step operational semantics for  $\lambda_{sp}$  parser expressions; the semantics of other terms in the language is standard. The semantics is defined via an evaluation relation ( $\Downarrow$ ) that is of the form  $(\mathcal{H}; e) \Downarrow (\mathcal{H}'; v)$ . The relation defines how a **Morpheus** expression  $e$  evaluates with respect to a heap  $\mathcal{H}$ , a store of locations to base-type values, to yield a value  $v$ , which can be a normal value or an exceptional one, the latter represented by the exception constant  $\text{Err}$ , and a new heap  $\mathcal{H}'$ .



$$\boxed{(\mathcal{H}; e) \Downarrow (\mathcal{H}'; v)}$$

$$\begin{array}{c}
\text{P-EPS} \frac{}{(\mathcal{H}; \text{eps}) \Downarrow (\mathcal{H}; ())} \quad \text{P-}\perp \frac{}{(\mathcal{H}; \perp) \Downarrow (\mathcal{H}; \text{Err})} \quad \text{P-FIX} \frac{(\mathcal{H}; [\mu x : \sigma.p/x]p) \Downarrow (\mathcal{H}'; v)}{(\mathcal{H}; \mu x : \sigma.p) \Downarrow (\mathcal{H}'; v)} \\
\\
\text{P-CHAR-TRUE} \frac{(\mathcal{H}; e) \Downarrow (\mathcal{H}; 'c') \quad \mathcal{H}(\text{inp}) = ('c' :: s) \quad \mathcal{H}' = \mathcal{H}[\text{inp} \mapsto s]}{(\mathcal{H}; \text{char } e) \Downarrow (\mathcal{H}'; 'c')} \\
\\
\text{P-CHAR-FALSE} \frac{(\mathcal{H}; e) \Downarrow (\mathcal{H}; 'c') \quad \mathcal{H}(\text{inp}) \neq ('c' :: s) \quad \mathcal{H}' = \mathcal{H}[\text{inp} \mapsto \text{inp}]}{(\mathcal{H}; \text{char } e) \Downarrow (\mathcal{H}'; \text{Err})} \\
\\
\text{P-BIND-SUCCESS} \frac{(\mathcal{H}; p) \Downarrow (\mathcal{H}'; v_1) \quad (\mathcal{H}'; e) \Downarrow (\mathcal{H}'; (\lambda x : \tau. e')) \quad (\mathcal{H}'; [v_1/x]e') \Downarrow (\mathcal{H}''; v_2)}{(\mathcal{H}; p \gg e) \Downarrow (\mathcal{H}''; v_2)} \\
\\
\text{P-BIND-ERR} \frac{(\mathcal{H}; p) \Downarrow (\mathcal{H}'; \text{Err})}{(\mathcal{H}; p \gg e) \Downarrow (\mathcal{H}'; \text{Err})} \\
\\
\text{P-CHOICE-L} \frac{(\mathcal{H}; p_1) \Downarrow (\mathcal{H}'; v_1)}{(\mathcal{H}; (p_1 <|> p_2)) \Downarrow (\mathcal{H}'; v_1)} \quad \text{P-CHOICE-R} \frac{(\mathcal{H}; p_2) \Downarrow (\mathcal{H}''; v_2)}{(\mathcal{H}; (p_1 <|> p_2)) \Downarrow (\mathcal{H}''; v_2)}
\end{array}$$

■ **Figure 5** Evaluation rules for  $\lambda_{sp}$  parser expressions.

The empty string parser (rule P-EPS) always succeeds, returning a value of type **unit**, without changing the heap. A “bottom” ( $\perp$ ) parser on the other hand always fails, producing an exception value, also without changing the heap. If the argument  $e$  to a character parser **char** yields value (a char ‘c’), and ‘c’ is the head of the input string (denoted by **inp**) being parsed, the parse succeeds (rule P-CHAR-TRUE), consuming the input and returning ‘c’, otherwise, the parse fails, with the input not consumed and the distinguished **Err** value being returned (rule P-CHAR-FALSE). The fixpoint parser  $\mu x.p$  (P-FIX) allows the construction of recursive parser expressions. The monadic bind parser primitive (rule P-BIND-SUCCESS) binds the result of evaluating its parser expression to the argument of the abstraction denoted by its second argument, returning the result of the evaluating the abstraction’s body (P-BIND-SUCCESS); the P-BIND-ERR rule deals with the case when the first expression fails. Evaluation of “choice” expressions, defined by rules P-CHOICE-L and P-CHOICE-R, introduce an unbiased choice semantics over two parsers allowing non-deterministic choices in parsers.

## 4 Typing $\lambda_{sp}$ Expressions

### 4.1 Specification Language

The syntax of Morpheus’s type system is shown in the bottom of Figure 4 and permits the expression of *base types* such as integers, booleans, strings, etc., as well as a special **heap** type to denote the type of abstract heap variables like **h**, **h’** found in the specifications described below. There are additionally user-defined datatypes **TN** (**list**, **tree**, etc.), a special sum type (**t result**) to define two options of a successful and exceptional result respectively, and a

special exception type. More interestingly, base types can be refined with *propositions* to yield monomorphic refinement types. Such types [41, 38, 43] are either *base refinement types*, refining a base typed term with a refinement; *dependent function types*, in which arguments and return values of functions can be associated with types that are refined by propositions; or a *computation type* specifying a type for an effectful computation.

Effectful computations are refined using an effect specification monad:

$$\text{PE}^\varepsilon \{ \forall h. \phi_1 \} \nu : \mathbf{t} \{ \forall h, \nu, h'. \phi_2 \}$$

that encapsulates a base type  $\mathbf{t}$ , parameterized by an effect label  $\varepsilon$ , with Hoare-style pre- ( $\{ \forall h. \phi_1 \}$ ) and post- ( $\{ \forall h, \nu, h'. \phi_2 \}$ ) conditions. This type captures the behavior of a computation that (a) when executed in a pre-state with input heap  $h$  satisfies proposition  $\phi_1$ ; (b) upon termination, returns a value denoted by  $\nu$  of base type  $\mathbf{t}$  along with output heap  $h'$ ; (c) satisfies a post-condition  $\phi_2$  that relates  $h$ ,  $\nu$ , and  $h'$ ; and (d) whose effect is over-approximated by effect label  $\varepsilon$  [20, 45]. An effect label  $\varepsilon$  is either (i) a **pure** effect that records an effect-free computation; (ii) a **state** effect that signifies a stateful computation over the program heap; (iii) an **exception** effect  $\text{exc}$  that denotes a computation that might trigger an exception; (iv) a **nondet** effect that records a computation that may have non-deterministic behavior; or (v) a *join* over these effects that reflect composite effectful actions. The need for the last is due to the fact that effectful computations are often defined in terms of a composition of effects, e.g. a parser oftentimes will define a computation that has a state effect along with a possible exception effect. To capture these composite effects, base effects can be joined to build a finite lattice that reflects the behavior of computations which perform multiple effectful actions, as we describe below.

Propositions ( $\phi$ ) are first-order predicate logic formulae over base-typed variables. Propositions also include a set of qualifiers which are applications of user-defined uninterpreted function symbols such as `mem`, `size` etc. used to encode properties of program objects, `sel` used to model accesses to the heap, and `dom` used to model membership of a location in the heap, etc. Proposition validity is checked by embedding them into a decidable logic that supports equality of uninterpreted functions and linear arithmetic (EUFLIA).

A type scheme ( $\sigma$ ) is either a monotype ( $\tau$ ) or a universally quantified polymorphic type over type variables expressed in prenex-normal form ( $\forall \alpha. \sigma$ ). A **Morpheus** specification is given as a type scheme.

There are two environments maintained by the **Morpheus** type-checker: (1) an environment  $\Gamma$  records the type of variables, which can include variables introduced by function abstraction as well as bindings to references introduced by `let` expressions, along with a set of propositions relevant to a specific context, and (2) an environment  $\Sigma$  maps datatype constructors to their signatures. Our typing judgments are defined with respect to a typing environment

$$\Gamma ::= . \mid \Gamma, x : \sigma \mid \Gamma, \ell : \tau \text{ ref}$$

that is either empty, or contains a list of bindings of variables to either type schemes or references. The rules have two judgment forms:  $(\Gamma \vdash e : \sigma)$  gives a type for a **Morpheus** expression  $e$  in  $\Gamma$ ; and  $(\Gamma \vdash \sigma_1 <: \sigma_2)$  defines a dependent subtyping rule under  $\Gamma$ .

Since our type expressions contain refinements, we generalize the usual notion of type substitution to reflect substitution within refined types:

$$\begin{aligned} [x_a/x] \{ \nu : \mathbf{t} \mid \phi \} &= \{ \nu : \mathbf{t} \mid [x_a/x] \phi \} \\ [x_a/x] (y : \tau) \rightarrow \tau' &= (y : [x_a/x] \tau) \rightarrow [x_a/x] \tau', y \neq x \\ [x_a/x] \text{PE}^\varepsilon \{ \phi_1 \} \{ \nu : \mathbf{t} \} \{ \phi_2 \} &= \text{PE}^\varepsilon \{ [x_a/x] \phi_1 \} \{ \nu : \mathbf{t} \} \{ [x_a/x] \phi_2 \} \end{aligned}$$

## 4.2 Typing Base Expressions

Figure 6 presents type rules for non-parser expressions. The type rules for non-reference variables, functions, and type abstractions (T-TYP-FUN) are standard. The syntax for function application restricts its argument to be a variable, allowing us to record the argument's (intermediate) effects in the typing environment when typing the application as a whole.

### Base Expression Typing $\boxed{\Gamma \vdash e : \sigma}$

$$\begin{array}{c}
\text{T-VAR} \frac{\Gamma(x) = \sigma}{\Gamma \vdash x : \sigma} \quad \text{T-FUN} \frac{\Gamma, x : \tau_1 \vdash e : \tau_2}{\Gamma \vdash \lambda(x : \tau_1).e : \tau_1 \rightarrow \tau_2} \quad \text{T-TYPAPP} \frac{\Gamma \vdash \Lambda\alpha.e : \forall\alpha.\sigma}{\Gamma \vdash \Lambda\alpha.e[t] : [t/\alpha]\sigma} \\
\\
\text{T-APP} \frac{\Gamma \vdash e_f : (x : \{\nu : t \mid \phi_x\}) \rightarrow \text{PE}^\varepsilon\{\phi\} \quad \nu : t \{\phi'\} \quad \Gamma \vdash x_a : \{\nu : t \mid \phi_x\}}{\Gamma \vdash e_f x_a : [x_a/x]\text{PE}^\varepsilon\{\phi\} \quad \nu : t \{\phi'\}} \\
\\
\text{T-TYPFUN} \frac{\Gamma \vdash e : \sigma \quad \alpha \notin FV(\Gamma)}{\Gamma \vdash \Lambda\alpha.e : \forall\alpha.\sigma} \quad \text{T-LET} \frac{\Gamma \vdash v : \forall\alpha.\sigma \quad \Gamma, x : \forall\alpha.\sigma \vdash e_2 : \sigma'}{\Gamma \vdash \text{let } x = v \text{ in } e_2 : \sigma'} \\
\\
\text{T-RETURN} \frac{\Gamma \vdash e : \{\nu : t \mid \phi\}}{\Gamma \vdash \text{return } e : \text{PE}^{\text{pure}}\{\forall h, \nu, h'. h' = h \wedge \phi\}} \\
\\
\text{T-CAPP} \frac{\Sigma(D_i) = \forall \overline{\alpha_k}. \overline{x_j} : \overline{\tau_j} \rightarrow \tau \quad \forall i, j. \Gamma \vdash v_j : [\overline{t_k}/\overline{\alpha_k}][\overline{v_j}/\overline{x_j}]\tau_j}{\Gamma \vdash D_i \overline{t_k} \overline{v_j} : [\overline{t}/\overline{\alpha}][\overline{v_j}/\overline{x_j}]\tau} \\
\\
\text{T-MATCH} \frac{\Sigma(D_i) = \forall \overline{\alpha_k}. \overline{x_j} : \overline{\tau_j} \rightarrow \tau_0 \quad \Gamma \vdash v : \tau_0 \quad \Gamma_i = \Gamma, \overline{\alpha_k}, \overline{x_j} : \overline{\tau_j} \quad \Gamma_i \vdash e_i : \text{PE}^\varepsilon\{\phi_i\} \quad \nu : t \{\phi_{i'}\}}{\Gamma \vdash \text{match } v \text{ with } D_i \overline{\alpha_k} \overline{x_j} \rightarrow e_i : \text{PE}^\varepsilon\{\forall h. \bigwedge_i (v = D_i \overline{\alpha_k} \overline{x_j}) \Rightarrow \phi_i\} \quad \nu : t \{\forall h, \nu', h'. \bigvee_i \phi_{i'}\}} \\
\\
\text{T-DEREF} \frac{\Gamma \vdash \ell : \text{PE}^{\text{state}}\{\phi_1\} \quad \nu : t \text{ ref } \{\phi_2\}}{\Gamma \vdash \text{deref } \ell : \text{PE}^{\text{state}}\{\forall h. \text{dom}(h, \ell)\} \quad \nu' : t \{\forall h, \nu', h'. \text{sel}(h, \ell) = \nu' \wedge h = h'\}} \\
\\
\text{T-ASSIGN} \frac{\Gamma \vdash e : \{\nu : t \mid \phi\}}{\Gamma \vdash \ell := e : \text{PE}^{\text{state}}\{\forall h. \text{dom}(h, \ell)\} \quad \nu' : t \{\forall h, \nu', h'. \text{sel}(h', \ell) = \nu' \wedge \phi(\nu')\}} \\
\\
\text{T-REF} \frac{\Gamma \vdash v : \{\nu : t \mid \phi\} \quad \Gamma, \ell : \text{PE}^{\text{state}}\{\forall h. \neg \text{dom}(h, \ell)\} \quad \nu' : t \text{ ref } \{\forall h, \nu', h'. \text{sel}(h', \ell) = v \wedge \phi(v) \wedge \text{dom}(h', \ell)\} \vdash e_b : \text{PE}^\varepsilon\{\text{dom}(h, \ell)\} \quad \nu'' : t \{\phi_b'\}}{\Gamma, h_i : \text{heap} \vdash \text{let } \ell = \text{ref } v \text{ in } e_b : \text{PE}^{\varepsilon \sqcup \text{state}}\{\forall h. \neg \text{dom}(h, \ell)\} \quad \nu'' : t \{\forall h, \nu'', h. \text{dom}(h_i, \ell) \wedge \text{sel}(h_i, \ell) = v \wedge \phi(v) \wedge \phi_b'\}}
\end{array}$$

■ **Figure 6** Typing Semantics for Morpheus Base Expressions.

The type rule for the return expression (T-RETURN) lifts its non-effectful expression argument  $e$  to have a computation effect with label **pure**, thereby allowing  $e$ 's value to be used in contexts where computational effects are required; a particularly important example of such contexts are bind expressions used to compose the effects of constituent parsers.

In the constructor application rule (T-CAPP), the expression’s type reflects the instantiation of the type and term variables in the constructor’s type with actual types and terms. A match expression is typed (rule T-MATCH) by typing each of the alternatives in a corresponding extended environment and returning a *unified type*. The pre-condition of the *unified* type is a conjunction of the pre-conditions for each alternative, while the post-condition over-approximates the behavior for each alternative by creating a disjunction of each of the possible alternative’s post-conditions. Location manipulating expressions (T-DEREF and T-ASSIGN) use qualifiers *sel* and *dom* to define constraints that reflect state changes on the underlying heap. The argument  $\ell$  of a dereferencing expression (rule T-DEREF) is associated with a computation type over a **tref** base type. Its pre-condition requires  $\ell$  to be in the domain of the input heap, and its post-condition establishes that  $\ell$ ’s contents is the value returned by the expression and that the heap state does not change. The assignment rule (T-ASSIGN) assigns the contents of a top-level reference  $\ell$  to the non-effectful value yielded by evaluating expression  $e$ . The pre-condition of its computation effect type requires that  $\ell$  is in the domain of the input heap and that  $\ell$ ’s contents in the output heap satisfies the refinement ( $\phi$ ) associated with its r-value. Finally, rule T-REF types a **let** expression that introduces a reference initialized to a value  $v$ . The body is typed in an environment in which  $\ell$  is given a computational effect type. The pre-condition of this type requires that the input heap, i.e., the heap extant at the point when the binding of  $\ell$  to **ref**  $v$  occurs, not include  $\ell$  in its domain; its postcondition constrains  $\ell$ ’s contents to be some value  $v'$  that satisfies the refinement  $\phi$  associated with  $v$ , its initialization expression. The body of the **let** expression is then typed in this augmented type environment.

### 4.3 Typing Parser Expressions

Figure 7 presents the type rules for Morpheus parser expressions. The (T-SUB) rule defines the standard type subsumption rule. The empty string parser typing rule (T-P-EPS) assigns a type with **pure** effect and **unit** return type, while the postcondition establishes the equivalence of the input and the output heaps. The T-P-BOT rule captures the always failing semantics of  $\perp$  with an exception effect **exc** and corresponding return types and return values while maintaining the stability of the input heap. The type rules governing a character parser (T-P-CHAR) is more interesting because it captures the semantics of the success and the failure conditions of the parser. We use a sum type ( $\alpha$  **result**) to define two options representing a successful and exceptional result, resp. (with the **Err** exception value in the latter case), using standard injection functions to differentiate among these alternatives. In the successful case, the returned value is equal to the consumed character, captured by an equality constraint over characters. In the successful case, the structure of the output heap with respect to the parse string **inp** must be the same as the input heap except for the absence of the 'c', the now consumed head-of-string character. In the failing case, the input remains unconsumed. Note that we also join the effect labels (**state**  $\sqcup$  **exc**), highlighting the state and exception effect. These effect labels form a standard join semi-lattice with an ordering relation ( $\leq$ )<sup>5</sup>.

Rule T-P-CHOICE defines the static semantics for a non-deterministic choice parser. It introduces a non-determinism effect to the parser’s composite type. The effect’s precondition requires that either of the choices can occur; we achieve this by restricting it to the conjunction of the two preconditions for the sub-parsers. The disjunctive post-condition requires that both the choices must imply the desired goal postcondition for a composite parser to be well-typed. The effect for the choice expression takes a join over the effects of the choices and the non-deterministic effect.

<sup>5</sup> Details of the effect-labels and their join semi-lattice is provided in the accompanied technical report [28]

Parser Expression Typing  $\boxed{\Gamma \vdash e : \sigma}$ 

$$\text{T-SUB} \frac{\Gamma \vdash e : \sigma_1 \quad \Gamma \vdash \sigma_1 <: \sigma_2}{\Gamma \vdash e : \sigma_2}$$

$$\text{T-P-EPS} \frac{}{\Gamma \vdash \text{eps} : \text{PE}^{\text{pure}} \{\forall h. \text{true}\} \nu : \text{unit} \{\forall h, \nu, h'. h' = h\}}$$

$$\text{T-P-BOT} \frac{}{\Gamma \vdash \perp : \text{PE}^{\text{exc}} \{\forall h. \text{true}\} \nu : \text{exc} \{\forall h, \nu, h'. h' = h \wedge \nu = \text{Err}\}}$$

$$\text{T-P-CHAR} \frac{\begin{array}{c} \Gamma \vdash e : \{\nu' : \text{char} \mid \nu' = 'c'\} \\ \phi_2 = \forall h, \nu, h'. \forall x, y. \\ (\text{Inl}(x) = \nu \implies x = 'c' \wedge \text{upd}(h', h, \text{inp}, \text{tail}(\text{inp}))) \wedge \\ (\text{Inr}(y) = \nu \implies y = \text{Err} \wedge \text{sel}(h, \text{inp}) = \text{sel}(h', \text{inp})) \end{array}}{\Gamma \vdash \text{char } e : \text{PE}^{\text{state} \sqcup \text{exc}} \{\forall h. \text{true}\} \nu : \text{char result} \{\phi_2\}}$$

$$\text{T-P-CHOICE} \frac{\Gamma \vdash p_1 : \text{PE}^{\varepsilon} \{\phi_1\} \nu_1 : \tau \{\phi'_1\} \quad \Gamma \vdash p_2 : \text{PE}^{\varepsilon} \{\phi_2\} \nu_2 : \tau \{\phi'_2\}}{\Gamma \vdash (p_1 < | > p_2) : \text{PE}^{\varepsilon \sqcup \text{nondet}} \{(\phi_1 \wedge \phi_2)\} \nu : \tau \{(\phi'_1 \vee \phi'_2)\}}$$

$$\text{T-P-FIX} \frac{\Gamma, x : (\text{PE}^{\varepsilon} \{\phi\} \nu : \mathbf{t}\{\phi'\}) \vdash p : \text{PE}^{\varepsilon} \{\phi\} \nu : \mathbf{t}\{\phi'\} \quad x \notin FV(\phi, \phi')}{\Gamma \vdash \mu x : (\text{PE}^{\varepsilon} \{\phi\} \nu : \mathbf{t}\{\phi'\}). p : \text{PE}^{\varepsilon} \{\phi\} \nu : \mathbf{t}\{\phi'\}}$$

$$\text{T-P-BIND} \frac{\begin{array}{c} \Gamma \vdash p : \text{PE}^{\varepsilon} \{\phi_1\} \nu : \mathbf{t}\{\phi_{1'}\} \quad \Gamma \vdash e : (x : \tau) \rightarrow \text{PE}^{\varepsilon} \{\phi_2\} \nu' : \mathbf{t}'\{\phi_{2'}\} \\ \Gamma' = \Gamma, x : \tau, h_i : \text{heap} \quad h_i \text{ fresh} \\ \Gamma' \vdash p \gg e : \\ \text{PE}^{\varepsilon} \{\forall h. \phi_1 h \wedge \phi_{1'}(h, x, h_i) \Rightarrow \phi_2 h_i\} \\ \nu' : \mathbf{t}' \text{ result} \\ \{\forall h, \nu', h', y. (x \neq \text{Err} \Rightarrow \nu' = \text{Inl } y \wedge \phi_{1'}(h, x, h_i) \wedge \phi_{2'}(h_i, y, h')) \wedge \\ (x = \text{Err} \Rightarrow \nu' = \text{Inr } \text{Err} \wedge \phi_{1'}(h, x, h_i))\} \end{array}}{\Gamma \vdash p \gg e : \text{PE}^{\varepsilon} \{\forall h. \phi_1 h \wedge \phi_{1'}(h, x, h_i) \Rightarrow \phi_2 h_i\} \nu' : \mathbf{t}' \{\phi_{2'}\}}$$

Subtyping  $\boxed{\Gamma \vdash \sigma_1 <: \sigma_2}$ 

$$\begin{array}{c} \text{T-SUB-BASE} \frac{\Gamma \vdash \{\nu : \mathbf{t} \mid \phi_1\} \quad \Gamma \vdash \{\nu : \mathbf{t} \mid \phi_2\}}{\Gamma \vDash \phi_1 \Rightarrow \phi_2} \quad \text{T-SUB-SCHEMA} \frac{\Gamma \vdash \sigma_1 <: \sigma_2}{\Gamma \vdash \forall \alpha. \sigma_1 <: \forall \alpha. \sigma_2} \\ \text{T-SUB-ARROW} \frac{\Gamma \vdash \tau_{21} <: \tau_{11} \quad \Gamma \vdash \tau_{12} <: \tau_{22}}{\Gamma \vdash (x : \tau_{11}) \rightarrow \tau_{12} <: (x : \tau_{21}) \rightarrow \tau_{22}} \quad \text{T-SUB-TVAR} \frac{}{\Gamma \vdash \alpha <: \alpha} \\ \text{T-SUB-COMP} \frac{\Gamma \vDash \phi_2 \Rightarrow \phi_1 \quad \Gamma \vdash \tau_1 <: \tau_2 \quad \Gamma \vdash \varepsilon_1 \leq \varepsilon_2 \quad \Gamma, \phi_2 \vDash (\phi_{1'} \Rightarrow \phi_{2'})}{\Gamma \vdash \text{PE}^{\varepsilon_1} \{\phi_1\} \tau_1 \{\phi_{1'}\} <: \text{PE}^{\varepsilon_2} \{\phi_2\} \tau_2 \{\phi_{2'}\}} \end{array}$$

■ **Figure 7** Typing semantics for primitive parser expressions and subtyping rules.

Rule (T-P-FIX) defines the semantics for the terminating recursive fix-point combinator. Given an annotated type  $\tau$  for the parameter  $x$ , if the type of the body  $p$  in an extended environment which has  $x$  mapping to  $\tau$ , is  $\tau$ , then  $\tau$  is also a valid type for a recursive

fixpoint parser expression. The T-P-BIND rule defines a typing judgement for the exceptional monadic composition of a parser expression  $p$  with an abstraction  $e$ . The composite parser is typed in an extended environment  $(\Gamma')$  containing a binding for the abstraction's parameter  $x$  and an intermediate heap  $h_i$  that acts as the output/post-state heap for the first parser and the input/pre-state for the second. The relation between these heaps is captured by the inferred pre-and post-conditions for the composite parser. There are two possible scenarios depending upon whether the first parser  $p$  results in a success (i.e.  $x \neq \text{Err}$ ) or a failure ( $x = \text{Err}$ ).

In the successful case, the inferred conditions capture the following properties: a) the output of the combined parser is a success; b) the post-condition for the first expression over the intermediate heap  $h_i$  and the output variable  $x$  should imply the precondition of the second expression (required for the evaluation of the second expression); and, c) the overall post-condition relates the post-condition of the first with the precondition of the second using the intermediate heap  $h_i$ . The case when  $p$  fails causes the combined parser to fail as well, with the post-condition after the failure of the first as the overall post-condition. Note that the core calculus is sub-optimal in size since  $\lambda_{sp}$  supports both `return` and `eps`, even though the latter could be modeled using `return`. However, this design choice enables decidable typechecking by limiting the combination of higher-order functions, combinators and states. This is achieved using a limited bind  $p \gg e$ , rather than the general  $e \gg e$ , allowing for the definition of semantic actions  $e$  that only perform limited state manipulation, i.e., reading and updating locations. Thus  $\gg$  and  $\langle | \rangle$  only take parser arguments; thus,  $\text{eps} \langle | \rangle p$  is not equivalent to  $(\text{return } () \langle | \rangle p)$ , in fact the latter is disallowed. Another such design restriction shows up in the typing rules, e.g., the typing rule for function application (T-APP) restricts the arguments to be of *basetype*, thus disallowing expressions returning abstractions or computations, like `return`  $(\lambda x. e)$  or `return`  $(x := e)$ . A more general definition for  $\gg$  will allow valid HO arguments, like  $\lambda x. e \gg e1$ , but translating such general HO stateful programs to decidable logic fragments is not always feasible, as is discussed in other fully dependent type systems [41].

The subtyping rules enable the propagation of refinement type information and relate the subtyping judgments to logical entailment. The subtyping rule for a base refinement (T-SUB-BASE) relates subtyping to the logical implication between the refinement of the subtype and the supertype. The (T-SUB-ARROW) rule defines subtyping between two function refinement types. The (T-SUB-COMP) rule for subtyping between computation types follows the standard Floyd-Hoare rule for *consequence*, coupled with the subtyping relation between result types and an ordering relation between effects ( $\leq$ ). The subtyping rule for type variables (T-Sub-TVar) relates each type variable to itself in a reflexive way, while the subtyping for a type-schema lifts the subtyping relation from a schema to another schema.

#### 4.4 Properties of the Type System

► **Definition 1** (Environment Entailment  $\Gamma \models \phi$ ). Given  $\Gamma = \dots, \overline{\phi_i}$ , the entailment of a formula  $\phi$  under  $\Gamma$  is defined as  $(\bigwedge_i \phi_i) \implies \phi$

In the following,  $\Gamma \models \phi(\mathcal{H})$  extends the notion of semantic entailment of a formula over an abstract heap  $\Gamma \models \phi(h)$  to a concrete heap using an interpretation of concrete heap  $\mathcal{H}$  to an abstract heap  $h$  and the standard notion of well-typed *stores*  $(\Gamma \vdash \mathcal{H})$ .<sup>6</sup>

To prove soundness of Morpheus typing, we first state a soundness lemma for pure expressions (i.e. expressions with non-computation type).

<sup>6</sup> Details are provided in the accompanied technical report [28].

► **Lemma 2** (Soundness Pure-terms). *If  $\Gamma \vdash e : \{ \nu : t \mid \phi \}$  then:*

- Either  $e$  is a value with  $\Gamma \models \phi (e)$
- OR Given there exists a  $v$ , such that  $(\mathcal{H}; e) \Downarrow (\mathcal{H}; v)$  then  $\Gamma \vdash v : t$  and  $\Gamma \models \phi (v)$

► **Theorem 3** (Soundness Morpheus). *Given a specification  $\sigma = \forall \bar{\alpha}. \text{PE}^\varepsilon \{ \phi_1 \} \nu : t \{ \phi_2 \}$  and a Morpheus expression  $e$ , such that under some  $\Gamma$ ,  $\Gamma \vdash e : \sigma$ , then if there exists  $\mathcal{H}$  such that  $\Gamma \models \phi_1(\mathcal{H})$  then:*

1. Either  $e$  is a value, and:  $\Gamma, \phi_1 \models \phi_2 (\mathcal{H}, e, \mathcal{H})$
2. Or, if there exists an  $\mathcal{H}'$  and  $v$  such that  $(\mathcal{H}; e) \Downarrow (\mathcal{H}'; v)$ , then
  - ◻  $\exists \Gamma', \Gamma \subseteq \Gamma'$  and (consistent  $\Gamma \Gamma'$ ), such that:
    - a.  $\Gamma' \vdash v : t$ .
    - b.  $\Gamma', \phi_1(\mathcal{H}) \models \phi_2(\mathcal{H}, v, \mathcal{H}')$

where (consistent  $\Gamma \Gamma'$ ) is a Boolean-valued function that ensures that  $\forall x \in (\text{dom}(\Gamma) \cap \text{dom}(\Gamma'))$ .  $\Gamma \vdash x : \sigma \implies \Gamma' \vdash x : \sigma$ . Additionally,  $\forall \phi. \Gamma \models \phi \implies \Gamma' \models \phi$ .

**Proof.** The soundness proof is by induction on typing rules in Figures 6 and 7, proving the soundness statement against the evaluation rules in Figures 5.<sup>7</sup> ◀

## 5 Evaluation

### 5.1 Implementation

Morpheus is implemented as a deeply-embedded DSL in OCaml<sup>8</sup> equipped with a refinement-type based verification system. It encodes the typing rules given in Section 4 and a parser translating an OCaml-based surface language of the kind presented in our motivating example to the Morpheus core, described in Section 3. To allow Morpheus programs to be easily used in an OCaml development, its specifications can be safely erased once the program has been type-checked. Note that a Morpheus program, verified against a safety specification is guaranteed to be safe when erased since verification takes place against a stricter memory abstraction; in particular, since Morpheus programs are free of aliasing by construction and thus remain so when evaluated as an ML program. This obviates the need for a separate interpreter/compilation phase and gives Morpheus-verified parsers efficiency comparable to the parsers written using OCaml parser-combinator libraries [30, 3].

Morpheus specifications typically require meaningful qualifiers over inductive data-types, beyond those discussed in our core language; in addition to the qualifiers discussed previously, typical examples include qualifiers to capture properties such as the length of a list, membership in a list, etc. Morpheus provides a way for users to write simple inductive propositions over inductive data types, translating them to axioms useful for the solver, in a manner similar to the use of *measures* and *predicates* in other refinement type works [38, 42]. For example, a qualifier for capturing the length property of a list can be written as:

```
qualifier len [] → 0 | len (x :: xs) → len (xs) + 1.
```

Morpheus generates the following axiom from this qualifier:

```
∀ xs : α list, x : α. len (x :: xs) = len (xs) + 1 ∧ len [] = 0
```

<sup>7</sup> The **decidability** theorem and proofs for all the theorems are provided in the technical report [28].

<sup>8</sup> <https://github.com/aegis-iisc/morpheus.git>



Morpheus is implemented in approximately 9K lines of OCaml code. The input to the verifier is a Morpheus program definition, correctness specifications, and any required qualifier definitions. Given this, Morpheus infers types for other expressions and component parsers, generates first-order verification conditions using the typing semantics discussed earlier, and checks the validity of these conditions.

## 5.2 Results and Discussions

We have implemented and verified the examples given in the paper, along with a set of benchmarks capturing interesting, real-world safety properties relevant to data-dependent parsing tasks. The goal of our evaluation is to consider the effectiveness of Morpheus with respect to generality, expressiveness and practicality. Table 1 shows a summary of the benchmark programs considered. Each benchmark is a Morpheus parser program affixed with a meaningful safety property (last column). The first column gives the name of the benchmark. The second column of the table describes benchmark size in terms of the number of lines of Morpheus code, without the specifications. The third column gives a pair D/P, showing the number of unique derived (D) combinators (like `count`, `many`, etc.) used in the benchmark from the Morpheus library, and the number of primitive (P) parsers (like `string`, `number`, etc.) from the Morpheus library used in the benchmark; the former provides some insight on the usability of our design choices in realizing extensibility. The fourth column lists the size of the grammar along with the number of production rules in the grammar. The fifth column gives the number of verification conditions generated, followed by the time taken to verify them (sixth column). The overall verification time is the time taken for generating verification conditions plus the time Z3 takes to solve these VCs. All examples were executed on a 2.7GHz, 64 bit Ubuntu. The next column quantifies the annotation effort for verification. It gives a ratio ( $\#A/\#Q$ ) of required user-provided specifications (in terms of the number of conjuncts in the specification) to the total specification size (annotated + inferred). User-provided specifications are required to specify a top-level safety property and to specify invariants for fix expressions akin to loop invariants that would be provided in a typical verification task. Finally, the last column gives a high-level description of the data-dependent safety property being verified.

Our benchmarks explore data-dependent parsers from several interesting categories.<sup>9</sup> The first category, represented by `ldris do-block`, `Haskell case-exp` and `Python while-statement`, capture parsing activities concerned with layout and indentation introduced earlier. Languages in which layout is used in the definition of their syntax require context-sensitive parser implementations [1, 2]. We encode a Morpheus parser for a sub-grammar for these languages whose specifications capture the layout-sensitivity property.

The second category, represented by `png` and `ppm` consider data-dependent image formats like PNG or PPM. Verifying data-dependence is non-trivial as it requires verifying an invariant over a monadic composition of the output of one parser component with that of a downstream parser component, interleaved with internal parsing logic.

The next category, captured by `xauction`, `xprotein`, and `health`, represent data-dependent parsing in data-processing pipelines over XML and CSV databases. For `xauction` and `xprotein`, we extend XPath expressions over XML to *dependent* XPath expressions. Given that XPath expressions are analogous to regular-expressions over structured XML data, *dependent* XPath expressions are analogous to dependent regular-expressions over XML.

---

<sup>9</sup> The grammar for each of our implementations is given in the technical report [28].

■ **Table 1** Summary of Benchmarks : #Loc Loc defines the size of the parser implementation in Morpheus; D/P gives the number of derived/primitive combinator uses in the parser implementation; grammar size  $G(\# \text{ prod})$  defines size of the grammar along with the number of production rules in the grammar; #VCs defines number of VCs generated; T(s) is the time for discharging these VCs in seconds;  $(\#A/\#Q)$  defines the ratio of number of conjuncts used in the specification provided by the user ( $\#A$ ) to the total number of conjuncts ( $\#Q$ ) across all files in the implementation; Property gives a high-level description of the data-dependent safety property.

Name	# Loc	D/P	G(#prod)	# VCs	T (s)	(#A/#Q)	data-dependence
haskell	110	5/4	20 (7)	17	8.11	9/39	layout-sensitivity
idris	115	5/5	22(8)	33	10.46	7/26	layout-sensitivity
python	47	3/3	25 (7)	23	7.44	6/20	layout-sensitivity
ppm	46	5/2	21 (7)	20	5.33	4/9	tag-length-data
png chunk	30	3/4	10 (2)	12	3.38	2/7	tag-length-data
xauction	54	4/4	31 (10)	19	6.70	2/8	data-dependent XPath expression
xprotein	45	3/3	24(6)	22	6.23	4/10	data-dependent XPath expression
health	40	4/3	15(5)	13	4.56	2/8	data-dependent CSV pattern-matching
c typedef	60	4/4	14 (5)	21	6.78	4/16	context-sensitive disambiguation
streams	51	4/2	12 (4)	16	5.21	2/9	safe stream manipulation

We use these expressions to encode a property of the XPath query over XML data for an online auction and protein database, resp. Note that verifying such properties over XPath queries is traditionally performed manually or through testing. In the case of **health**, we extend regular custom pattern-matching over CSV files to stateful custom pattern-matching, writing a data-dependent custom pattern matcher. We verify that the parser correctly checks relational properties between different columns in the database.

The next two categories have one example each: we introduced the **c typedef** parser in Section 2 that uses data dependence and effectful data structures to disambiguate syntactic categories (e.g., *typenamees* and *identifiers*) in a language definition. Benchmark **streams** defines a parser over streams (i.e. input list indexed with natural numbers).

### 5.2.1 Annotation overhead vs inference

There are some interesting things to note in the second to last column( $\#A/\#Q$ ); First, as the benchmarks (grammars) become more complex, i.e., have a greater number of functions (sub-parsers), the ratio decreases (small is better). In other words, the gains of type-inference become more visible (e.g., **haskell**, **idris**, **c typedef**). The worst (highest) ratio is for the PPM parser. This parser is interesting because, even though the grammar is small, it makes multiple calls to fixpoint combinators. Thus, the user must provide specifications for the top-level parser and each fix-point combinator, thus increasing ( $\#A$ ). Additionally, given a small number of functions (sub-parsers) due to small grammar size, the gains due to inference are also low. In summary, these trends show that the efforts needed for verification are at par with other Refinement typed languages like, Liquid Types [38], FStar [41], etc, and as the parsers become bigger, the benefits of inference become more prominent.

<pre> DoBlock ::= 'do' OpenBlock Do*           CloseBlock; Do ::=   'let' Name TypeSig '=' Expr     'let' Expr '=' Expr     Name '←' Expr     Expr '←' Expr     Ext Expr     Expr </pre>	<pre> expr = do   t ← term   symbol "+"   e ← expr   pure t + e   symbol '*' </pre>
--	---

(a) An Idris grammar rule for a do block.

(b) An input to the parser.

■ **Figure 8** An Idris grammar rule for a do block and an example input.

### 5.3 Case Study: Indentation Sensitive Parsers

As a case study to illustrate Morpheus’s capabilities, we consider a particular class of stateful parsers that are *indentation-sensitive*. These parsers are characterized by having indentation or layout as an essential part of their grammar. Because indentation sensitivity cannot be specified using a context-free grammar, their specification is often specified via an orthogonal set of rules, for example, the offside rule in Haskell.<sup>10</sup> Haskell language specifications define these rules in a complex routine found in the lexing phase of the compiler [26]. Other indentation-sensitive languages like Idris [4] use parsers written using a parser combinator libraries like Parsec or its variants [25, 19] to enforce indentation constraints.

Consider the Idris grammar fragment shown in Figure 8a. The grammar defines the rule to parse a `do`-block. Such a block begins with the `do` keyword, and is followed by zero or more `do` statements that can be `let` expressions, a binding operation (`←`) over names and expressions, an external expression, etc. The Idris documentation specifies the indentation rule in English governing where these statements must appear, saying that the “*indentation of each do statement in a do-block Do\* must be greater than the current indentation from which the rule is invoked* [13].” Thus, in the Idris code fragment shown in Figure 8b, indentation sensitivity constraints require that the last statement is not a part of the `do`-block, while the inner four statements are. A correct Idris parser must ensure that such indentation rules are preserved.

Figure 9 presents a fragment of the parser implementation in Haskell for the above grammar, taken from the Idris language implementation source, and simplified for ease of explanation. The implementation uses Haskell’s Parsec library, it implements indentation rules using a state abstraction (called `IState`) that stores the current indentation level as parsing proceeds. The parser then manually performs reads and updates to this state and performs indentation checks at appropriate points in the code (e.g. line 24, 53). The `IdrisParser` (line 8) is defined in terms of Parsec’s parser monad over an Idris state (here, `IState`), which along with other fields has an integer field (`ist`) storing the current indentation value. A typical indentation check (e.g. see lines 22 - 24) fetches the current value of `ist` using `getIst`, fetches the indentation of the next lexeme using the Parsec library function `indent`, and compares these values.

The structure of the implementation follows the grammar (Figure 8a): the `doBlock` parser parses a reserved keyword “`do`” followed by a block of `do_` statement lists. The indentation is enforced using the parser `indentedDoBlock` (defined at line 49) that gets the

<sup>10</sup><https://www.haskell.org/onlinereport/haskell2010/haskellch10.html>

```

1  data IState = IState {
2    ist :: Int
3    ...
4  } deriving (Show)
5  data PTerm = PDoBlock [PDo]
6  data PDo t = DoExp t | DoExt t
7             | DoLet t t | ...
8  type IdrisParser a = Parser IState a
9
10 getIst :: IdrisParser IState
11 getIst = get
12 putIst :: (i : Int) → IdrisParser ()
13 putIst i = put {ist = i}
14
15 doBlock :: IdrisParser PTerm
16 doBlock = do
17     reserved "do"
18     ds ← indentedDoBlock
19     return (PDoBlock ds)
20 indentedDo :: IdrisParser (PDo PTerm)
21 indentedDo = do
22     allowed ← ist getIst
23     i ← indent
24     if (i <= allowed)
25         then fail ("end of block")
26         else do_
27 indent :: IdrisParser Int
28 indent =
29     do
30     if (lookAheadMatches (operator)) then
31         do
32             operator
33             return (sourceColumn.getSourcePos)
34     else
35         return (sourceColumn.getSourcePos)
36 do_ :: IdrisParser (PDo PTerm)
37 do_ = do
38     reserved "let"
39     i ← name
40     reservedOp "="
41     e ← expr
42     return (DoLet i e)
43 <|> do
44     e ← expr
45     return (DoExt i e)
46 <|> do e ← expr
47     return (DoExp e)
48 indentedDoBlock :: IdrisParser [PDo PTerm]
49 indentedDoBlock =
50     do
51         allowed ← ist getIst
52         lvl' ← indent
53         if (lvl' > allowed) then
54             do
55                 putIst lvl'
56                 res ← many (indentedDo)
57                 putIst allowed
58                 return res
59             else fail "Indentation error"
60
61 lookAheadMatches :: IdrisParser a →
62 IdrisParser Bool
63 lookAheadMatches p =
64     do
65         match ← lookAhead (optional p)
66         return (isJust match)

```

■ **Figure 9** A fragment of a Parsec implementation for Idris `do`-blocks with indentation checks.

current indentation value (`allowed`) and the indentation for the next lexeme using `indent`, checks that the indentation is greater than the current indentation (line 53) and updates the current indentation so that each `do` statement is indented with respect to this new value. It then calls a parser combinator `many` (line 56), which is the Parsec combinator for the Kleene-star operation, over the result of `indentedDo`, i.e., `indentedDo*`. The `indentedDo` parser again performs a manual indentation check, comparing the indentation value for the next lexeme against the block-start indentation (set earlier by `indentedDoBlock` at line 55) and, if successful, runs the actual `do_` parser (line 26). Finally, `indentedDoBlock` resets the indentation value to the value before the block (line 57).

Unfortunately, it is non-trivial to reason that these manual checks suffice to enforce the indentation sensitivity property we desire. Since they are sprinkled throughout the implementation, it is easy to imagine missing or misplacing a check, causing the parser to misbehave. More significantly, the implementation make incorrect assumptions about the

## 20:20 Morpheus: Automated Safety Verification for Parser Combinators

```
1  expr = do
2      t ← term
3      do
4          symbol "+"
5          e ← expr
6          pure t + e
7      "mplus" pure t
```

■ **Figure 10** An input expression that is incorrectly parsed by the implementation shown in Figure 9.

effective actions performed by the library that are reflected in API signatures. In fact, the logic in the above code has a subtle bug [1] that manifests in the input example shown in Figure 10.

Note that the indentation of the token “mplus” is such that it is not a part of either do block; the implementation, however, parses the last statement as a part of the inner do-block, thereby violating the indentation rule, leading to the program being incorrectly parsed.

The problem lies in a mismatch between the contract provided by the library’s indent function and the assumptions made about its behavior at the check at line 24 in the indentedDo parser (or similarly at line 53). Since checking indentation levels for each character is costly, indent is implemented (line 28) in a way that causes certain lexemes (user defined operators like “mplus”) to be ignored during the process of computing the next indentation level. It uses a lookAheadMatches parser to skip all lexemes that are defined as operators. In this example, indent does not check the indentation of lexeme “mplus”, returning the indentation of the token pure instead. Thus, the indentation of the last statement is considered to start at pure, which incorrectly satisfies the checks at line 24 or line 53, and thus causes this statement to be accepted as part of indentedDoBlock. Unearthing and preventing such bugs is challenging. We show how implementing the same parser in Morpheus allows us to catch the bug and verify a correct version of the parser. Figure 11 shows a Morpheus implementation for a portion of the Idris doBlock parser from Figure 9 showing the implementation of three parsers for brevity, doBlock, indentedDo, and indent, along with other helper functions. The structure is similar to the original Haskell implementation.

To specify an *indentation-sensitivity* safety property, we first define an inductive type for a parse-tree (tree) and refine this type using a dependent function type, (offsideTree i), that specifies an indentation value for each parsed result.

```
type tree = Tree {term : pterm; indentT : int; children : tree list}
type offsideTree i = Tree {term : pterm; indentT : { v : int | v > i }; children : (offsideTree i) list}
```

This type defines a tree with three fields:

- A term of type pterm.
- The indentation (indentT) of a returned parse tree, the refinement constraints on indentT requires its value to be greater than i.
- A list of sub-parse trees (children) for each of the terminals and non-terminals in the current grammar rule’s right-hand side, each of which must also satisfy this refinement.

Morpheus additionally automatically generates *qualifiers* like, indentT, children, etc, for each of the datatype’s constructors and fields with the same name that can be used in type refinements. However, this type is not sufficiently expressive to specify the required safety property for doBlock that requires “the indentation of the parse tree returned by doBlock must be greater than the current value of ist” because ist is an effectful heap variable.

<pre> 1  type <math>\alpha</math> pdo = DoExp of <math>\alpha</math> 2                  DoExt of <math>\alpha</math>   ... 3  type pterm = 4  PDoBlock of ((pterm pdo) list) 5  let ist = ref 0 ... 6 </pre>	<pre> 21 </pre>
<pre> doBlock : PE<sup>stexc</sup> {<math>\forall</math> h, I. sel(h, ist) = I} <math>\nu</math> : (offsideTree I) result {<math>\forall</math> h, <math>\nu</math>, h', I, I'. (<math>\nu</math> = Inl ( _ ) =&gt; (sel (h, ist) = I <math>\wedge</math> sel (h', ist) = I') =&gt; I' = I) <math>\wedge</math> <math>\nu</math> = Inr (Err) =&gt; (sel (h', inp) <math>\subseteq</math> sel (h, inp)) }</pre>	<pre> indentedDo : PE<sup>stexc</sup> {<math>\forall</math> h, I. sel(h, ist) = I } <math>\nu</math> : tree result {<math>\forall</math> h, <math>\nu</math>, h', I, I'. <math>\forall</math> i :int. (i &lt;= I =&gt; sel (h', inp) <math>\subseteq</math> sel (h, inp)) <math>\wedge</math> (i &gt; I =&gt; indentT (<math>\nu</math>) = pos (sel (h, inp)) <math>\wedge</math> children (<math>\nu</math>) = nil}</pre>
<pre> 7  let doBlock = 8      do<sub>m</sub> 9      dot <math>\leftarrow</math> reserved "do" 10     ds <math>\leftarrow</math> indentedDoBlock 11     return Tree {term = PDoBlock ds; 12                  indentT = indentT (dot); 13                  children = (dot :: ds) } 14 </pre>	<pre> 22 23 24 25 26 27 28 29 30 </pre>
<pre> do_ : PE<sup>stexc</sup> {<math>\forall</math> h, I. sel(h, inp) = I} <math>\nu</math> : tree result {<math>\forall</math> h, <math>\nu</math>, h', I, I'. (<math>\nu</math> = Inl ( _ ) =&gt; indentT(<math>\nu</math>) = pos (sel (h, inp)) children (<math>\nu</math>) = nil ) <math>\wedge</math> <math>\nu</math> = Inr (Err) =&gt; (sel (h', inp) <math>\subseteq</math> sel (h, inp)) }</pre>	<pre> sourceColumn : (char * int) list -&gt; int </pre>
<pre> 15  let do_ = ... 16 </pre>	<pre> 31 32 </pre>
<pre> lookAheadMatches : PE<sup>pure</sup> {true} <math>\nu</math> : bool { [h'=h] }</pre>	<pre> indent : PE<sup>state</sup> {true} <math>\nu</math> : int {<math>\forall</math> h, <math>\nu</math>, h'. sel (h', inp) <math>\subseteq</math> sel (h, inp) }</pre>
<pre> 17  lookAheadMatches p = 18      do<sub>m</sub> 19      match <math>\leftarrow</math> lookAhead (optional p) 20      return (isJust match) </pre>	<pre> 33 34 35 36 37 38 39 40 </pre>

■ **Figure 11** Morpheus implementation and specifications for a portion of an Idris Do-block with indentation checks, `dom` is a syntactic sugar for Morpheus’s monadic bind. Specifications given in Blue are provided by the parser writer; Gray specifications are inferred by Morpheus. Line number 21 represents the complete multiline type specification.

We can specify a safety property for a `doBlock` parser as shown on line 6 in Figure 11. Again, the type specification in blue are provided by the programmer. The type should be understood as follows: The effect label (`stexc`) defines that the possible effects produced by the parser include `state` and `exc`. The precondition binds the value of the mutable state variable `ist`, a reference to the current indentation level, to `l` via the use of the built-in qualifier `sel` that defines a select operation on the heap [27]. The return type (`offsideTree l result`) obligates the computation to return a parse tree (or a failure) whose indentation must be greater than `l`. The postcondition constraints that the final value of the indentation is to be reset to its value prior to the parse (a *reset* property) when the parser succeeds (case  $\nu = \text{Inl } (\_)$ ) or that the input stream `inp` is monotonically consumed when the parser fails (case  $\nu = \text{Inr } (\text{Err})$ ). The types for other parsers in the figure can be specified as shown at lines 14, 21, 32, etc.; these types shown in gray are automatically inferred by Morpheus’s type inference algorithm.

### 5.3.1 Revisiting the Bug in the Example

The **bug** described in the previous paragraph is unearthed while typechecking the `indentedDo` implementation or the `indentedDoBlock` implementation. We discuss the case for `indentedDo` case here. To verify that `doBlock` satisfies its specification, Morpheus needs to prove that the type inferred for the body of `indentedDo` (lines 22- 29):

1. has a return type that is of the form, `offsideTree I`. Concretely, the indentation of the returned tree must be greater than the initial value of `ist` (i.e. `indentT ( $\nu$ ) > I`).
2. asserts that the final value of `ist` is equal to the initial value.

Goal (1) is required because `indentedDo` is used by `indentedDoBlock` (see Figure 9), which is then invoked by `doBlock`, where its result constructs the value for `children`, whose type is `offsideTree I list`. Goal (2) is required because `doBlock`'s specified post-condition demands it. Type-checking the body for `indentedDo` yields the type shown at line 21. The two conjuncts in the post-condition correspond to the *then* (failure case) and *else* (success case) branch in the parser's body.

The failure conjunct asserts that the input stream is consumed monotonically if the indentation level is greater than `ist`. The success conjunct is the post-condition of the `do_` parser. This inferred type is, however, too weak to prove goal (1) given above, which requires the combinator to return a parse tree that respects the offside rule. The problem is that `indent`'s type (line 32), inferred as:

```
indent : PEstate{true}  $\nu$  : int { $\forall h, \nu, h'. \text{sel}(h', \text{inp}) \subseteq \text{sel}(h, \text{inp})$ }
```

does not allow us to conclude that `indentedDo` satisfies the indentation condition demanded by `doBlock`, i.e., that it returns a well-typed (`offsideTree I`). This is because the type imposes no constraint between the integer `indent` returns and the function's input heap, and thus offers no guarantees that its result gives the position of the first lexeme of the input list.

We can revise `indent`'s implementation such that it does not skip any reserved operators and always returns the position of the first element of the input list, allowing us to track the indentation of every lexeme:

```
indent : PEstate{true}  $\nu$  : int{ $\forall h, \nu, h'. \nu = \text{pos}(\text{sel}(h, \text{inp})) \wedge \text{sel}(h', \text{inp}) \subseteq \text{sel}(h, \text{inp})$ }
let indent = dom s ← !inp
return (sourceColumn s)
```

This type defines a stronger constraint, sufficient to type-check the revised implementation and raise a type error for the original. For this example, Morpheus generated 33 Verification Conditions (VCs) for the revised successful case and 6 VCs for the failing case. We were able to discharge these VCs to the SMT Solver Z3 [7], yielding a total overall verification time of 10.46 seconds in the successful case, and 2.06 seconds in the case when type-checking failed.

This example highlights several key properties of Morpheus verification: The specification language and the type system allows verifying interesting properties over inductive data types (e.g., the `offsideTree` property over the parse trees). It also allows verifying properties dependent on state and other effects such as the *input consumption* property over input streams (`inp`). Secondly, the annotation burden on the programmer is proportional to the complexity of the top-level safety property that needs to be checked. Finally, the similarities between the Haskell implementation and the Morpheus implementation minimize the idiomatic burden placed on Morpheus users.



## 6 Related Work

**Parser Verification.** Traditional approaches to parser verification involve mechanization in theorem provers like Coq or Agda [29, 6, 10, 21, 39, 23, 15]. These approaches trade-off both automation and expressiveness of the grammar they verify to prove full correctness. Consequently, these approaches cannot verify safety properties of data-dependent parsers, the subject of study in this paper. For instance, RockSalt [29] focuses on regular grammars, while [21, 10] present interpreters for parsing expression grammars (without nondeterminism) and limited semantic actions without data dependence. Jourdan et al. [16] gives a certifying compiler for LR(1) grammars, which translates the grammar into a pushdown automaton and a certificate of language equivalence between the grammar and the automaton. More recently CoStar [23] presents a fully verified parser for the ALL(\*) fragment mitigating some of the limitations of the above approaches. However, unlike *Morpheus*, CoStar does not handle data-dependent grammars or user-defined semantic actions.

Deductive synthesis techniques for parsers like Narcissus [8] and [37] focus mainly on tag-length-payload, binary data formats. Narcissus [8] provides a Coq framework (an `encode_decode` tactic) that can automatically generate correct-by-construction encoders and decoders from a given user format input, albeit for a restricted class of parsers. Notably, the system is not easily extensible to complex user-defined data-dependent formats such as the examples we discuss in *Morpheus*. This can be attributed to the fact that the underlying `encode_decode` Coq tactic is complex and brittle and may require manual proofs to verify a new format. In contrast, *Morpheus* enables useful verification capabilities for a larger class of parsers, albeit at the expense of automatic code generation and full correctness. Writing a safe parser implementation for a user-defined format in *Morpheus* is no more difficult than manually building the parser in any combinator framework with the user only having to provide an additional safety specification. EverParse [37] likewise focuses mainly on binary data formats, guaranteeing full-parser correctness, albeit with some expressivity limitations. For example, it does not support user-defined semantic actions or global data-dependences for general data formats. Compared to these efforts, the properties *Morpheus* can validate are more high-level and general. E.g., “non-overlapping of two lists of strings” in a C-decl parser; “layout-sensitivity properties”, etc.. Verifying these properties requires reasoning over a challenging combination of rich algebraic data types, mutable states, and higher-order functions.

[22] also explore types for parsing, defining a core type-system for context-free expressions. However, their goals are orthogonal to *Morpheus* and are targeted towards identifying expressions that can be parsed unambiguously.

**Data-dependent and Stateful Parsers.** *Morpheus* allows writing parsers for data-dependent and stateful parsers. There is a long line of work aimed at writing such parsers [14, 1, 2, 24]. None of these efforts, however, provide a mechanism to reason about the parsers they can express. Further, many of these systems are specialized for a particular class/domain of problems, such as [14] for data-dependent grammars with trivial semantic actions, or [1] for indentation sensitive grammars, etc. *Morpheus* is sufficiently expressive to both write parsers and grammars discussed in many of these approaches, as well as verifying interesting safety properties. Indeed, several of our benchmarks are selected from these works. In contrast, systems such as [14] argue about the correctness of the input parsed against the underlying CFG, a property challenging to define and verify as a *Morpheus* safety property, beyond simple string-patterns and regular expressions. We leave the expression of such grammar-related properties in *Morpheus* as a subject for future work.

**Refinement Types.** Our specification language and type system builds over a refinement type system developed for functional languages like Liquid Types [38] or Liquid Haskell [43]. Extending Liquid Types with *bounds* [42] provides some of the capabilities required to realize data-dependent parsing actions, but it is non-trivial to generalize such an abstraction to complex parser combinators found in Morpheus with multiple effects and local reasoning over states and effects.

**Effectful Verification.** Our work is also closely related to dependent-type-based verification approaches for effectful programs based on monads indexed with either pre- and post-conditions [31, 32] or more recently, predicate monads capturing the weakest pre-condition semantics for effectful computations [41]. As we have illustrated earlier, the use of expressive and general dependent types, while enabling the ability to write rich specifications (certainly richer than what can be expressed in Morpheus), complicates the ability to realize a fully automated verification pathway.

Verification using natural proofs [36] is based on a mechanism in which a fixed set of proof tactics are used to reason about a set of safety properties; automation is achieved via a search procedure over in this set. This idea is orthogonal to our approach where we rather utilize the restricted domain of parsers to remain in a decidable realm. Both our effort and these are obviously incomplete. Another line of work verifying effectful specifications use characteristic formulae [5]; although more expressive than Morpheus types, these techniques do not lend themselves to automation.

**Local Reasoning over Heaps.** Our approach to controlling aliasing is distinguished from substructural typing techniques such as the ownership type system found in Rust [17]. Such type systems provide a much richer and more expressive framework to reason about memory and effects, and can provide useful guarantees like memory safety and data-race freedom etc. Since our DSL is targeted at parser combinator programs which generally operate over a much simplified memory abstraction, we found it unnecessary to incorporate the additional complexity such systems introduce. The integration of these richer systems within a refinement type framework system of the kind provided in Morpheus is a subject we leave for future work.

**Parser Combinators.** There is a long line of work implementing Parser Combinator Libraries and DSLs in different languages [11]. These also include those which provide a principled way for writing stateful parsers using these libraries [1, 24]. As we have discussed, none of these libraries provide an automated verification machinery to reason about safety properties of the parsers. However, since they allow the full expressive power of the host language, they may, in some instances, be more expressive than Morpheus. For example, Morpheus does not allow arbitrary user-defined higher-order functions and builds only on the core API discussed earlier. This may require a more intricate definition for some parsers compared to traditional libraries. For example, traditional parser combinator libraries typically define a higher-order combinator like `many_fold_apply` with the following signature and use this combinator to concisely define a *Kleene-star* parser:

```
many_fold_apply : f : ('b → 'a → 'b) → (a : 'a) → (g : 'a → 'a) → p : ('a, 's) t → ('b, 's) t
let many p = many_fold_apply (fun xs x → x :: xs) [] List.rev p
```

Contrary to this, in Morpheus, we need to define Kleene-star using a more complex, lower-level fixpoint combinator.

## 7 Conclusions

This paper presents Morpheus, a deeply-embedded DSL in OCaml that offers a restricted language of composable effectful computations tailored for parsing and semantic actions and a rich specification language used to define safety properties over the constituent parsers comprising a program. Morpheus is equipped with a rich refinement type-based automated verification pathway. We demonstrate Morpheus's utility by using it to implement a number of challenging parsing applications, validating its ability to verify non-trivial correctness properties in these benchmarks.

---

### References

- 1 Michael D. Adams and Ömer S. Ağacan. Indentation-sensitive parsing for parsec. In *SIGPLAN Notices*, volume 49(12), pages 121–132, New York, NY, USA, September 2014. Association for Computing Machinery. doi:10.1145/2775050.2633369.
- 2 Ali Afrozeh and Anastasia Izmaylova. One parser to rule them all. In *2015 ACM International Symposium on New Ideas, New Paradigms, and Reflections on Programming and Software (Onward!)*, Onward! 2015, pages 151–170, New York, NY, USA, 2015. Association for Computing Machinery. doi:10.1145/2814228.2814242.
- 3 Angstrom. Angstrom parser-combinator library, 2021. URL: <https://github.com/inhabitedtype/angstrom>.
- 4 Edwin Brady. Idris: Implementing a dependently typed programming language. In *Proceedings of the 2014 International Workshop on Logical Frameworks and Meta-Languages: Theory and Practice*, LFMTP '14, New York, NY, USA, 2014. Association for Computing Machinery. doi:10.1145/2631172.2631174.
- 5 Arthur Charguéraud. Characteristic formulae for the verification of imperative programs. *SIGPLAN Not.*, 46(9):418–430, September 2011. doi:10.1145/2034574.2034828.
- 6 Nils Anders Danielsson. Total parser combinators. In *Proceedings of the 15th ACM SIGPLAN International Conference on Functional Programming, ICFP '10*, pages 285–296, New York, NY, USA, 2010. Association for Computing Machinery. doi:10.1145/1863543.1863585.
- 7 Leonardo de Moura and Nikolaj Bjørner. Z3: An efficient smt solver. In C. R. Ramakrishnan and Jakob Rehof, editors, *Tools and Algorithms for the Construction and Analysis of Systems*, pages 337–340, Berlin, Heidelberg, 2008. Springer Berlin Heidelberg.
- 8 Benjamin Delaware, Sorawit Suriyakarn, Clément Pit-Claudel, Qianchuan Ye, and Adam Chlipala. Narcissus: Correct-by-construction derivation of decoders and encoders from binary formats. *Proc. ACM Program. Lang.*, 3(ICFP), July 2019. doi:10.1145/3341686.
- 9 DNS. Domain names – Implementation and specification, 1987. Network Working Group. URL: <https://www.rfc-editor.org/rfc/rfc1035>.
- 10 J. Gross and Adam Chlipala. Parsing parsers a pearl of ( dependently typed ) programming and proof, 2015.
- 11 HaskellWiki. Parsec – HaskellWiki, 2021. [Online; accessed 7-July-2022]. URL: <https://wiki.haskell.org/index.php?title=Parsec&oldid=64649>.
- 12 Graham Hutton and Erik Meijer. Monadic parser combinators, September 1999.
- 13 *Documentation for the Idris Language*, 2017. URL: <https://docs.idris-lang.org/en/latest/index.html>.
- 14 Trevor Jim, Yitzhak Mandelbaum, and David Walker. Semantics and algorithms for data-dependent grammars. In *Proceedings of the 37th Annual ACM SIGPLAN-SIGACT Symposium on Principles of Programming Languages*, POPL '10, pages 417–430, New York, NY, USA, 2010. Association for Computing Machinery. doi:10.1145/1706299.1706347.
- 15 Jacques-Henri Jourdan, François Pottier, and Xavier Leroy. Validating lr(1) parsers. In Helmut Seidl, editor, *Programming Languages and Systems*, pages 397–416, Berlin, Heidelberg, 2012. Springer Berlin Heidelberg.

- 16 Jacques-Henri Jourdan, François Pottier, and Xavier Leroy. Validating lr(1) parsers. In Helmut Seidl, editor, *Programming Languages and Systems*, pages 397–416, Berlin, Heidelberg, 2012. Springer Berlin Heidelberg.
- 17 Ralf Jung, Jacques-Henri Jourdan, Robbert Krebbers, and Derek Dreyer. Rustbelt: Securing the foundations of the rust programming language. *Proc. ACM Program. Lang.*, 2(POPL), December 2017. doi:10.1145/3158154.
- 18 Gowtham Kaki and Suresh Jagannathan. A relational framework for higher-order shape analysis. In *Proceedings of the 19th ACM SIGPLAN International Conference on Functional Programming*, ICFP '14, pages 311–324, New York, NY, USA, 2014. Association for Computing Machinery. doi:10.1145/2628136.2628159.
- 19 Mark Karpov. Megaparsec: Monadic Parser Combinators, 2022. URL: <https://github.com/mrkkp/megaparsec>.
- 20 Shin-ya Katsumata. Parametric effect monads and semantics of effect systems. In *Proceedings of the 41st ACM SIGPLAN-SIGACT Symposium on Principles of Programming Languages*, POPL '14, pages 633–645, New York, NY, USA, 2014. Association for Computing Machinery. doi:10.1145/2535838.2535846.
- 21 Adam Koprowski and Henri Binszok. Trx: A formally verified parser editor=Gordon, Andrew D., interpreter. In *Programming Languages and Systems*, pages 345–365, Berlin, Heidelberg, 2010. Springer Berlin Heidelberg.
- 22 Neelakantan Krishnaswami and Jeremy Yallop. A typed, algebraic approach to parsing. In *Proceedings of the 40th ACM SIGPLAN Conference on Programming Language Design and Implementation*, pages 379–393, June 2019. doi:10.1145/3314221.3314625.
- 23 Sam Lasser, Chris Casinghino, Kathleen Fisher, and Cody Roux. Costar: A verified all(\*) parser. In *Proceedings of the 42nd ACM SIGPLAN International Conference on Programming Language Design and Implementation*, PLDI 2021, pages 420–434, New York, NY, USA, 2021. Association for Computing Machinery. doi:10.1145/3453483.3454053.
- 24 Nicolas Laurent and Kim Mens. Taming context-sensitive languages with principled stateful parsing. In *Proceedings of the 2016 ACM SIGPLAN International Conference on Software Language Engineering*, SLE 2016, pages 15–27, New York, NY, USA, 2016. Association for Computing Machinery. doi:10.1145/2997364.2997370.
- 25 Daan Leijen and Erik Meijer. Parsec: Direct style monadic parser combinators for the real world. Technical Report UU-CS-2001-27, University of Utrecht, July 2001. User Modeling 2007, 11th International Conference, UM 2007, Corfu, Greece, June 25-29, 2007. URL: <https://www.microsoft.com/en-us/research/publication/parsec-direct-style-monadic-parser-combinators-for-the-real-world/>.
- 26 Simon Marlow. Haskell 2010 language report, 2010. URL: <https://www.haskell.org/onlinereport/haskell2010/>.
- 27 J. McCarthy. *Towards a Mathematical Science of Computation*, pages 35–56. Springer Netherlands, Dordrecht, 1993. doi:10.1007/978-94-011-1793-7\_2.
- 28 Ashish Mishra and Suresh Jagannathan. Morpheus: Automated safety verification of data-dependent parser combinator programs, 2023. doi:10.48550/arXiv.2305.07901.
- 29 Greg Morrisett, Gang Tan, Joseph Tassarotti, Jean-Baptiste Tristan, and Edward Gan. Rocksalt: Better, faster, stronger sfi for the x86. In *Proceedings of the 33rd ACM SIGPLAN Conference on Programming Language Design and Implementation*, PLDI '12, pages 395–404, New York, NY, USA, 2012. Association for Computing Machinery. doi:10.1145/2254064.2254111.
- 30 Max Murato. MParser, A Simple Monadic Parser Combinator Library, 2021. URL: <https://github.com/murmour/mparser>.
- 31 Aleksandar Nanevski, Greg Morrisett, and Lars Birkedal. Polymorphism and separation in hoare type theory. *SIGPLAN Not.*, 41(9):62–73, September 2006. doi:10.1145/1160074.1159812.

- 32 Aleksandar Nanevski, Greg Morrisett, Avraham Shinnar, Paul Govereau, and Lars Birkedal. Ynot: Dependent types for imperative programs. *SIGPLAN Not.*, 43(9):229–240, September 2008. doi:10.1145/1411203.1411237.
- 33 Meredith L. Patterson. Hammer primer, 2015. URL: <https://github.com/sergeybratus/HammerPrimer>.
- 34 PDF. Iso 32000 (pdf), 2008. PDF Association. URL: <https://www.pdfa.org/resource/iso-32000-pdf/pdf-2>.
- 35 PKWare. zip file format specification, 2020. URL: <https://pkware.cachefly.net/webdocs/casestudies/APPNOTE.TXT>.
- 36 Xiaokang Qiu, Pranav Garg, Andrei Ştefănescu, and Parthasarathy Madhusudan. Natural proofs for structure, data, and separation. *SIGPLAN Not.*, 48(6):231–242, June 2013. doi:10.1145/2499370.2462169.
- 37 Tahina Ramananandro, Antoine Delignat-Lavaud, Cédric Fournet, Nikhil Swamy, Tej Chajed, Nadim Kobeissi, and Jonathan Protzenko. Everparse: Verified secure zero-copy parsers for authenticated message formats. In *Proceedings of the 28th USENIX Conference on Security Symposium, SEC'19*, pages 1465–1482, USA, 2019. USENIX Association.
- 38 Patrick M. Rondon, Ming Kawaguci, and Ranjit Jhala. Liquid types. In *Proceedings of the 29th ACM SIGPLAN Conference on Programming Language Design and Implementation, PLDI '08*, pages 159–169, New York, NY, USA, 2008. Association for Computing Machinery. doi:10.1145/1375581.1375602.
- 39 Kathleen Fisher Sam Lasser, Chris Casinghino and Cody Roux. A verified ll(1) parser generator. In *ITP*, 2019.
- 40 Wolfram Schulte. VCC: Contract-based modular verification of concurrent C. In *31st International Conference on Software Engineering, ICSE 2009*. IEEE Computer Society, January 2008.
- 41 Nikhil Swamy, Joel Weinberger, Cole Schlesinger, Juan Chen, and Benjamin Livshits. Verifying higher-order programs with the dijkstra monad. In *Proceedings of the 34th ACM SIGPLAN Conference on Programming Language Design and Implementation, PLDI '13*, pages 387–398, New York, NY, USA, 2013. Association for Computing Machinery. doi:10.1145/2491956.2491978.
- 42 Niki Vazou, Alexander Bakst, and Ranjit Jhala. Bounded refinement types. In *Proceedings of the 20th ACM SIGPLAN International Conference on Functional Programming, ICFP 2015*, pages 48–61, New York, NY, USA, 2015. Association for Computing Machinery. doi:10.1145/2784731.2784745.
- 43 Niki Vazou, Eric L. Seidel, Ranjit Jhala, Dimitrios Vytiniotis, and Simon L. Peyton Jones. Refinement types for Haskell. In Johan Jeuring and Manuel M. T. Chakravarty, editors, *Proceedings of the 19th ACM SIGPLAN international conference on Functional programming, Gothenburg, Sweden, September 1-3, 2014*, pages 269–282. ACM, 2014. doi:10.1145/2628136.2628161.
- 44 Philip Wadler. Monads for functional programming. In Manfred Broy, editor, *Program Design Calculi*, pages 233–264, Berlin, Heidelberg, 1993. Springer Berlin Heidelberg.
- 45 Philip Wadler and Peter Thiemann. The marriage of effects and monads. *ACM Trans. Comput. Logic*, 4(1):1–32, January 2003. doi:10.1145/601775.601776.