# PARSER PARSER COMBINATORS FOR PROGRAM TRANSFORMATION

Rijnard van Tonder

@rvtond

Carnegie Mellon University

sourcegraph

# About Me

Automated

Program {Transformation, Analysis, Repair}

Carnegie Mellon University

sourcegraph

# About Me

Automated

Program {Transformation, Analysis, Repair}

Carnegie
Mellon
University

sourcegraph

# About Me

Automated

Program {Transformation, Analysis, Repair}

Code {Search, Intelligence, Review}

Carnegie Mellon University

sourcegraph

# **About Me**

I <3 developer tools

My favorite color is OCaml

# About You
# (and this talk)

Code changes.

# About You
# (and this talk)

# Code changes.

## All the time.

# About You
# (and this talk)



HERACLITVS.

change is the only constant

# About You
# (and this talk)



HERACLITVS.

code
^ change is the only constant

# See the link for demo video

https://drive.google.com/open?id=1ziCXghgdNwheDCnLy1ml5us8fGBzW3CU

# Have this for any language

# Have this for any language

# What I'd like you to take away from this talk:
## A new way to change code

1. A deeper understanding of program syntax, manipulation, and challenges

2. My solution for manipulating syntax
   – Exposure to neat ideas in a functional paradigm

3. Practical tooling & demos

19

# Language Design For Program Manipulation

Eduardus A. T. Merks, J. Michael Dyck, and Robert D. Cameron, *Member, IEEE*

*Abstract*—The design of procedural and object-oriented programming languages is considered with respect to how easily programs written in those languages can be formally manipulated. Current procedural languages such as Pascal, Modula-2, and Ada generally support such program manipulations, except for some annoying anomalies and special cases. Three main areas of language design are identified as being of concern from a manipulation viewpoint—namely, the interface between concrete and abstract syntax, the relationship between the abstract syntax and static semantics (naming, scoping, and typing), and the ability to express basic transformations (folding and unfolding). Design principles are suggested so that problems identified for current languages can be avoided in the future.

*Index Terms*— Language design, program manipulation, language environment, syntax, semantics.

## I. INTRODUCTION

PROGRAMS that manipulate other programs are becoming increasingly important in providing automated assistance for program development and maintenance. In particular, there has been long-standing interest in the concept of interactive program manipulation systems [1], including everything from program transformation systems [2] to language-based editors [3]–[5]. Furthermore, smaller scale program manipulation tools are also of interest, including program instrumenters [6], [7], program restructurers [8], program slicers [9], source-to-source translators [10], [11], and even program-specific manipulation tools that are designed to process a single (presumably large) program or class of programs [12], [13].

The ease of manipulating programs as data objects is strongly influenced by the nature of the *target language* in which programs are written. Two programming lan-

[14], C [15], Modula-2 [16], Ada [17], Eiffel [18], and Modula-3 [19]. Unfortunately, little attention has been paid to manipulation issues in the design of such languages; consequently, they contain a number of avoidable problems in this regard. In contrast, manipulability has been of considerable importance in the design of modern functional languages such as HOPE [20] and Miranda [21], inspired by early work on program transformation [22]. Nevertheless, many useful program manipulations can be carried out in procedural languages. Furthermore, simple design changes for these languages could have considerably alleviated their limitations with respect to manipulation.

To an extent, appropriate language-processing technology can make up for deficiencies in target languages. For example, language-based editors generated by systems such as the Synthesizer Generator [23] or PSG [4] can easily deal with syntactic ambiguity by reference to incrementally maintained semantic attributes. For many other types of program-manipulation application, however, such technologies are either unavailable or inappropriate. In general, our viewpoint is that there is little benefit in using complex technology to solve manipulation problems when those problems could have been avoided altogether by careful language design.

The paper will proceed by proposing various principles of language design that are aimed at ensuring desirable manipulation properties. It is important to emphasize that these principles must be weighed carefully against other concerns that arise during language design, and certainly cannot be considered a recipe for success. Language design involves trade-offs between various desirable properties, and it is our aim in this paper to focus specifically on those that involve manipulation. The importance that a language designer places

14

# Language Design For Program Manipulation

Eduardus A. T. Merks, J. Michael Dyck, and Robert D. Cameron, *Member, IEEE*

*Abstract*—The design of procedural and object-oriented programming languages is considered with respect to how easily gramming languages is considered with respect to how easily written in those languages can be formally manipu-

[14], C [15], Modula-2 [16], Ada [17], Eiffel [18], and Modula-3 [19]. Unfortunately, little attention has been paid to manipulation issues in the design of such languages; conse-

"What are sensible choices for program syntax and semantics if the No. 1 concern is changing code?"

has been long-standing interest in program manipulation systems [1], including everything from program transformation systems [2] to language-based editors [3]–[5]. Furthermore, smaller scale program manipulation tools are also of interest, including program instrumenters [6], [7], program restructurers [8], program slicers [9], source-to-source translators [10], [11], and even program-specific manipulation tools that are designed to process a single (presumably large) program or class of programs [12], [13].

The ease of manipulating programs as data objects is strongly influenced by the nature of the *target language* in

that there is little benefit manipulation problems when those problems could have been avoided altogether by careful language design.

The paper will proceed by proposing various principles of language design that are aimed at ensuring desirable manipulation properties. It is important to emphasize that these principles must be weighed carefully against other concerns that arise during language design, and certainly cannot be considered a recipe for success. Language design involves trade-offs between various desirable properties, and it is our aim in this paper to focus specifically on those that involve

15

*Abstract*—The design of procedural and object-oriented programming languages is considered with respect to how easily programs written in those languages can be formally manipulated. Current procedural languages such as Pascal, Modula-2, and Ada generally support such program manipulations, except for some annoying anomalies and special cases. Three main areas of language design are identified as being of concern from a manipulation viewpoint—namely, the interface between concrete and abstract syntax, the relationship between the abstract syntax and static semantics (naming, scoping, and typing), and the ability to express basic transformations (folding and unfolding). **Design principles are suggested so that problems identified for current languages can be avoided in the future.**

16

*Abstract*—The design of procedural and object-oriented programming languages is considered with respect to how easily programs written in those languages can be formally manipulated. Current procedural languages such as Pascal, Modula-2, and Ada generally support such program manipulations, except for some annoying anomalies and special cases. Three main areas of language design are identified as being of concern from a manipulation viewpoint—namely, the interface between concrete and abstract syntax, the relationship between the abstract syntax and static semantics (naming, scoping, and typing), and the ability to express basic transformations (folding and unfolding). **Design principles are suggested so that problems identified for current languages can be avoided in the future.**

Narrator: they were not avoided

# Just one example:

**concrete syntax** diverges **abstract syntax tree**

```
if (condition) {
    return;
}
```

parse →

```
if (condition)
    return;
```

parse →

if
├─ condition
└─ return

# Just one example:

## concrete syntax ◀ diverges ▶ abstract syntax tree

```
if (condition) {
    return;
}
```

*parse*

```
if (condition)
    return;
```

*parse*



**More work for analysis tools**

# **Multiple languages?**

# Well...

But wouldn't it be nice if...

No one wrote a tool yet

# See the link for demo video

https://drive.google.com/open?id=1xu0Vt_XXyY_9iVT7wmYhSWRfc0cyvltM

No one wrote a tool yet

Syntax extensions



PEP 572

```
if (match := pattern.search(data)) ...
```

No one wrote a tool yet

Syntax extensions

Inline assembly?

```
if (ZEND_CONST_COND(offset == 0, 0))
{
    __asm__ ("mul" LP_SUFF  " %3\n\t"
        "adc $0,%1"
        : "=&a"(res), "=&d" (m_overflow)
        : "%0"(res),
        "rm"(size));
}
```

This is why we can't have nice things

# Lightweight Syntax Transformations & Tooling

# Example:
# Remove redundant nil checks in Go

```
if s != nil {
    for _, x : = range s {
     …
    }
}
```

# Example:
# Remove redundant nil checks in Go

Omit redundant nil check around loop

```go
if s != nil {
    for _, x : = range s {
     ...
    }
}
```

# Example:
# Remove redundant nil checks in Go

Omit redundant nil check around loop

```go
if s != nil {
    for _, x : = range s {
        …
    }
}
```

→

```go
for _, x : = range s {
    …
}
```

# Implementation for redundant checks

```go
func (c *Checker) LintNilCheckAroundRange(j *lint.Job) {
        fn := func(node ast.Node) bool {
                ifstmt, ok := node.(*ast.IfStmt)
                if !ok {
                        return true
                }
                cond, ok := ifstmt.Cond.(*ast.BinaryExpr)
                if !ok {
                        return true
                }
                if cond.Op != token.NEQ || !IsNil(j, cond.Y) || len(ifstmt.Body.List) != 1 {
                        return true
                }
                loop, ok := ifstmt.Body.List[0].(*ast.RangeStmt)
                if !ok {
                        return true
                }
                ifXIdent, ok := cond.X.(*ast.Ident)
                if !ok {
                        return true
                }
                rangeXIdent, ok := loop.X.(*ast.Ident)
                if !ok {
                        return true
                }
                if ifXIdent.Obj != rangeXIdent.Obj {
                        return true
                }
                switch j.Program.Info.TypeOf(rangeXIdent).(type) {
                case *types.Slice, *types.Map:
                        j.Errorf(node, "unnecessary nil check around range")
                }
                return true
        }
        for _, f := range c.filterGenerated(j.Program.Files) {
                ast.Inspect(f, fn)
        }
}
```

# Implementation for redundant checks

```go
func (c *Checker) LintNilCheckAroundRange(j *lint.Job) {
    fn := func(node ast.Node) bool {
        ifstmt, ok := node.(*ast.IfStmt)
        if !ok {

            return true
        }
        if cond.Op != token.NEQ || !IsNil(j, cond.Y) || len(ifstmt.Body.List) != 1 {
            return true
        }
        loop, ok := ifstmt.Body.List[0].(*ast.RangeStmt)
        if !ok {
            return true
        }
        ifXIdent, ok := cond.X.(*ast.Ident)
        if !ok {
            return true
        }
        rangeXIdent, ok := loop.X.(*ast.Ident)
        if !ok {
            return true
        }
        if ifXIdent.Obj != rangeXIdent.Obj {
            return true
        }
        switch j.Program.Info.TypeOf(rangeXIdent).(type) {
        case *types.Slice, *types.Map:
            j.Errorf(node, "unnecessary nil check around range")
        }
        return true
    }
    for _, f := range c.filterGenerated(j.Program.Files) {
        ast.Inspect(f, fn)
    }
}
```

ast.BinaryExpr

Know the AST data structure

35

# Implementation for redundant checks

```go
func (c *Checker) LintNilCheckAroundRange(j *lint.Job) {
    fn := func(node ast.Node) bool {
        ifstmt, ok := node.(*ast.IfStmt)
        if !ok {

                    return true
        }
        if cond.Op != token.NEQ || !IsNil(j, cond.Y) || len(ifstmt.Body.List) != 1 {
            return true
        }
        loop, ok := ifstmt.Body.List[0].(*ast.RangeStmt)
        if !ok {

            return true

        }
        ifXIdent, ok := cond.X.(*ast.Ident)
        if !ok {

            return true

        }
        rangeXIdent, ok := loop.X.(*ast.Ident)
        if !ok {

            return true

        }
        if ifXIdent.Obj != rangeXIdent.Obj {
            return true
        }
        switch j.Program.Info.TypeOf(rangeXIdent).(type) {
        case *types.Slice, *types.Map:
            j.Errorf(node, "unnecessary nil check around range")
        }
        return true

    }
                                                    {

}
```
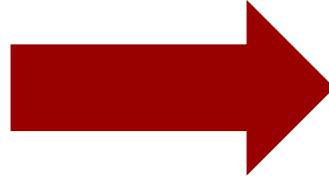
**ast.BinaryExpr**

**Know the AST data structure**

**ast.Inspect**

**Learn the visitor API**

# Implementation for redundant checks

```go
func (c *Checker) LintNilCheckAroundRange(j *lint.Job) {
    fn := func(node ast.Node) bool {
        ifstmt, ok := node.(*ast.IfStmt)
        if !ok {
```

**ast.BinaryExpr**

```
            return true
        }
        if cond.Op != token.NEQ || !IsNil(j, cond.Y) || len(ifstmt.Body.List) != 1 {
            return true
        }
        loop, ok := ifstmt.Body.List[0].(*ast.RangeStmt)
        if !ok {
            return true
        }
        ifXIdent, ok := cond.X.(*ast.Ident)
        if !ok {
            return true
        }
        rangeXIdent, ok := loop.X.(*ast.Ident)
        if !ok {
            return true
        }
        if ifXIdent.Obj != rangeXIdent.Obj {
            return true
        }
        switch j.Program.Info.TypeOf(rangeXIdent).(type) {
        case *types.Slice, *types.Map:
            j.Errorf(node, "unnecessary nil check around range")
        }
        return true
    }
```

**ast.Inspect**

```
                                                            {
}
```

**Know the AST data structure**

**Implement it in your language**

**Learn the visitor API**

# Implementation for redundant checks

```go
func (c *Checker) LintNilCheckAroundRange(j *lint.Job) {
    fn := func(node ast.Node) bool {
        ifstmt, ok := node.(*ast.IfStmt)
        if !ok {
            ...
            return true
        }
        if cond.Op != token.NEQ || !IsNil(j, cond.Y) || len(ifstmt.Body.List) != 1 {
            return true
        }
        loop, ok := ifstmt.Body.List[0].(*ast.RangeStmt)
        if !ok {
            return true
        }
        ifXIdent, ok := cond.X.(*ast.Ident)
        if !ok {
            return true
        }
        rangeXIdent, ok := loop.X.(*ast.Ident)
        if !ok {
            return true
        }
        if ifXIdent.Obj != rangeXIdent.Obj {
            return true
        }
        switch j.Program.Info.TypeOf(rangeXIdent).(type) {
        case *types.Slice, *types.Map:
            j.Errorf(node, "unnecessary nil check around range")
        }
        return true
    }
    {
    }
}
```

**ast.BinaryExpr**

**ast.Inspect**

Know the AST data structure

Implement it in your language

Learn the visitor API

Now do the same for Rust, C, Haskell…

```
if s != nil {
    for _, x : = range s {
     …
    }
}
```

→

```
for _, x : = range s {
 …
}
```

```
if s != nil {
    for _, x : = range s {
      …
    }
}
```

```
for _, x : = range s {
  …
}
```

# Solution: syntactically close templates

```
if :[var] != nil {
  for :[x] : = range :[var] {
      :[body]
  }
}
```

```
for :[x] : = range :[var] {
  :[body]
}
```

# See the link for demo video

[https://drive.google.com/open?id=19X9YL2tZmfOCvK8GxL8OEnUkUB88SC3n](https://drive.google.com/open?id=19X9YL2tZmfOCvK8GxL8OEnUkUB88SC3n)

# Syntax only

```
if :[var] != nil {
   for :[x] : = range :[var] {
       :[body]
   }
}
```

→

```
for :[x] : = range :[var] {
   :[body]
}
```

# Nothing about this is Go specific

```
if :[var] != nil {
  for :[x] : = range :[var] {
      :[body]
  }
}
```

→

```
for :[x] : = range :[var] {
  :[body]
}
```

# Nothing about this is Go specific

## (syntactically)

```
if :[var] != nil {
  for :[x] : = range :[var] {
      :[body]
  }
}
```

→

```
for :[x] : = range :[var] {
  :[body]
}
```

if (:[x].length != 0) → if (:[x].isNotEmpty)

(= :[x] nil)  →  (?nil :[x])

49

.filter(:[x]).size → .count(:[x])

## 34 Answers

active    oldest    **votes**

[1]  [2]  [next]

▲

4420

▼

✓

You can't parse [X]HTML with regex. Because HTML can't be parsed by regex. Regex is not a tool that can be used to correctly parse HTML. As I have answered in HTML-and-regex questions here so many times before, the use of regex will not allow you to consume HTML. Regular expressions are a tool that is insufficiently sophisticated to understand the constructs employed by HTML. HTML is not a regular language and hence cannot be parsed by regular expressions. Regex queries are not equipped to break down HTML into its meaningful parts. so many times but it is not getting to me. Even enhanced irregular regular expressions as used by Perl are not up to the task of parsing HTML. You will never make me crack. HTML is a language of sufficient complexity that it cannot be parsed by regular expressions. Even Jon Skeet cannot parse HTML using regular expressions. Every time you attempt to parse HTML with regular expressions, the unholy child weeps the blood of virgins, and Russian hackers pwn your webapp. Parsing HTML with regex summons tainted souls into the realm of the living. HTML and regex go together like love, marriage, and ritual infanticide. The <center> cannot hold it is too late. The force of regex and HTML together in the same conceptual space will destroy your mind like so much watery putty. If you parse HTML with regex you are giving in to Them and their blasphemous ways which doom us all to inhuman toil for the One whose Name cannot be expressed in the Basic Multilingual Plane, he comes. HTML-plus-regexp will liquify the nerves of the sentient whilst you observe, your psyche withering in the onslaught of horror. Regex-based HTML parsers are the cancer that is killing StackOverflow *it is too late it is too late we cannot be saved* the trangession of a child ensures regex will consume all living tissue (except for HTML which it cannot, as previously prophesied) *dear lord help us how can anyone survive this scourge* using regex to parse HTML has doomed humanity to an eternity of dread torture and security holes *using regex* as a tool to process HTML establishes a brea*ch between this world* and the dread realm of c̃orrupt entities (like SGML entities, but *more corrupt*) *a mere glimps*e of the world of reg**ex parsers for HTML will ins**tantly transport a p*rogrammer's consciousness i*nto a wo*rld* of ceaseless screaming, he comes, the pestilent slithy regex-infection will **devour your** HT̃ML parser, application and existence for all time like Visual̲ Basic only worse *he comes he comes* do not f̧ight h**e com̧es,** ̃h̨is unholy radianćé de*stro̸ying all enlightenment,* HTML tags leak̶ing̡ from̲ your eyes̲͟ ̷like liquid p̲ain, the song of regular expr̶ession parsing̲ will exti̧nguish the voices of mor̶ta̧l man from the sph̨ere I can see it can you see ̷it it is beautiful the f͟inal snuf fing of the lie̶s of Man ALL IS L̶OST̲ ̲A̲LL IS L̲OST the pony he̷ comes he com̶es̶ he comes the ich̷or permeates a̶ll MY FACE M̲Y̲ FACE ̶oh god no NO NOO̷OO NӨ stop the an̸*gĺes ̸are* no̵t rea̵l ZA̲LGO IS̢ ̯TO̵N̡Y̲ THE̅ ̯PONY̡ HE ̱COMES

## 34 Answers

active    oldest    **votes**

| 1 | 2 | next |

▲

4420

▼

✓

You can't parse [X]HTML with regex. Because HTML can't be parsed by regex. Regex is not a tool that can be used to correctly parse HTML. As I have answered in HTML-and-regex questions here so many times before, the use of regex will not allow you to consume HTML. Regular expressions are a tool that is insufficiently sophisticated to understand the constructs employed by HTML. HTML is not a regular language and hence cannot be parsed by regular expressions. Regex queries are not equipped to break down HTML into its meaningful parts. so many times but it is not getting to me. Even enhanced irregular regular expressions as used by Perl are not up to the task of parsing

# Have you tried using an XML parser instead?

late it is too late we cannot be saved the trangession of a child ensures regex will consume all living tissue (except for HTML which it cannot, as previously prophesied) *dear lord help us how can anyone survive this scourge* using regex to parse HTML has doomed humanity to an eternity of dread torture and security holes *using regex* as a tool to process HTML establishes a brea*ch* *between this world* and the dread realm of c̃orrupt entities (like SGML entities, but *more corrupt*) a *mere glimp*se of the world of reg**ex parsers for HTML will ins**tantly transport a p*rogrammer's* *consciousness i*nto a w*orl*d of ceaseless screaming, he comes, the pestilent slithy regex-infection will **devour your** HTML parser, application and existence for all time like Visual Basic only worse *he comes he comes* do not f*i*ght h**e comes, h**is unholy radiańcé de*stroy*ing all enlightenment, HTML tags lea**king from your eyes l**ike liqu*id* p*ain, the song of re*gular expr*ession parsing* will exti*nguish* the voices of mor*ta*l **man from the sph**ere I can see it can you see i͟t it is beautiful the final snuf fing of *the lie***s of Man ALL IS LOST ALL IS L**OST the *pony he comes* he com̲es *he comes* th̲e͟ichor permeates all MY FACE MY FACE ҉h god no NO NOOOO NΘ stop the an*gles are n**ot re̲a̅l ZALGO IS TONY THE PONY HE COMES**

javascript    ruby-on-rails

### Linked

| 29 | Writing regular expression in PHP to wrap <img> with <a> |
| 0 | Regular expression for remove html links |
| 9 | regular expression to remove links |
| 8 | Regexp for html |
| 5 | Regular Expression to remove Div tags |

# A Parser for Multiple Languages

# A parser for multiple languages

- Shared context-free language properties
  - Balanced delimiters
  - Delineate trees

# A parser for multiple languages

- Shared context-free language properties
  - Balanced delimiters
  - Delineate trees

- Take a parenthesis language ➡ extend it

$$S \rightarrow \epsilon \mid SS \mid (\ S\ )$$

# Lightweight Multi-Language Syntax Transformation with Parser Parser Combinators

Rijnard van Tonder
School of Computer Science
Carnegie Mellon University
USA
rvt@cs.cmu.edu

Claire Le Goues
School of Computer Science
Carnegie Mellon University
USA
clegoues@cs.cmu.edu

## Abstract

Automatically transforming programs is hard, yet critical for automated program refactoring, rewriting, and repair. Multi-language syntax transformation is especially hard due to heterogeneous representations in syntax, parse trees, and abstract syntax trees (ASTs). Our insight is that the problem can be decomposed such that (1) a common grammar expresses the central context-free language (CFL) properties shared by many contemporary languages and (2) open extension points in the grammar allow customizing syntax (e.g., for balanced delimiters) and hooks in smaller parsers to handle language-specific syntax (e.g., for comments). Our key contribution operationalizes this decomposition using a Parser Parser combinator (PPC), a mechanism that generates parsers for matching syntactic fragments in source code by parsing declarative user-supplied templates. This allows our approach to detach from translating input pro-

## 1 Introduction

Automatically transforming programs is hard, yet critical for automated program refactoring [1, 2, 45], rewriting [8, 44], and repair [37, 43, 52, 54]. The complexity of automatically transforming code has yielded a plethora of approaches

## 34 Answers

active     oldest     **votes**

1     2     next

▲

**4420**

▼

✓

You can't parse [X]HTML with regex. Because HTML can't be parsed by regex. Regex is not a tool that can be used to correctly parse HTML. As I have answered in HTML-and-regex questions here so many times before, the use of regex will not allow you to consume HTML. Regular expressions are a tool that is insufficiently sophisticated to understand the constructs employed by HTML. HTML is not a regular language and hence cannot be parsed by regular expressions. Regex queries are not equipped to break down HTML into its meaningful parts. so many times but it is not getting to me. Even enhanced irregular regular expressions as used by Perl are not up to the task of parsing HTML. You will never make me crack. HTML is a language of sufficient complexity that it cannot be parsed by regular expressions. Even Jon Skeet cannot parse HTML using regular expressions. Every time you attempt to parse HTML with regular expressions, the unholy child weeps the blood of virgins, and Russian hackers pwn your webapp. Parsing HTML with regex summons tainted souls into the realm of the living. HTML and regex go together like love, marriage, and ritual infanticide. The <center> cannot hold it is too late. The force of regex and HTML together in the same conceptual space will destroy your mind like so much watery putty. If you parse HTML with regex you are giving in to Them and their blasphemous ways which doom us all to inhuman toil for the One whose Name cannot be expressed in the Basic Multilingual Plane, he comes. HTML-plus-regexp will liquify the nerves of the sentient whilst you observe, your psyche withering in the onslaught of horror. Rege̶x-based HTML parsers are the cancer that is killing StackOverflow *it is too late it is too late we cannot be saved* the trangession of a chi̷ld ensures regex will consume all living tissue (except for HTML which it cannot, as previously prophesied) *dear lord help us how can anyone survive this scourge* using regex to parse HTML has doomed humanity to an eternity of dread torture and security holes *using regex* as a tool to process HTML establishes a breach *between this world* and the dread realm of c̕orrupt entities (like SGML entities, but *more corrupt*) a *mere glimp*se of the world of reg**ex parsers for HTML will ins**tantly transport a p*rogrammer's consciousness i*nto a world of ceaseless screaming, he comes, the pestilent slithy regex-infection wil**l devour your H**TML parser, application and existence for all time like Visual Basic only worse *he comes he comes do not fi*ght h**e comps, h**is unho̵ly radia͑nc̕é destro҉ying all enlightenment, HTML tags leͨak̴ing fr̶om ̶your eͮye͠s̸ ͤlike liquid p͠ain, the song of re̸gular expression parsing will exti̧nguish the voices of mortal man from the sphere I can see it can you see ̷it it is beautiful the final snuf fing of the lies of Man ALL IS LOŚT ALL IS LOST the pony he comes he comĕs he comes the ̨ichor permeates all MY FACE MY FAC̯E ͞oh god no NO NOΟOO NΘ stop the an*gͫles ͭare not rea̸l ZA̡LGO IS ͡TONY THE̗ PONY HE ̷COMES

# A simple grammar

grammar ::= term* EOF
term ::= '(' term ')' | '{' term '}' | '[' term ']' | term term | token
token ::= …

# A simple grammar

grammar ::= term* EOF
term ::= '(' term ')' | '{' term '}' | '[' term ']' | term term | token
token ::= anything_else

# Example parse tree for Go code
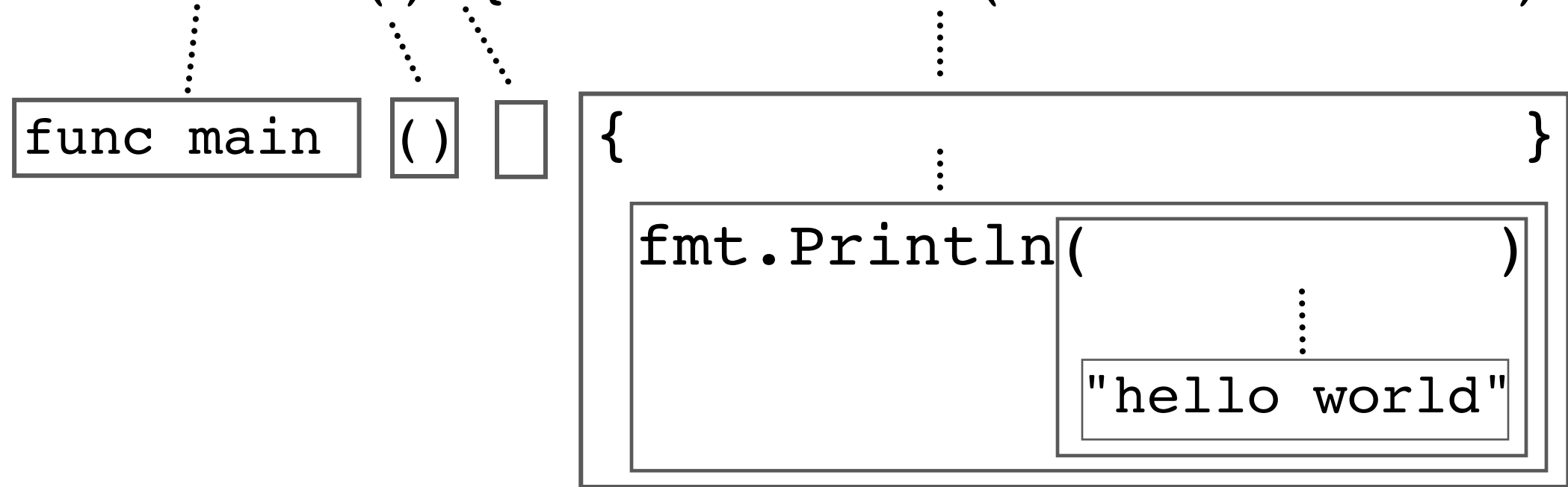
```
func main() { fmt.Println("hello world") }
```

# Example parse tree for Go code

```
func main() { fmt.Println("hello world") }
```
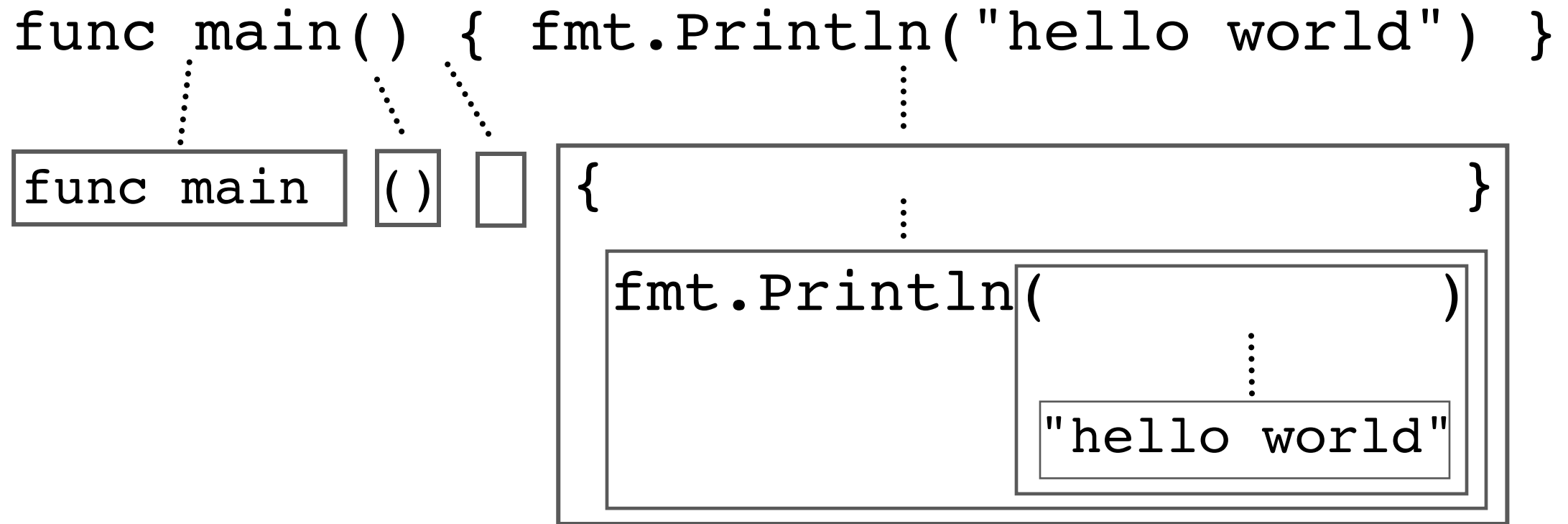
Decompose with respect to delimiters

# Example parse tree for Go code

```
func main() { fmt.Println("hello world") }
```

```
func main    ()        {                                    }
             fmt.Println(                    )
                         "hello world"
```

# Example parse tree for Go code

```
func main() { fmt.Println("hello world") }
```

| func main |   | () |   |   | { ... } |
|-----------|---|----|---|---|---------|

fmt.Println( ... )

"hello world"

grammar ::= term EOF
term ::= '(' term ')' | '{' term '}' | '[' term ']' | term term | token
token ::= anything_else

# How to match
# declarative templates ⬌ source code?

```
if :[var] != nil {
  for :[x] : = range :[var] {
      :[body]
  }
}
```

# Parser Combinators

Model parsers as functions that can be composed using higher-order functions (combinators) to implement grammar constructions.

# Parser Combinators

A parser for an int…

# Parser Combinators

A parser for a string…

# Parser Combinators

A parser for an expression...

# Parser Combinators are polymorphic in their production

A parser for an expression…

# Parser Combinators are polymorphic in their production

A parser  for parsers

# How to match
# declarative templates ⟷ source code?

```
if :[var] != nil {
  for :[x] : = range :[var] {
      :[body]
  }
}
```

# How to match
# declarative templates ⬌ source code?

```
if :[var] != nil {
  for :[x] : = range :[var] {
      :[body]
  }
}
```

## This defines a parser

# How to match
# declarative templates ⟷ source code?

Parse an if

```
if :[var] != nil {
    for :[x] : = range :[var] {
        :[body]
    }
}
```

# This defines a parser

# How to match
# declarative templates ⬌ source code?

Parse an if

```
if :[var] != nil {
    for :[x] : = range :[var] {
        :[body]
    }
}
```

Produce a parser for "if"

## This defines a parser

# How to match
# declarative templates ⬌ source code?

Parse whitespace

```
if :[var] != nil {
    for :[x] : = range :[var] {
        :[body]
    }
}
```

## This defines a parser

# How to match
# declarative templates ⬅➡ source code?

Parse whitespace

```
if :[var] != nil {
    for :[x] : = range :[var] {
        :[body]
    }
}
```

Produce a parser for whitespace

## This defines a parser

# How to match
# declarative templates ⬌ source code?



```
if :[var] != nil {
    for :[x] : = range :[var] {
        :[body]
    }
}
```

Parse a hole

## This defines a parser

# How to match
# declarative templates ⬌ source code?

```
if :[var] != nil {
    for :[x] : = range :[var] {
        :[body]
    }
}
```

Parse a hole

Produce a parser to
match & store text

## This defines a parser

# How to match
# declarative templates ⬌ source code?

Parse balanced {}

```
if :[var] != nil {
    for :[x] : = range :[var] {
        :[body]
    }
}
```

# This defines a parser

# How to match
# declarative templates ⬌ source code?

Parse balanced {}

```
if :[var] != nil {
    for :[x] : = range :[var] {
        :[body]
    }
}
```

Produce a parser
for balanced {}

# This defines a parser

# How to match
# declarative templates ⬌ source code?

```
if :[var] != nil {
    for :[x] : = range :[var] {
        :[body]
    }
}
```

Chain all the parsers…

## This defines a parser

# How to match
# declarative templates ⬌ source code?

```
if :[var] != nil {
    for :[x] : = range :[var] {
        :[body]
    }
}
```

Chain all the parsers…

# ~~Parser~~ Parser Combinators

# Where is the complexity?

```
// 1) More comments more problems
printf(/* arg 1) */ "1) unbalanced \"parens\" (");
```

// 1) More comments more problems
printf(/* arg 1) */ "1) unbalanced \"parens\" (");

```
// 1) More comments more problems
printf(/* arg 1) */ "1) unbalanced \"parens\" (");
```

# See the link for demo video

[https://drive.google.com/open?id=1EGhrBfxw_GQqeH5aWEFLfHSzYBKCjRyz](https://drive.google.com/open?id=1EGhrBfxw_GQqeH5aWEFLfHSzYBKCjRyz)

# Small parsers handle syntax idiosyncracies across languages

# Small parsers handle syntax idiosyncracies across languages

```
let user_defined_delimiters = [ "(" , ")"; "{" , "}"; "[" , "]"]

let string_literals = ["\""; """]

let raw_string_literals = ["`", "`"]

let comment_parser = [ Multiline ("/*" , "*/") ; Until_newline "//"]
```

# Small parsers handle syntax idiosyncracies across languages

```
let user_defined_delimiters = [ "(" , ")"; "{" , "}"; "[" , "]"]

let string_literals = ["\""; """]

let raw_string_literals = ["`", "`"]

let comment_parser = [ Multiline ("/*" , "*/") ; Until_newline "//"]
```

**And embed into a parser skeleton**

# Real world application

# **Large scale application**

- Top 100 GitHub repos for 12 languages
  - 1,200 repos

  Go, Dart, Julia, JS, Rust, Scala, Elm, OCaml, C, Clojure, Erlang, Python

- One to three rewrite rules per language

- 280 million lines of code parsed

- 42 minutes (20 cores)

# **Large scale application**

- Pull requests to 50 unique repositories
  - Merged ~40 PRs

https://catalog.comby.dev/

https://github.com/squaresLab/pldi-artifact-2019/blob/master/PullRequests.md

# Large scale application

- Pull requests to 50 unique repositories
  - Merged ~40 PRs



👁 Watch ▾ | 1,363 ★ Star | 32,159 ⑂ Fork | 5,344

## refactor: use shorthand fields #55734

⑂ Merged   bors merged 1 commit into `rust-lang:master` from `teresy:shorthand-fields`   24 days ago

```
84    84           fn has_type_flags(&self, flags: TypeFlags) -> bool {
85        −            self.visit_with(&mut HasTypeFlagsVisitor { flags: flags })
      85    +            self.visit_with(&mut HasTypeFlagsVisitor { flags })
86    86           }
```

# Demo: end-to-end with nested rewrite

# See the link for demo video

https://drive.google.com/open?id=14Up
dLtYA-2YD71AUDawt_zh0D_SCSZ7C

# Summary

Language Design For Program Manipulation

Eduardus A. T. Merks, J. Michael Dyck, and Robert D. Cameron, *Member, IEEE*

```
func main() { fmt.Println("hello world") }
```

```
func main   ()      {                          }

              fmt.Println(                    )

                           "hello world"
```

grammar ::= term EOF
term ::= '(' term ')' | '{' term '}' | '[' term ']' | term term | token
token ::= anything_else

## How to match
## declarative templates ⟷ source code?

```
if :[var] != nil {
  for :[x] : = range :[var] {
    :[body]
  }
}
```

**This defines a parser**

## Large scale application

• Pull requests to 50 unique repositories
  – Merged ~40 PRs

Watch ⌄  1,363   ★ Star  32,159   Fork  5,344

refactor: use shorthand fields #55734

Merged   bors merged 1 commit into rust-lang:master from teresy:shorthand-fields 24 days ago

```
84   84        fn has_type_flags(&self, flags: TypeFlags) -> bool {
85    -             self.visit_with(&mut HasTypeFlagsVisitor { flags: flags })
     85   +             self.visit_with(&mut HasTypeFlagsVisitor { flags })
86   86        }
```

https://github.com/comby-tools/comby

🐦 @rvtond