# GLL Parsing with Flexible Combinators

L. Thomas van Binsbergen,
Elizabeth Scott, and Adrian Johnstone

Royal Holloway, University of London
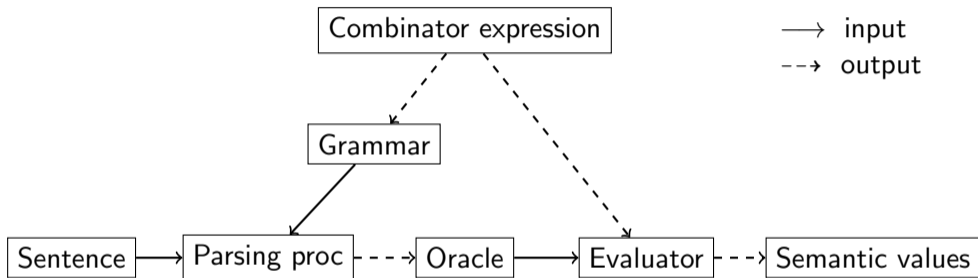ltvanbinsbergen@acm.org

5 November, 2018

http://hackage.haskell.org/package/gll

# Simple, efficient, sound and complete combinator parsing for all context-free grammars, using an oracle

Tom Ridge

University of Leicester

# Library Architecture (Ridge 2014)



- Parser is replaceable
- Similar suggestion by [Ljunglöf, 2002]

## Contributions

1. Functional description and implementation of GLL parsing:
   - All datastructures are basic sets/relations
   - Recursive descent extended to GLL
2. Grammar combinators without grammar binarisation:
   - Combinator expressions evaluate to a grammar object
   - This grammar is an argument to parsing procedure
3. Emperical evaluation on real-world grammars:
   - Demonstrates "acceptable" runtimes on ANSI-C, Caml Light, CBS
   - Significant speed-ups achieved by avoiding binarisation

## 1) Recursive descent parsing

- Every nonterminal is implemented by a *parse function*

- Every parse function has a *branch* for every alternate of the nonterminal

- Every branch is a sequence of:
  - calls to parse functions
  - code matching terminals

## 1) Recursive descent parsing

- Every nonterminal is implemented by a *parse function*

- Every parse function has a *branch* for every alternate of the nonterminal

- Every branch is a sequence of:
  - calls to parse functions
  - code matching terminals

## 2) We consider parse functions that:

- have a parameter holding an index $k$ into the input string (pivot)
- have a local variable remembering the initial pivot value $l$ (left extent)
- return the value $r$ (right extent) held by the parameter at the end of a branch

## 2) We consider parse functions that:

- have a parameter holding an index $k$ into the input string (pivot)
- have a local variable remembering the initial pivot value $l$ (left extent)
- return the value $r$ (right extent) held by the parameter at the end of a branch

## 3) Abstract representation

A *descriptor* $(X ::= \alpha \cdot \beta, l, k)$ models the state of a parse

A *commencement* $(X, l)$ models a (parse) function call

A *continuation* $(X ::= \alpha Y \cdot \beta, l)$ models a return context

# GLL parsing

## 3) Abstract representation

A *descriptor* $(X ::= \alpha \cdot \beta, l, k)$ models the state of a parse

A *commencement* $(X, l)$ models a (parse) function call

A *continuation* $(X ::= \alpha Y \cdot \beta, l)$ models a return context

## 4) Descriptor processing

Process encountered descriptors in any order, exactly once, starting with $(S ::= \cdot \alpha, 0, 0)$

There are three forms of descriptors:

- $(X ::= \alpha \cdot t\beta, l, k)$ with $t$ terminal                                        **match** action
- $(X ::= \alpha \cdot Y\beta, l, k)$ with $Y$ nonterminal             **descend**/**skip** action
- $(Y ::= \delta \cdot, l, r)$                                                      **ascend** action

### 4) Descriptor processing

Process encountered descriptors in any order, exactly once, starting with $(S ::= \cdot\alpha, 0, 0)$

There are three forms of descriptors:

- $(X ::= \alpha \cdot t\beta, l, k)$ with $t$ terminal                            **match** action
- $(X ::= \alpha \cdot Y\beta, l, k)$ with $Y$ nonterminal         **descend**/**skip** action
- $(Y ::= \delta\cdot, l, r)$                                         **ascend** action

### 5) GLL datatypes

The set $\mathcal{U}$ contains all descriptors processed so far

The relation $\mathcal{P}$ pairs commencements with right extents

The relation $\mathcal{G}$ pairs commencements with continuations

## 5) GLL datatypes

The set $\mathcal{U}$ contains all descriptors processed so far

The relation $\mathcal{P}$ pairs commencements with right extents

The relation $\mathcal{G}$ pairs commencements with continuations

## match

$$(X ::= \alpha \cdot t\beta, l, k) \quad \rightarrow \quad (X ::= \alpha t \cdot \beta, l, k+1)$$

*pre-conditions*:

- $t$ is the $k$'th terminal in the input string

*post-conditions*:

- $(X ::= \alpha \cdot t\beta, l, k) \in \mathcal{U}$

**descend**

$$(X ::= \alpha \cdot Y\beta, l, k) \quad \rightarrow \quad (Y ::= \cdot\delta_1, k, k)$$
$$\ldots$$
$$\rightarrow \quad (Y ::= \cdot\delta_i, k, k)$$

*pre-conditions*:

- $Y ::= \delta_i$ is in the grammar, for all $i$
- There is no $r$ such that $((Y, k), r) \in \mathcal{P}$

*post-conditions*:

- Possible new continuation: $((Y, k), (X ::= \alpha Y \cdot \beta, l)) \in \mathcal{G}$
- $(X ::= \alpha \cdot Y\beta, l, k) \in \mathcal{U}$

# GLL parsing

## skip

$$(X ::= \alpha \cdot Y\beta, l, k) \quad \rightarrow \quad (X ::= \alpha Y \cdot \beta, l, r_1)$$
$$\cdots$$
$$\rightarrow \quad (X ::= \alpha Y \cdot \beta, l, r_j)$$

*pre-conditions*:

- For all $1 \leqslant i \leqslant j$, we have $((Y, k), r_i) \in \mathcal{P}$ (at least one)

*post-conditions*:

- Possible new continuation: $((Y, k), (X ::= \alpha Y \cdot \beta, l)) \in \mathcal{G}$
- $(X ::= \alpha \cdot Y\beta, l, k) \in \mathcal{U}$

**ascend**

$$(Y ::= \delta \cdot, l, r) \quad \rightarrow \quad (X ::= \alpha_1 Y \cdot \beta_1, l_1, r)$$
$$\cdots$$
$$\rightarrow \quad (X ::= \alpha_i Y \cdot \beta_i, l_j, r)$$

*pre-conditions*:

- For all $1 \leqslant i \leqslant j$, we have $((Y, l), (X ::= \alpha_i Y \cdot \beta_i, l_i)) \in \mathcal{G}$

*post-conditions*:

- Possible new right extent: $((Y, l), r) \in \mathcal{P}$
- $(Y ::= \delta \cdot, l, r) \in \mathcal{U}$

$$(X ::= \alpha \cdot s\beta, l, k) \in \mathcal{U} \quad \& \quad (X ::= \alpha s \cdot \beta, l, r) \in \mathcal{U}$$

gives

$$(X ::= \alpha s \cdot \beta, l, k, r) \in \mathcal{O}$$

$$(Y ::= \delta \cdot, l, r) \quad \textbf{with } l = r, \delta = \epsilon$$

gives

$$(Y ::= \delta \cdot, l, l, l) \in \mathcal{O}$$

## Contributions

1. Functional description and implementation of GLL parsing:
   - All datastructures are basic sets/relations
   - Recursive descent extended to GLL
2. Grammar combinators without grammar binarisation:
   - Combinator expressions evaluate to a grammar object
   - This grammar is an argument to parsing procedure
3. Emperical evaluation on real-world grammars:
   - Demonstrates "acceptable" runtimes on ANSI-C, Caml Light, CBS
   - Significant speed-ups achieved by avoiding binarisation

## Parser combinators

$term$     $:: Eq\ t \Rightarrow t \rightarrow Parser$
$epsilon :: Parser$
$(\langle * \rangle)$   $:: Parser \rightarrow Parser \rightarrow Parser$
$(\langle | \rangle)$   $:: Parser \rightarrow Parser \rightarrow Parser$

## Example    $T ::= (\ A\ )$   $A ::= \epsilon\ |\ M\ a$   $M ::= \epsilon\ |\ M\ a$,

$pT = term$ '(' $\langle * \rangle\ pA\ \langle * \rangle\ term$ ')'
$pA = epsilon\ \langle | \rangle\ pM\ \langle * \rangle\ term$ 'a'
$pM = epsilon\ \langle | \rangle\ pM\ \langle * \rangle\ term$ 'a' $\langle * \rangle\ term$ ','

## Grammar combinators

$$
\begin{array}{lll}
term & :: Eq\ t \Rightarrow t & \rightarrow Grammar \\
epsilon & :: Grammar & \\
(\langle * \rangle) & :: Grammar \rightarrow Grammar \rightarrow Grammar \\
(\langle | \rangle) & :: Grammar \rightarrow Grammar \rightarrow Grammar
\end{array}
$$

## Grammar extraction

- Expressions yield at most two productions with at most two symbols in rhs

$nt(x) = $ "(" $+\!\!\!+ \ nt(l) \ +\!\!\!+ $ "*" $+\!\!\!+ \ nt(r) \ +\!\!\!+ $ ")"        **if** $x = l \ \langle * \rangle \ r$

$nt(y) = $ "(" $+\!\!\!+ \ nt(p) \ +\!\!\!+ $ "|" $+\!\!\!+ \ nt(q) \ +\!\!\!+ $ ")"        **if** $y = p \ \langle | \rangle \ q$

productions: $nt(x) ::= nt(l)nt(r)$, $nt(y) ::= nt(p)$, $nt(y) ::= nt(q)$

## Grammar combinators

$$nterm :: String \rightarrow Grammar$$
$$term :: Eq\ t \Rightarrow t \rightarrow Grammar$$
$$epsilon :: Grammar$$
$$(\langle * \rangle) :: Grammar \rightarrow Grammar \rightarrow Grammar$$
$$(\langle | \rangle) :: Grammar \rightarrow Grammar \rightarrow Grammar$$

## Grammar extraction

• Expressions yield at most two productions with at most two symbols in rhs

$$nt(x) = \text{"("} + nt(l) + \text{"*"} + nt(r) + \text{")"} \qquad \textbf{if } x = l \langle * \rangle r$$

$$nt(y) = \text{"("} + nt(p) + \text{"|"} + nt(q) + \text{")"} \qquad \textbf{if } y = p \langle | \rangle q$$

productions: $nt(x) ::= nt(l)nt(r)$, $nt(y) ::= nt(p)$, $nt(y) ::= nt(q)$

## BNF combinators

$$( \langle ::= \rangle ) \ :: \ String \rightarrow Choice_{\text{EX}} \rightarrow Symb_{\text{EX}}$$

$$term \quad :: \ Eq \ t \ \Rightarrow t \qquad \rightarrow Symb_{\text{EX}}$$

$$( \langle ** \rangle ) \ :: \ Seq_{\text{EX}} \rightarrow Symb_{\text{EX}} \rightarrow Seq_{\text{EX}}$$

$$seqStart :: Seq_{\text{EX}}$$

$$( \langle || \rangle ) \quad :: \ Choice_{\text{EX}} \rightarrow Seq_{\text{EX}} \rightarrow Choice_{\text{EX}}$$

$$altStart :: Choice_{\text{EX}}$$

## Example $T ::= ( \ A \ )$

$$gT = \texttt{"T"} \ \langle ::= \rangle \ altStart \ \langle || \rangle \ seqStart \ \langle ** \rangle \ term \ \texttt{'('} \ \langle ** \rangle \ gA \ \langle ** \rangle \ term \ \texttt{')'}$$

$$gA = ...$$

## Flexible BNF combinators

$(IsSeq\ seq, IsCh\ ch, IsSymb\ symb) \Rightarrow$

| | | |
|---|---|---|
| $(\langle ::= \rangle)$ | $:: String \rightarrow ch \rightarrow Symb_{\text{EX}}$ |
| $term$ | $:: Eq\ t \Rightarrow t \rightarrow Symb_{\text{EX}}$ |
| $(\langle ** \rangle)$ | $:: seq \rightarrow symb \rightarrow Seq_{\text{EX}}$ |
| $seqStart$ | $:: Seq_{\text{EX}}$ |
| $(\langle || \rangle)$ | $:: ch \rightarrow seq \rightarrow Choice_{\text{EX}}$ |
| $altStart$ | $:: Choice_{\text{EX}}$ |

## Example $T ::= (\ A\ )$

$gT = $ "T" $\langle ::= \rangle\ term\ $ '(' $\langle ** \rangle\ gA\ \langle ** \rangle\ term\ $ ')'
$gA = ...$

## Conversions



|  |  |  |  |
|---|---|---|---|
| **instance** *IsSeq Seq*$_{\text{EX}}$ | **where** ... | -- id |  |
| **instance** *IsSeq Symb*$_{\text{EX}}$ | **where** ... | -- (a) |  |
| **instance** *IsSeq Choice*$_{\text{EX}}$ | **where** ... | -- (a) ∘ (c) |  |
| **instance** *IsCh Choice*$_{\text{EX}}$ | **where** ... | -- id |  |
| **instance** *IsCh Seq*$_{\text{EX}}$ | **where** ... | -- (b) |  |
| **instance** *IsCh Symb*$_{\text{EX}}$ | **where** ... | -- (b) ∘ (a) |  |
| **instance** *IsSymb Symb*$_{\text{EX}}$ | **where** ... | -- id |  |
| **instance** *IsSymb Choice*$_{\text{EX}}$ | **where** ... | -- (c) |  |
| **instance** *IsSymb Seq*$_{\text{EX}}$ | **where** ... | -- (c) ∘ (b) |  |

## Contributions

1. Functional description and implementation of GLL parsing:
   - All datastructures are basic sets/relations
   - Recursive descent extended to GLL
2. Grammar combinators without grammar binarisation:
   - Combinator expressions evaluate to a grammar object
   - This grammar is an argument to parsing procedure
3. Emperical evaluation on real-world grammars:
   - Demonstrates "acceptable" runtimes on ANSI-C, Caml Light, CBS
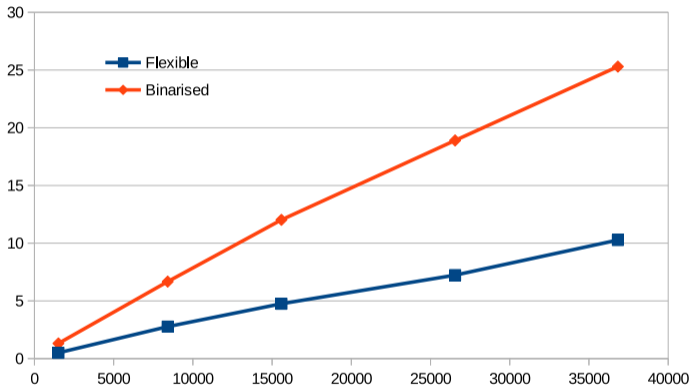   - Significant speed-ups achieved by avoiding binarisation

## Claims

- Running times show that the approach is practical
    *Although the emphasis has been on correctness*

- Avoiding binarisation improves running times
    *Syntax descriptions have not been manipulated to benefit evaluation*

$$( \ \langle ::= \rangle_{\text{BIN}} \ ) :: String \rightarrow Symb_{\text{EX}} \rightarrow Symb_{\text{EX}}$$

$$( \ \langle ::= \rangle_{\text{BIN}} \ ) = ( \ \langle ::= \rangle \ )$$

$$( \ \langle || \rangle_{\text{BIN}} \ ) :: Symb_{\text{EX}} \rightarrow Symb_{\text{EX}} \rightarrow Symb_{\text{EX}}$$

$$p \ \langle || \rangle_{\text{BIN}} \ q = toSymb \ (p \ \langle || \rangle \ q)$$

$$( \ \langle ** \rangle_{\text{BIN}} \ ) :: Symb_{\text{EX}} \rightarrow Symb_{\text{EX}} \rightarrow Symb_{\text{EX}}$$

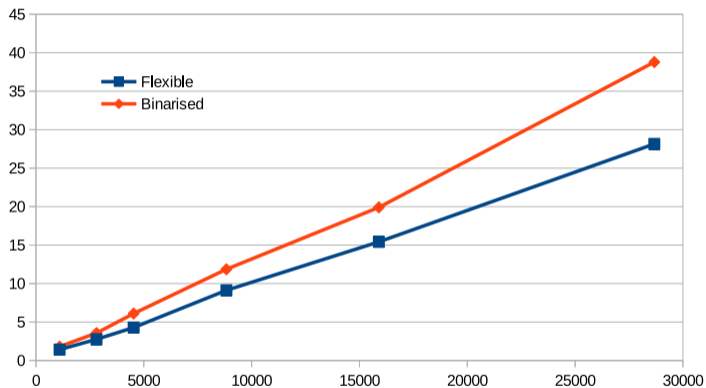$$p \ \langle ** \rangle_{\text{BIN}} \ q = toSymb \ (p \ \langle ** \rangle \ q)$$

Binarised: 690 nonterminals, 848 alternates
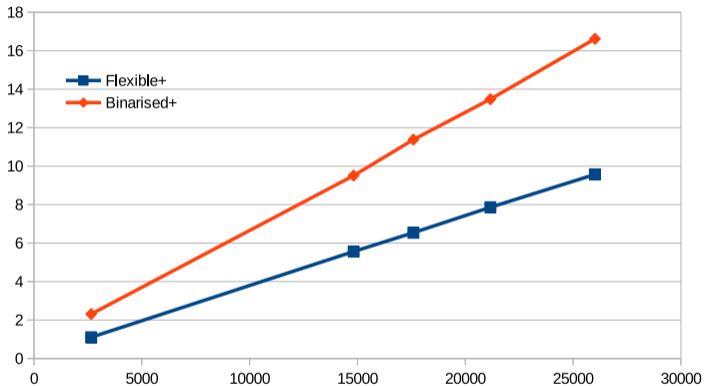Flexible: 71 nonterminals, 229 alternates          2.4–2.6x speed-up (with lookahead)

Binarised: 580 nonterminals, 731 alternates

Flexible: 134 nonterminals, 285 alternates          1.3-1.4x speed-up (with lookahead)

Binarised: 640 nonterminals, 771 alternates
Flexible: 126 nonterminals, 257 alternates          1.7-2.1x speed-up (with lookahead)

- An EDSL for describing context-free grammars based on 'BNF combinators'

- Parsers with on-the-fly semantics available for described grammars

- Generalised parsing certainly simplifies SLE

- Library suitable for our purpose: reference interpreters for programming languages

- Caveats:
  - Disambiguation mostly ad-hoc
  - Manual nonterminal insertion problematic in some cases

# GLL Parsing with Flexible Combinators

L. Thomas van Binsbergen,
Elizabeth Scott, and Adrian Johnstone

Royal Holloway, University of London
ltvanbinsbergen@acm.org

5 November, 2018

```
http://hackage.haskell.org/package/gll
http://ltvanbinsbergen.nl/thesis
```

| | |
|---|---|
| Parser Combinators | `parsec`<br>`UU-lib`<br>... |
| Explicit Nonterminals | Scheme recognisers (Johnson 1995)<br>`Meerkat` (Izmaylova/Afroozeh 2015/16) |
| Grammar Combinators | P3 (Ridge 2014)<br>GLL.Combinators (2015/16)<br>`grammar-combinators` (Devriese 2011/12) |
| Meta-Programming | BNFC-meta (Duregard 2011) |
| Parser Generators | `Bison`<br>`yacc`<br>`Happy`<br>... |