

# Generalised parsing with Binary Subtree Representations (BSR)

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

Royal Holloway, University of London  
ltvanbinsbergen@acm.org

24 April, 2019

Research output: SLE 2018, SCP 2019 vol. 175  
Hackage packages: g11, happy (in development)  
Thesis: <http://ltvanbinsbergen.nl/thesis>

## Past research

- Alternative descriptions of GLL parsing:
  - Purely functional GLL parsing
  - GLL parsing in parser combinators

All developed parser algorithms output binary subtree representation (BSR) sets

## Future work

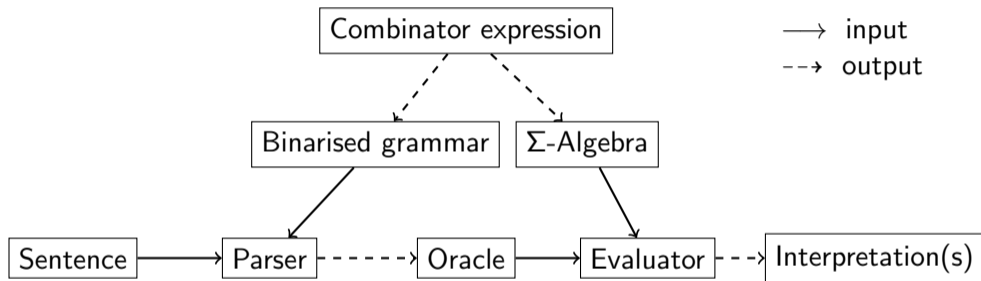
- Generating GLL parsers from:
  - Parameterised grammar descriptions (Happy)
  - Generating GLL parsers from CBS syntax descriptions
- Expressing disambiguation strategies on BSR sets

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

Tom Ridge

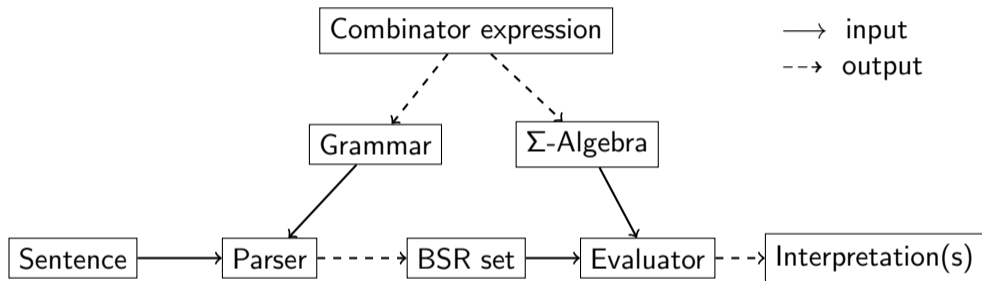
University of Leicester

# Library architecture (Ridge 2014)



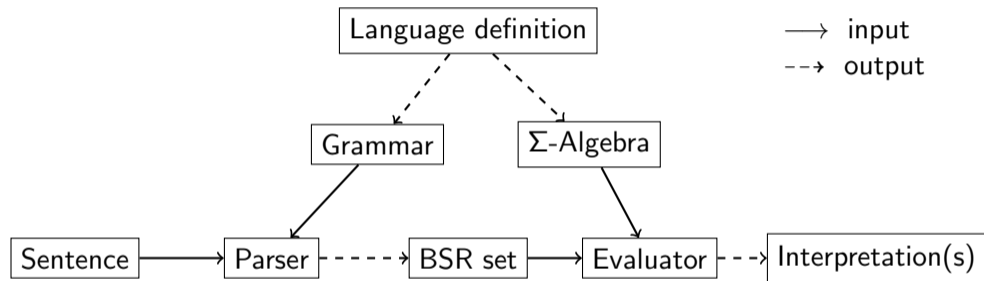
- Parser is replaceable
- Similar suggestion by [Ljunglöf, 2002]

# Library architecture (SLE 2018)



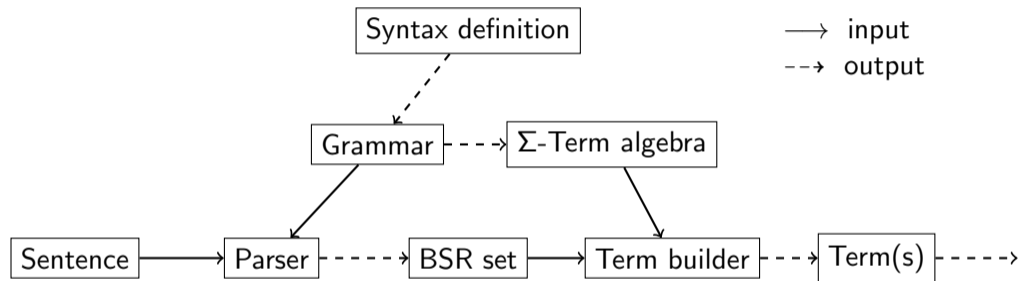
- Parser is replaceable
- Binary Subtree Representation (BSR) sets for arbitrary CFGs

# General architecture



- Parser is replaceable
- Binary Subtree Representation (BSR) sets for arbitrary CFGs

# Term construction



- Parser is replaceable
- The term builder might construct terms that share common sub-terms

*Let parsers produce traces that contain sufficient information to reproduce derivations and delay term construction or interpretation to a separate phase*

## Benefits

- Separation of concerns
- Simplifies definition of the parser
- Parser improvements/optimisations do not influence subsequent phases



Grammar:  $\Gamma = \langle N, T, P, Z \rangle$

Input sentence:  $t_0 \dots t_{m-1}$

Grammar slots:  $G(\Gamma) = \{X \rightarrow \alpha \bullet \beta \mid X \rightarrow \alpha \beta \in P\}$

Items:  $\mathcal{I}(\Gamma, t_0 \dots t_{m-1}) = \{\langle g, l, k \rangle \mid g \in G(\Gamma), 0 \leq l \leq k \leq m\}$

Valid items:  $\mathcal{U}(\Gamma, t_0 \dots t_{m-1}) = \dots$

## Parsing schemata (Sikkel 1998)

Valid items are inferable according to some inference relation

The inference relation captures a parsing algorithm conceptually (not operationally)

# Generalised top-down parsing (canonical Earley)

$$\frac{Z \rightarrow \delta \in P}{\langle Z \rightarrow \bullet \delta, 0, 0 \rangle \in \mathcal{U}} \quad \text{(init)}$$

$$\frac{\langle X \rightarrow \alpha \bullet t \beta, l, k \rangle \in \mathcal{U} \quad t = t_k}{\langle X \rightarrow \alpha t \bullet \beta, l, k + 1 \rangle \in \mathcal{U}} \quad \text{(scanner)}$$

$$\frac{\langle X \rightarrow \alpha \bullet Y \beta, l, k \rangle \in \mathcal{U} \quad Y \rightarrow \delta \in p}{\langle Y \rightarrow \bullet \delta, k, k \rangle \in \mathcal{U}} \quad \text{(predictor)}$$

$$\frac{\langle X \rightarrow \alpha \bullet Y \beta, l, k \rangle \in \mathcal{U} \quad \langle Y \rightarrow \delta \bullet, k, r \rangle \in \mathcal{U}}{\langle X \rightarrow \alpha Y \bullet \beta, l, r \rangle \in \mathcal{U}} \quad \text{(completer)}$$

$$\langle X \rightarrow \alpha \bullet s \beta, l, k \rangle \in \mathcal{U} \quad \& \quad \langle X \rightarrow \alpha s \bullet \beta, l, r \rangle \in \mathcal{U}$$

gives

$$\langle X \rightarrow \alpha s \bullet \beta, l, k, r \rangle \in \Upsilon$$

$$\langle Y \rightarrow \bullet, l, l \rangle \in \mathcal{U}$$

gives

$$\langle Y \rightarrow \bullet, l, l, l \rangle \in \Upsilon$$

```
program chart parser
begin
  create initial chart and agenda;
  while agenda is not empty
  do
    delete some (arbitrarily chosen) current item from agenda
    add current to chart;
    for each item that can be recognized by current in
    combination with other items in chart
    do
      if item is neither in chart nor in agenda
      then add item to agenda fi
    od
  od
end.
```

The chart parser algorithm.

(Sikkel 1998).

# Implementing canonical Earley

Earley's algorithm: process items in ascending order of  $k$  (by executing *scanner* last)

The *completer* on grammars without nullable nonterminals

$$\frac{\langle X \rightarrow \alpha \bullet Y \beta, l, k \rangle \in \mathcal{U} \quad \langle Y \rightarrow \delta \bullet, k, r \rangle \in \mathcal{U}}{\langle X \rightarrow \alpha Y \bullet \beta, l, r \rangle \in \mathcal{U}} \quad (\text{completer})$$

When processing  $\langle X \rightarrow \alpha \bullet Y \beta, l, k \rangle$ , no items  $\langle Y \rightarrow \delta \bullet, k, r \rangle$  will be discovered yet

When processing  $\langle Y \rightarrow \delta \bullet, k, r \rangle$ , all  $\langle X \rightarrow \alpha \bullet Y \beta, l, k \rangle$  are discovered when  $k = r$

GLL parsing: process items in any order with special data structures that are used instead of 'searching through'  $\mathcal{U}$

GLL introduces a GSS and a pop-set for both directions of the *completer* action

$$\frac{\langle Y \rightarrow \delta \bullet, k, r \rangle \in \mathcal{U} \quad \langle X \rightarrow \alpha Y \bullet \beta, l \rangle \in \text{gss}(Y, k)}{\langle X \rightarrow \alpha Y \bullet \beta, l, r \rangle \in \mathcal{U} \quad [r \in \text{pop}(Y, k)]} \quad \text{(completer (1)/ascend)}$$

$$\frac{\langle X \rightarrow \alpha \bullet Y \beta, l, k \rangle \in \mathcal{U} \quad r \in \text{pop}(Y, k)}{\langle X \rightarrow \alpha Y \bullet \beta, l, r \rangle \in \mathcal{U} \quad [\langle X \rightarrow \alpha Y \bullet \beta, l \rangle \in \text{gss}(Y, k)]} \quad \text{(completer (2)/skip)}$$

$$\frac{\langle X \rightarrow \alpha \bullet Y \beta, l, k \rangle \in \mathcal{U} \quad Y \rightarrow \delta \in p \quad \emptyset = \text{pop}(Y, k)}{\langle Y \rightarrow \bullet \delta, k, k \rangle \in \mathcal{U} \quad [\langle X \rightarrow \alpha Y \bullet \beta, l \rangle \in \text{gss}(Y, k)]} \quad \text{(predictor/descend)}$$

GLL introduces a GSS and a pop-set for both directions of the *completer* action

$$\frac{\langle Y \rightarrow \delta \bullet, k, r \rangle \in \mathcal{U}}{\langle X \rightarrow \alpha Y \bullet \beta, l, r \rangle \in \mathcal{U}} \quad \text{(completer (1)/ascend)}$$

$$\frac{\langle X \rightarrow \alpha \bullet Y \beta, l, k \rangle \in \mathcal{U}}{\langle X \rightarrow \alpha Y \bullet \beta, l, r \rangle \in \mathcal{U}} \quad \text{(completer (2)/skip)}$$

$$\frac{\langle X \rightarrow \alpha \bullet Y \beta, l, k \rangle \in \mathcal{U} \quad Y \rightarrow \delta \in p}{\langle Y \rightarrow \bullet \delta, k, k \rangle \in \mathcal{U}} \quad \text{(predictor/descend)}$$

GLL introduces a GSS and a pop-set for both directions of the *completer* action

$$\frac{\langle Y \rightarrow \delta \bullet, k, r \rangle \in \mathcal{U} \quad \langle X \rightarrow \alpha Y \bullet \beta, l \rangle \in \text{gss}(Y, k)}{\langle X \rightarrow \alpha Y \bullet \beta, l, r \rangle \in \mathcal{U}} \quad (\text{completer (1)/ascend})$$

$$\frac{\langle X \rightarrow \alpha \bullet Y \beta, l, k \rangle \in \mathcal{U}}{\langle X \rightarrow \alpha Y \bullet \beta, l, r \rangle \in \mathcal{U}} \quad (\text{completer (2)/skip})$$

$$\frac{\langle X \rightarrow \alpha \bullet Y \beta, l, k \rangle \in \mathcal{U} \quad Y \rightarrow \delta \in p}{\langle Y \rightarrow \bullet \delta, k, k \rangle \in \mathcal{U}} \quad (\text{predictor/descend})$$



GLL introduces a GSS and a pop-set for both directions of the *completer* action

$$\frac{\langle Y \rightarrow \delta \bullet, k, r \rangle \in \mathcal{U} \quad \langle X \rightarrow \alpha Y \bullet \beta, l \rangle \in \text{gss}(Y, k)}{\langle X \rightarrow \alpha Y \bullet \beta, l, r \rangle \in \mathcal{U} \quad [r \in \text{pop}(Y, k)]} \quad \text{(completer (1)/ascend)}$$

$$\frac{\langle X \rightarrow \alpha \bullet Y \beta, l, k \rangle \in \mathcal{U}}{\langle X \rightarrow \alpha Y \bullet \beta, l, r \rangle \in \mathcal{U}} \quad \text{(completer (2)/skip)}$$

$$\frac{\langle X \rightarrow \alpha \bullet Y \beta, l, k \rangle \in \mathcal{U} \quad Y \rightarrow \delta \in p}{\langle Y \rightarrow \bullet \delta, k, k \rangle \in \mathcal{U}} \quad \text{(predictor/descend)}$$

GLL introduces a GSS and a pop-set for both directions of the *completer* action

$$\frac{\langle Y \rightarrow \delta \bullet, k, r \rangle \in \mathcal{U} \quad \langle X \rightarrow \alpha Y \bullet \beta, l \rangle \in \text{gss}(Y, k)}{\langle X \rightarrow \alpha Y \bullet \beta, l, r \rangle \in \mathcal{U} \quad [r \in \text{pop}(Y, k)]} \quad \text{(completer (1)/ascend)}$$

$$\frac{\langle X \rightarrow \alpha \bullet Y \beta, l, k \rangle \in \mathcal{U} \quad r \in \text{pop}(Y, k)}{\langle X \rightarrow \alpha Y \bullet \beta, l, r \rangle \in \mathcal{U}} \quad \text{(completer (2)/skip)}$$

$$\frac{\langle X \rightarrow \alpha \bullet Y \beta, l, k \rangle \in \mathcal{U} \quad Y \rightarrow \delta \in p}{\langle Y \rightarrow \bullet \delta, k, k \rangle \in \mathcal{U}} \quad \text{(predictor/descend)}$$

GLL introduces a GSS and a pop-set for both directions of the *completer* action

$$\frac{\langle Y \rightarrow \delta \bullet, k, r \rangle \in \mathcal{U} \quad \langle X \rightarrow \alpha Y \bullet \beta, l \rangle \in \text{gss}(Y, k)}{\langle X \rightarrow \alpha Y \bullet \beta, l, r \rangle \in \mathcal{U} \quad [r \in \text{pop}(Y, k)]} \quad \text{(completer (1)/ascend)}$$

$$\frac{\langle X \rightarrow \alpha \bullet Y \beta, l, k \rangle \in \mathcal{U} \quad r \in \text{pop}(Y, k)}{\langle X \rightarrow \alpha Y \bullet \beta, l, r \rangle \in \mathcal{U} \quad [\langle X \rightarrow \alpha Y \bullet \beta, l \rangle \in \text{gss}(Y, k)]} \quad \text{(completer (2)/skip)}$$

$$\frac{\langle X \rightarrow \alpha \bullet Y \beta, l, k \rangle \in \mathcal{U} \quad Y \rightarrow \delta \in p}{\langle Y \rightarrow \bullet \delta, k, k \rangle \in \mathcal{U}} \quad \text{(predictor/descend)}$$

GLL introduces a GSS and a pop-set for both directions of the *completer* action

$$\frac{\langle Y \rightarrow \delta \bullet, k, r \rangle \in \mathcal{U} \quad \langle X \rightarrow \alpha Y \bullet \beta, l \rangle \in \text{gss}(Y, k)}{\langle X \rightarrow \alpha Y \bullet \beta, l, r \rangle \in \mathcal{U} \quad [r \in \text{pop}(Y, k)]} \quad \text{(completer (1)/ascend)}$$

$$\frac{\langle X \rightarrow \alpha \bullet Y \beta, l, k \rangle \in \mathcal{U} \quad r \in \text{pop}(Y, k)}{\langle X \rightarrow \alpha Y \bullet \beta, l, r \rangle \in \mathcal{U} \quad [\langle X \rightarrow \alpha Y \bullet \beta, l \rangle \in \text{gss}(Y, k)]} \quad \text{(completer (2)/skip)}$$

$$\frac{\langle X \rightarrow \alpha \bullet Y \beta, l, k \rangle \in \mathcal{U} \quad Y \rightarrow \delta \in p \quad \emptyset = \text{pop}(Y, k)}{\langle Y \rightarrow \bullet \delta, k, k \rangle \in \mathcal{U}} \quad \text{(predictor/descend)}$$

GLL introduces a GSS and a pop-set for both directions of the *completer* action

$$\frac{\langle Y \rightarrow \delta \bullet, k, r \rangle \in \mathcal{U} \quad \langle X \rightarrow \alpha Y \bullet \beta, l \rangle \in \text{gss}(Y, k)}{\langle X \rightarrow \alpha Y \bullet \beta, l, r \rangle \in \mathcal{U} \quad [r \in \text{pop}(Y, k)]} \quad \text{(completer (1)/ascend)}$$

$$\frac{\langle X \rightarrow \alpha \bullet Y \beta, l, k \rangle \in \mathcal{U} \quad r \in \text{pop}(Y, k)}{\langle X \rightarrow \alpha Y \bullet \beta, l, r \rangle \in \mathcal{U} \quad [\langle X \rightarrow \alpha Y \bullet \beta, l \rangle \in \text{gss}(Y, k)]} \quad \text{(completer (2)/skip)}$$

$$\frac{\langle X \rightarrow \alpha \bullet Y \beta, l, k \rangle \in \mathcal{U} \quad Y \rightarrow \delta \in p \quad \emptyset = \text{pop}(Y, k)}{\langle Y \rightarrow \bullet \delta, k, k \rangle \in \mathcal{U} \quad [\langle X \rightarrow \alpha Y \bullet \beta, l \rangle \in \text{gss}(Y, k)]} \quad \text{(predictor/descend)}$$

## GLL parsing summary

- Items can be processed through a worklist in any order
- For each item being processed, there is at most one applicable rule  
Possible actions: *descend*, *ascend*, *skip*, *ascend*
- If a rule is applicable for multiple instantiations, apply it for all instantiations 'simultaneously'

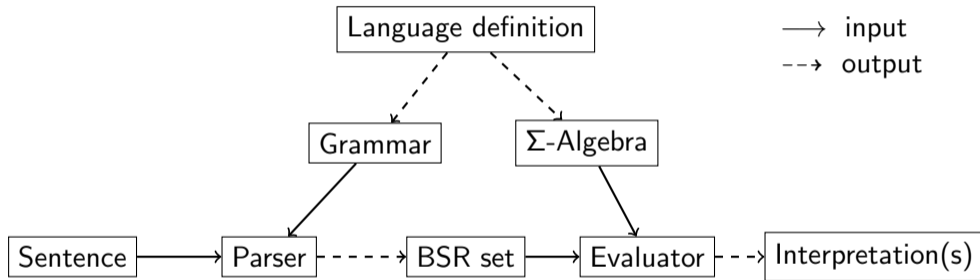
## Correctness

Termination	There is a finite number of items: $O(m^2)$
Soundness	Derived from canonical Earley
Completeness	Derived from: <i>skip</i> + <i>ascend</i> = <i>completer</i>

## Summary - Part 1

- Any parsing algorithm that specialises canonical Earley can yield BSR sets
- GLL parsing specialises Earley with data structures for *completer*

# General architecture (recap)





## Definition

A *language definition* consists of a CFG  $\Gamma = \langle N, T, P, Z \rangle$ , a family of sets  $\{\tau_s\}_{s \in N \cup T}$  with  $\tau_t = \{t\}$  for  $t \in T$ , and a family of higher-order functions  $\{\mathcal{A}_p\}_{p \in P}$  with

$$\mathcal{A}_{s_{n+1} \rightarrow s_1 \dots s_n} : \tau_{s_1} \rightarrow \dots \rightarrow \tau_{s_n} \rightarrow \tau_{s_{n+1}} \quad (\text{type also denoted } \tau_{s_1 \dots s_{n+1}})$$

## Definition

A *language definition* consists of a CFG  $\Gamma = \langle N, T, P, Z \rangle$ , a family of sets  $\{\tau_s\}_{s \in N \cup T}$  with  $\tau_t = \{t\}$  for  $t \in T$ , and a family of higher-order functions  $\{\mathcal{A}_p\}_{p \in P}$  with

$$\mathcal{A}_{s_{n+1} \rightarrow s_1 \dots s_n} : \tau_{s_1} \rightarrow \dots \rightarrow \tau_{s_n} \rightarrow \tau_{s_{n+1}} \quad (\text{type also denoted } \tau_{s_1 \dots s_{n+1}})$$

## Partial application

The family  $\mathcal{A}$  is extended with a function for each slot  $g \in G(\Gamma)$

$$\mathcal{A}_{X \rightarrow \bullet \beta} : \tau_{\beta X}$$

$$\mathcal{A}_{X \rightarrow \bullet \beta} = \mathcal{A}_{X \rightarrow \beta}$$

$$\mathcal{A}_{X \rightarrow \alpha s \bullet \beta} : (\tau_{s \beta X} \times \tau_s) \rightarrow \tau_{\beta X}$$

$$\mathcal{A}_{X \rightarrow \alpha s \bullet \beta}(f, v) = f(v)$$

## Canonical BSR-based Evaluation

For each  $s \in N \cup T$ , a function  $\mathcal{E}_s : \hat{\Upsilon} \times \mathbb{N} \times \mathbb{N} \rightarrow \mathcal{P}(\tau_s)$  is defined and for each  $g \in G(\Gamma)$  with  $g = X \rightarrow \alpha \bullet \beta$  a function  $\mathcal{E}_g : \hat{\Upsilon} \times \mathbb{N} \times \mathbb{N} \rightarrow \mathcal{P}(\tau_{\beta X})$

For each  $s \in N \cup T$ , a function  $\mathcal{E}_s : \hat{\Upsilon} \times \mathbb{N} \times \mathbb{N} \rightarrow \mathcal{P}(\tau_s)$  is defined and for each  $g \in G(\Gamma)$  with  $g = X \rightarrow \alpha \bullet \beta$  a function  $\mathcal{E}_g : \hat{\Upsilon} \times \mathbb{N} \times \mathbb{N} \rightarrow \mathcal{P}(\tau_{\beta X})$

$$\mathcal{E}_t(\Upsilon, l, r) = \{t \mid r = l + 1\}$$

## Canonical BSR-based Evaluation

For each  $s \in N \cup T$ , a function  $\mathcal{E}_s : \hat{\Upsilon} \times \mathbb{N} \times \mathbb{N} \rightarrow \mathcal{P}(\tau_s)$  is defined and for each  $g \in G(\Gamma)$  with  $g = X \rightarrow \alpha \bullet \beta$  a function  $\mathcal{E}_g : \hat{\Upsilon} \times \mathbb{N} \times \mathbb{N} \rightarrow \mathcal{P}(\tau_{\beta X})$

$$\mathcal{E}_t(\Upsilon, l, r) = \{t \mid r = l + 1\}$$

$$\mathcal{E}_X(\Upsilon, l, r) = \{v \mid X \rightarrow \beta \in P, v \in \mathcal{E}_{X \rightarrow \beta \bullet}(\Upsilon, l, r)\}$$

## Canonical BSR-based Evaluation

For each  $s \in N \cup T$ , a function  $\mathcal{E}_s : \hat{\Upsilon} \times \mathbb{N} \times \mathbb{N} \rightarrow \mathcal{P}(\tau_s)$  is defined and for each  $g \in G(\Gamma)$  with  $g = X \rightarrow \alpha \bullet \beta$  a function  $\mathcal{E}_g : \hat{\Upsilon} \times \mathbb{N} \times \mathbb{N} \rightarrow \mathcal{P}(\tau_{\beta X})$

$$\mathcal{E}_t(\Upsilon, l, r) = \{t \mid r = l + 1\}$$

$$\mathcal{E}_X(\Upsilon, l, r) = \{v \mid X \rightarrow \beta \in P, v \in \mathcal{E}_{X \rightarrow \beta \bullet}(\Upsilon, l, r)\}$$

$$\mathcal{E}_{X \rightarrow \bullet \beta}(\Upsilon, l, r) = \{\mathcal{A}_{X \rightarrow \bullet \beta} \mid l = r\}$$

# Canonical BSR-based Evaluation

For each  $s \in N \cup T$ , a function  $\mathcal{E}_s : \hat{\Upsilon} \times \mathbb{N} \times \mathbb{N} \rightarrow \mathcal{P}(\tau_s)$  is defined and for each  $g \in G(\Gamma)$  with  $g = X \rightarrow \alpha \bullet \beta$  a function  $\mathcal{E}_g : \hat{\Upsilon} \times \mathbb{N} \times \mathbb{N} \rightarrow \mathcal{P}(\tau_{\beta X})$

$$\mathcal{E}_t(\Upsilon, l, r) = \{t \mid r = l + 1\}$$

$$\mathcal{E}_X(\Upsilon, l, r) = \{v \mid X \rightarrow \beta \in P, v \in \mathcal{E}_{X \rightarrow \beta \bullet}(\Upsilon, l, r)\}$$

$$\mathcal{E}_{X \rightarrow \bullet \beta}(\Upsilon, l, r) = \{\mathcal{A}_{X \rightarrow \bullet \beta} \mid l = r\}$$

$$\begin{aligned} \mathcal{E}_{X \rightarrow \alpha s \bullet \beta}(\Upsilon, l, r) = \{ & \mathcal{A}_{X \rightarrow \alpha s \bullet \beta}(f, v) \mid \langle X \rightarrow \alpha s \bullet \beta, l, k, r \rangle \in \Upsilon \\ & , v \in \mathcal{E}_s(\Upsilon, k, r) \\ & , f \in \mathcal{E}_{X \rightarrow \alpha \bullet s \beta}(\Upsilon, l, k)\} \end{aligned}$$

# Term construction (without sharing)

## Terms

Induced by a CFG is a family of sets of terms  $\{\mathcal{T}_s\}_{s \in N \cup T}$  where  $\mathcal{T}_t = \{t\}$  (for  $t \in T$ ) and  $\mathcal{T}_X$  (for  $X \in N$ ) is defined as the smallest set such that:

$$\frac{X \rightarrow s_1 \dots s_n \in P \quad t_1 \in \mathcal{T}_{s_1} \dots t_n \in \mathcal{T}_{s_n}}{\langle X, t_1 \dots t_n \rangle \in \mathcal{T}_X}$$

## Term construction

Given a CFG, form a language definition by taking  $\tau_s = \mathcal{T}_s$  (for all  $s \in N \cup T$ ) and by defining  $\mathcal{A}_p$  as follows (for all  $p \in P$  with  $p = X \rightarrow s_1 \dots s_n$ )

$$\mathcal{A}_{X \rightarrow s_1 \dots s_n}(t_1) \dots (t_n) = \langle X, t_1 \dots t_n \rangle$$



# Evaluator modifications for ambiguity reduction

- 1 General filtering of nonterminal results
- 2 Removing 'invalid' interpretations (invalid according to some static semantics)  
...
- 3 Removing 'non-minimal' terms (or their interpretations) for cyclic grammars
- 4 Associativity of operators
- 5 Precedence of operators
- 6 Follow-restrictions
- 7 Longest/shortest-match

## 2. Removing invalid interpretations

Let  $\tau_s^o = \{\mathbf{none}\} \cup \{\mathbf{some}(v) \mid v \in \tau_s\}$  and let  $\tau_{s_1 \dots s_n}^o$  denote  $\tau_{s_1} \rightarrow \dots \rightarrow \tau_{s_{n-1}} \rightarrow \tau_{s_n}^o$

### Language definitions

The actions  $\mathcal{A}_p$  of a language definition may now fail, as is reflected in the types:

$$\mathcal{A}_{s_{n+1} \rightarrow s_1 \dots s_n} : \tau_{s_1 \dots s_{n+1}}^o$$

$$\mathcal{A}_{s_{n+1} \rightarrow \bullet s_1 \dots s_n} : \tau_{s_1 \dots s_{n+1}}^o$$

$$\mathcal{A}_{s_{n+1} \rightarrow \alpha s_0 \bullet s_1 \dots s_n} : (\tau_{s_0 \dots s_{n+1}}^o \times \tau_{s_0}) \rightarrow \tau_{s_1 \dots s_{n+1}}^o$$

Update the definition of  $\mathcal{E}_X$ :

$$\mathcal{E}_X(\Upsilon, l, r) = \{v \mid \langle X \rightarrow \beta \bullet, l, k, r \rangle \in \Upsilon, \mathbf{some}(v) \in \mathcal{E}_{X \rightarrow \beta \bullet}(\Upsilon, l, r)\}$$

Arithmetic expression example

# 1. Filtering nonterminal results

## Language definitions

Let language definitions provide a family of functions  $\{\mathcal{F}_X\}_{X \in N}$  with  $\mathcal{F}_X : \mathcal{P}(\tau_X) \rightarrow \mathcal{P}(\tau_X)$  s.t. for all  $\phi_X \subseteq \mathcal{P}(\tau_X)$  it holds that  $\mathcal{F}_X(\phi_X) \subseteq \phi_X$ .

Update the definition of  $\mathcal{E}_X$ :

$$\mathcal{E}_X(\Upsilon, l, r) = \mathcal{F}_X(\{v \mid \langle X \rightarrow \beta \bullet, l, k, r \rangle \in \Upsilon, \text{some}(v) \in \mathcal{E}_{X \rightarrow \beta \bullet}(\Upsilon, l, r)\})$$

Overloading example

A grammar is cyclic if it has a nonterminal that can derive itself via one or more steps

## Definition

A term is non-minimal if a subterm  $t$  can be replaced by a term  $t'$  for which holds that  $t$  and  $t'$  have the same label and yield, and  $t'$  is a subterm of  $t$ .

- A sentence that admits a non-minimal term has infinitely many interpretations
- Cyclic grammars are the only grammars that admit non-minimal terms

$$E ::= EEE \mid "1" \mid \epsilon$$

# Cyclic grammars - Example

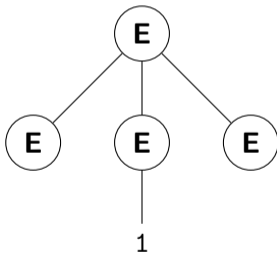
$E ::= EEE \mid "1" \mid \epsilon$





# Cyclic grammars - Example

$E ::= EEE \mid "1" \mid \epsilon$



### 3. Removing non-minimal terms/interpretations

Remember the nonterminal labels of ancestor nodes with the same extents in  $\mathcal{C}$

$$\mathcal{E}_t(\Upsilon, l, r, \mathcal{C}) = \{t \mid r = l + 1\}$$

$$\mathcal{E}_X(\Upsilon, l, r, \mathcal{C}) = \{v \mid X \notin \mathcal{C}, X \rightarrow \beta \in P, v \in \mathcal{E}_{X \rightarrow \beta \bullet}(\Upsilon, l, r, \mathcal{C} \cup \{X\})\}$$

$$\mathcal{E}_{X \rightarrow \bullet \beta}(\Upsilon, l, r, \mathcal{C}) = \{\mathcal{A}_{X \rightarrow \bullet \beta} \mid l = r\}$$

$$\begin{aligned} \mathcal{E}_{X \rightarrow \alpha s \bullet \beta}(\Upsilon, l, r, \mathcal{C}) = \{ & \mathcal{A}_{X \rightarrow \alpha s \bullet \beta}(f, v) \mid \langle X \rightarrow \alpha s \bullet \beta, l, k, r \rangle \in \Upsilon \\ & , v \in \mathcal{E}_s(\Upsilon, k, r, \{Y \mid k = l, Y \in \mathcal{C}\}) \\ & , f \in \mathcal{E}_{X \rightarrow \alpha \bullet s \beta}(\Upsilon, l, k, \{Y \mid k = r, Y \in \mathcal{C}\})\} \end{aligned}$$

# Evaluator modifications for ambiguity reduction

- 1 General filtering of nonterminal results
- 2 Removing 'invalid' interpretations (invalid according to some static semantics)  
...
- 3 Removing 'non-minimal' terms (or their interpretations) for cyclic grammars
- 4 Associativity of operators
- 5 Precedence of operators
- 6 Follow-restrictions
- 7 Longest/shortest-match

- The binary subtree representation can be the basis of a separate evaluation phase
- General and specific disambiguation strategies can manifest themselves in the evaluation phase

### Research Questions

- Is it feasible to disambiguate with the more general strategies only?
- How to avoid hard-coding strategies into the evaluation algorithm?
- With what sort of technique could we describe strategies at BSR element level, without interfering (much) with the evaluator?

Scales example

# Generalised parsing with Binary Subtree Representations (BSR)

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

Royal Holloway, University of London  
ltvanbinsbergen@acm.org

24 April, 2019

Research output: SLE 2018, SCP 2019 vol. 175  
Hackage packages: g11, happy (in development)  
Thesis: <http://ltvanbinsbergen.nl/thesis>

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