

Disambiguating Grammars with Tree Automata¹

Michael D. Adams Matthew Might

University of Utah

Parsing@SLE
November 1, 2015

¹This material is based in part upon work supported by the National Science Foundation under Grant Number 1248464 and the Defense Advanced Research Projects Agency under agreement no. AFRL FA8750-15-2-0092. Any opinions, findings, and conclusions or recommendations expressed in this material are those of the author(s) and do not necessarily reflect the views of the National Science Foundation.

Disambiguation Problems

- Precedence and Associativity
- Dangling else
- ML `if`
- JavaScript constructors

- Precedence and Associativity
 - Domain specific annotations
 - Doesn't cover C:
 - Invalid: ++ (int) x
 - Valid: * (int*) x
 - Valid: ++ * (int*) x
- Dangling else
- ML if
- JavaScript constructors

- Precedence and Associativity
 - Domain specific annotations
 - Doesn't cover C:
 - Invalid: ++ (int) x
 - Valid: * (int*) x
 - Valid: ++ * (int*) x
- Dangling else
 - Shift/reduce bias
- ML if
- JavaScript constructors

- Precedence and Associativity
 - Domain specific annotations
 - Doesn't cover C:
 - Invalid: ++ (int) x
 - Valid: * (int*) x
 - Valid: ++ * (int*) x
- Dangling else
 - Shift/reduce bias
- ML if
 - SDF annotations
- JavaScript constructors

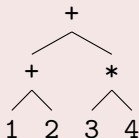
- Precedence and Associativity
 - Domain specific annotations
 - Doesn't cover C:
 - Invalid: ++ (int) x
 - Valid: * (int*) x
 - Valid: ++ * (int*) x
- Dangling else
 - Shift/reduce bias
- ML if
 - SDF annotations
- JavaScript constructors
 - Complicated grammar

Disambiguation Problems

- Precedence and Associativity
 - ~~Domain-specific annotations~~ Tree automata
 - ~~Doesn't cover C:~~ Tree automata cover C:
 - Invalid: ++ (int) x
 - Valid: * (int*) x
 - Valid: ++ * (int*) x
- Dangling else
 - ~~Shift/reduce bias~~ Tree automata
- ML if
 - ~~SDF-annotations~~ Tree automata
- JavaScript constructors
 - ~~Complicated grammar~~ Tree automata

Basic Idea

- String automata: 1234
- Tree automata:



Properties

- Closed under union, intersection and negation
- Never adds ambiguities
- LR(k) closed
- Conjecture: Produces optimal grammars

1+2+3*4

Grammar

Exp \rightarrow nat

Exp \rightarrow Exp "+" Exp

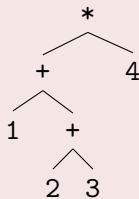
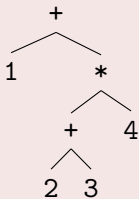
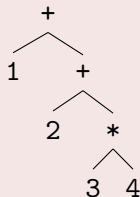
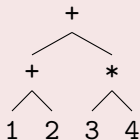
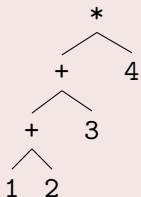
Exp \rightarrow Exp "*" Exp

Exp \rightarrow "(" Exp ")"

Tree Automata

1+2+3*4

Trees



1+2+3*4

Grammar

Exp \rightarrow nat

Exp \rightarrow Exp "+" Exp

Exp \rightarrow Exp "*" Exp

Exp \rightarrow "(" Exp ")"

1+2+3*4

Grammar

Exp $\rightarrow_{\mathbb{N}}$ nat

Exp \rightarrow_{+} Exp "+" Exp

Exp \rightarrow_{*} Exp "*" Exp

Exp $\rightarrow_{()}$ "(" Exp ")"

1+2+3*4

Grammar

$\text{Exp} \rightarrow_{\mathbb{N}} \epsilon$

$\text{Exp} \rightarrow + \text{Exp} \text{Exp}$

$\text{Exp} \rightarrow * \text{Exp} \text{Exp}$

$\text{Exp} \rightarrow (\text{Exp}$

1+2+3*4

Grammar

$\text{Exp} \rightarrow_{\mathbb{N}} \epsilon$

$\text{Exp} \rightarrow + \text{Exp} \text{Exp}$

$\text{Exp} \rightarrow * \text{Exp} \text{Exp}$

$\text{Exp} \rightarrow (\text{Exp})$

Exp

Tree Automata

1+2+3*4

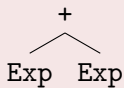
Grammar

$\text{Exp} \rightarrow_{\mathbb{N}} \epsilon$

$\text{Exp} \rightarrow + \text{Exp} \text{Exp}$

$\text{Exp} \rightarrow * \text{Exp} \text{Exp}$

$\text{Exp} \rightarrow (\text{Exp}$



Tree Automata

1+2+3*4

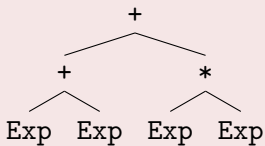
Grammar

$\text{Exp} \rightarrow_{\mathbb{N}} \epsilon$

$\text{Exp} \rightarrow + \text{Exp} \text{Exp}$

$\text{Exp} \rightarrow * \text{Exp} \text{Exp}$

$\text{Exp} \rightarrow (\text{Exp})$



Tree Automata

1+2+3*4

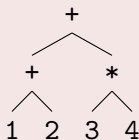
Grammar

$\text{Exp} \rightarrow_{\mathbb{N}} \epsilon$

$\text{Exp} \rightarrow + \text{Exp} \text{Exp}$

$\text{Exp} \rightarrow * \text{Exp} \text{Exp}$

$\text{Exp} \rightarrow (\text{Exp}$



Tree Automata

1+2+3*4

Grammar

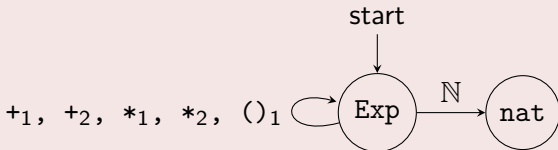
$\text{Exp} \rightarrow_{\mathbb{N}} \epsilon$

$\text{Exp} \rightarrow + \text{Exp} \text{Exp}$

$\text{Exp} \rightarrow * \text{Exp} \text{Exp}$

$\text{Exp} \rightarrow () \text{Exp}$

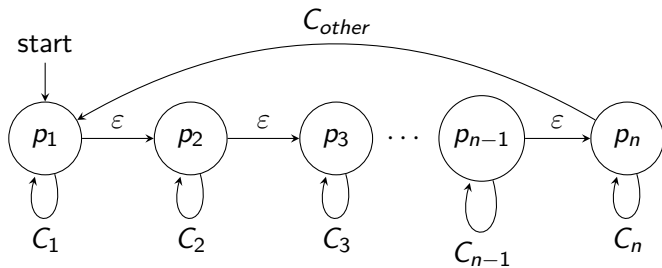
Automaton



Regular Expression

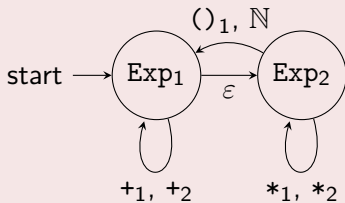
$(+(\square, \square) \mid *(\square, \square) \mid () (\square))^* \cdot \mathbb{N}$

Example: Precedence



Example: Precedence

Automaton



Grammar

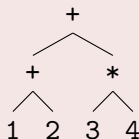
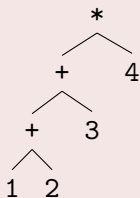
$Exp_1 \rightarrow + Exp_1 Exp_1$

$Exp_1 \rightarrow \epsilon Exp_2$

$Exp_2 \rightarrow N \epsilon$

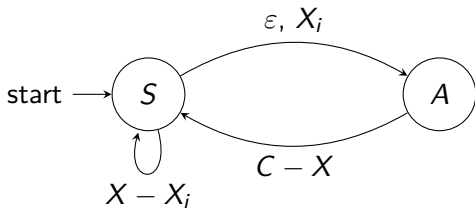
$Exp_2 \rightarrow * Exp_2 Exp_2$

$Exp_2 \rightarrow () Exp_1$

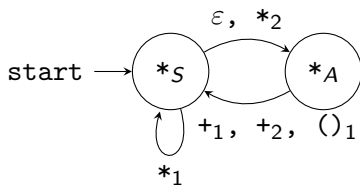
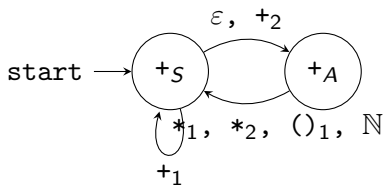


Example: Associativity

Forbidden



Example: Associativity



$\text{Exp} \rightarrow_+ \text{Exp Exp}^+$

$\text{Exp} \rightarrow_\epsilon \text{Exp}^+$

$\text{Exp}^+ \rightarrow_N \epsilon$

$\text{Exp}^+ \rightarrow_* \text{Exp Exp}$

$\text{Exp}^+ \rightarrow_{()} \text{Exp}$

$\text{Exp} \rightarrow_* \text{Exp Exp}^*$

$\text{Exp} \rightarrow_\epsilon \text{Exp}^*$

$\text{Exp}^* \rightarrow_N$

$\text{Exp}^* \rightarrow_+ \text{Exp Exp}$

$\text{Exp}^* \rightarrow_{()} \text{Exp}$

Example: Combined grammar

Combined

Exp ₁	→ _N	ε
Exp ₁	→ ₊	Exp ₁ Exp ₁ ₊
Exp ₁	→ _*	Exp ₂ Exp ₂ _*
Exp ₁	→ _()	Exp ₁
Exp ₁ ₊	→ _N	ε
Exp ₁ ₊	→ _*	Exp ₂ Exp ₂ ₊
Exp ₁ ₊	→ _()	Exp ₁
Exp ₂	→ _N	ε
Exp ₂	→ _*	Exp ₂ Exp ₂ _*
Exp ₂	→ _()	Exp ₁
Exp ₂ _*	→ _N	ε
Exp ₂ _*	→ _()	Exp ₁

Simplified

Exp	→ ₊	Exp Term
Exp	→ _ε	Term
Term	→ _*	Term Factor
Term	→ _ε	Factor
Factor	→ _N	ε
Factor	→ _()	Exp

Grammar

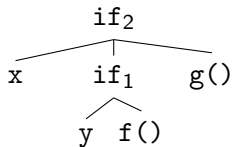
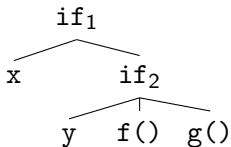
```
Exp  $\rightarrow_{\mathbb{N}}$  nat  
Exp  $\rightarrow_{+}$  Exp "+" Exp  
Exp  $\rightarrow_{*}$  Exp "*" Exp  
Exp  $\rightarrow_{()}$  "(" Exp ")"
```

Restrictions

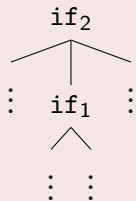
```
Precedence: ((+(□, □))*  
             . *(□, □))*  
             . (nat | ()(□))**)*  
Assoc. for +: not EF(+(_, +( _, _)))  
Assoc. for *: not EF(*(_, *( _, _)))
```


Example: Dangling else

```
if (x) if (y) f(); else g();
```



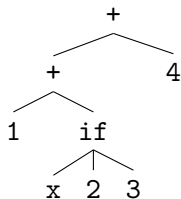
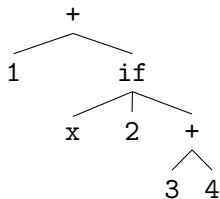
Forbidden



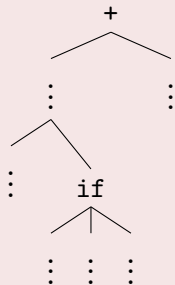
```
not (EF (if2(_, □, _)
        . if1(_, _)))
```

Example: ML's if

1 + if x then 2 else 3 + 4



Forbidden



not (EF (+(\square , $_$)

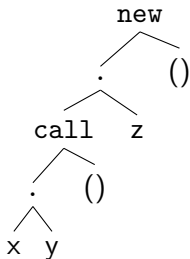
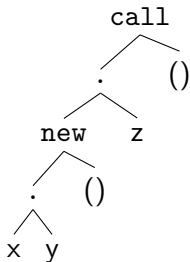
. R^*

. if($_$, $_$, $_$)))

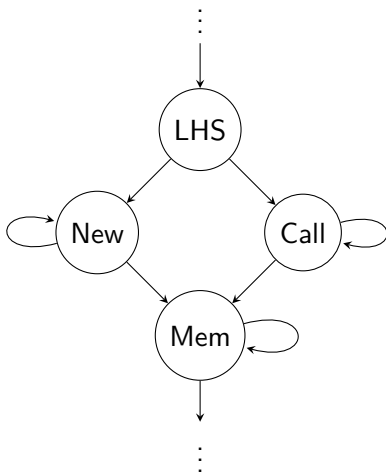
where $R = +(_$, \square) | $-(_$, \square) | \dots

Example: JavaScript

`new x.y().z()`

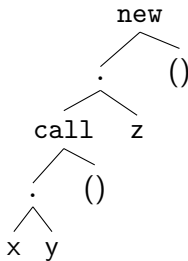
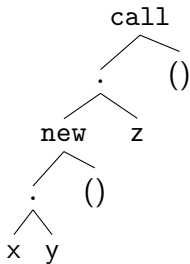


Example: JavaScript



Tree Automata

```
not (EF (New(□, _)
        . (Field(□, _) | Array(□, _))*
        . Call(_, _)))
```



Theorem

Intersecting a TA with a CFG does not introduce ambiguities.

Theorem

LR(k) is closed under intersection with a TA.

Conjecture

The minimization and bottom-up determinization of a CFG results in the fewest LR(k) states of any CFG producing the same parse trees.

Using TA to restrict grammars

- Write your grammar the “natural” way
- Encode grammatical restrictions as TA
- Intersect the grammar with the TA
- Feed the result to any parsing system you want

Open Question

Does minimization and determinization produce an optimal LR parser?

