

# On LR Parsing with Selective Delays

E. Bertsch

Ruhr-U. Bochum



M.-J. Nederhof

U. St Andrews



S. Schmitz

ENS Cachan



CC 2013, March 22nd

Rome, Italy

# SELECTIVE ML

NAMED AFTER **MARCUS (1980)** AND **LEERMAKERS (1992)**

Conflicts

Grammar transformation

Parser construction

Results

# C++ QUALIFIED IDENTIFIERS

(ISO STANDARD, 1998)

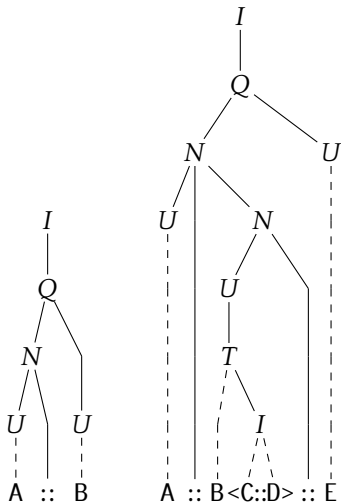
Grammar:

$$I \rightarrow U \mid Q$$

$$U \rightarrow i \mid T$$

$$Q \rightarrow NU$$

$$N \rightarrow U :: N \mid U ::$$

$$T \rightarrow i \langle I \rangle$$


# C++ QUALIFIED IDENTIFIERS

(ISO STANDARD, 1998)

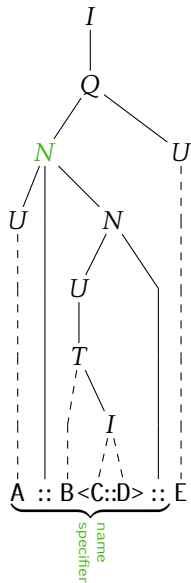
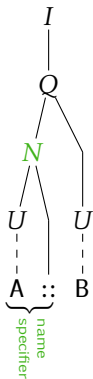
Grammar:

$$I \rightarrow U \mid Q$$

$$U \rightarrow i \mid T$$

$$Q \rightarrow NU$$

$$N \rightarrow U :: N \mid U ::$$

$$T \rightarrow i \langle I \rangle$$


# C++ QUALIFIED IDENTIFIERS

(ISO STANDARD, 1998)

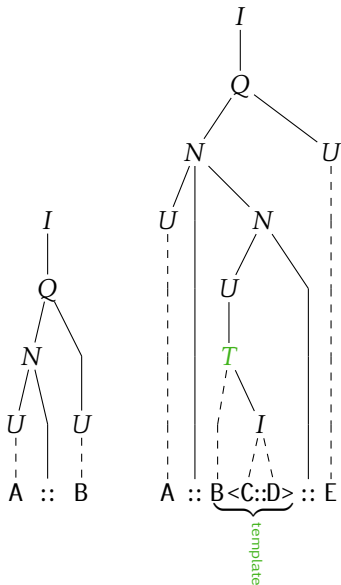
Grammar:

$$I \rightarrow U \mid Q$$

$$U \rightarrow i \mid T$$

$$Q \rightarrow NU$$

$$N \rightarrow U :: N \mid U ::$$

$$T \rightarrow i \langle I \rangle$$


# C++ QUALIFIED IDENTIFIERS

(ISO STANDARD, 1998)

Grammar:

$$I \rightarrow U \mid Q$$
$$U \rightarrow i \mid T$$
$$Q \rightarrow NU$$
$$N \rightarrow U :: N \mid U ::$$
$$T \rightarrow i \langle I \rangle$$

GNU/Bison output:

- ▶ conflicts: 1 shift/reduce
- ▶ warning: rule useless in parser due to conflicts:

$$N \rightarrow U ::$$

# C++ QUALIFIED IDENTIFIERS

(ISO STANDARD, 1998)

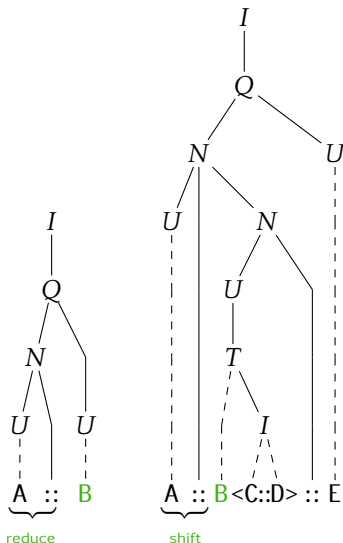
Grammar:

$$I \rightarrow U \mid Q$$

$$U \rightarrow i \mid T$$

$$Q \rightarrow NU$$

$$N \rightarrow U :: N \mid U ::$$

$$T \rightarrow i \langle I \rangle$$


# C++ QUALIFIED IDENTIFIERS

(ISO STANDARD, 1998)

Grammar:

$$I \rightarrow U | Q$$

$$U \rightarrow i | T$$

$$Q \rightarrow NU$$

$$N \rightarrow U :: N | U ::$$

$$T \rightarrow i \langle I \rangle$$

## Solutions?

- ▶ priority mechanisms  
(Bouwers et al., 2008)
- ▶ GLR
- ▶ refactor the grammar
  - ▶ language-preserving
  - ▶ large space of possible transformations  
(see e.g. Lämmel, 2001)



# C++ QUALIFIED IDENTIFIERS

(ISO STANDARD, 1998)

Grammar:

$$I \rightarrow U \mid Q$$
$$U \rightarrow i \mid T$$
$$Q \rightarrow NU$$
$$N \rightarrow U :: N \mid U ::$$
$$T \rightarrow i \langle I \rangle$$

## Solutions?

- ▶ **priority mechanisms**  
(Bouwers et al., 2008)
- ▶ GLR
- ▶ refactor the grammar
  - ▶ language-preserving
  - ▶ large space of possible transformations  
(see e.g. Lämmel, 2001)

# C++ QUALIFIED IDENTIFIERS

(ISO STANDARD, 1998)

Grammar:

$$I \rightarrow U \mid Q$$
$$U \rightarrow i \mid T$$
$$Q \rightarrow NU$$
$$N \rightarrow U :: N \mid U ::$$
$$T \rightarrow i \langle I \rangle$$

## Solutions?

- ▶ **priority mechanisms**  
(Bouwers et al., 2008)
- ▶ **GLR**
- ▶ refactor the grammar
  - ▶ language-preserving
  - ▶ large space of possible transformations  
(see e.g. Lämmel, 2001)

# C++ QUALIFIED IDENTIFIERS

(ISO STANDARD, 1998)

Grammar:

$$I \rightarrow U | Q$$

$$U \rightarrow i | T$$

$$Q \rightarrow NU$$

$$N \rightarrow U :: N | U ::$$

$$T \rightarrow i \langle I \rangle$$

## Solutions?

- ▶ priority mechanisms  
(Bouwers et al., 2008)
- ▶ GLR
- ▶ refactor the grammar
  - ▶ language-preserving
  - ▶ large space of possible transformations  
(see e.g. Lämmel, 2001)

# C++ QUALIFIED IDENTIFIERS

(ISO STANDARD, 1998)

left-recursive Grammar:

$$I \rightarrow U \mid Q$$

$$U \rightarrow i \mid T$$

$$Q \rightarrow NU$$

$$N \rightarrow NU :: \mid U ::$$

$$T \rightarrow i \langle I \rangle$$

## Solutions?

- ▶ priority mechanisms  
(Bouwers et al., 2008)
- ▶ GLR
- ▶ refactor the grammar
  - ▶ language-preserving
  - ▶ large space of possible transformations  
(see e.g. Lämmel, 2001)

# C++ QUALIFIED IDENTIFIERS

(ISO STANDARD, 1998)

combed Grammar:

$$I \rightarrow U \mid Q$$

$$U \rightarrow i \mid T$$

$$Q \rightarrow [NU]$$

$$N \rightarrow U :: N \mid U ::$$

$$T \rightarrow i \langle I \rangle$$

## Solutions?

- ▶ priority mechanisms  
(Bouwers et al., 2008)
- ▶ GLR
- ▶ refactor the grammar
  - ▶ language-preserving
  - ▶ large space of possible transformations  
(see e.g. Lämmel, 2001)

# C++ QUALIFIED IDENTIFIERS

(ISO STANDARD, 1998)

combed Grammar:

$$I \rightarrow U \mid Q$$

$$U \rightarrow i \mid T$$

$$Q \rightarrow [NU]$$

$$[NU] \rightarrow U :: [NU] \mid U :: U$$

$$T \rightarrow i \langle I \rangle$$

## Solutions?

- ▶ priority mechanisms  
(Bouwers et al., 2008)
- ▶ GLR
- ▶ refactor the grammar
  - ▶ language-preserving
  - ▶ large space of possible transformations  
(see e.g. Lämmel, 2001)

# C++ QUALIFIED IDENTIFIERS

(ISO STANDARD, 1998)

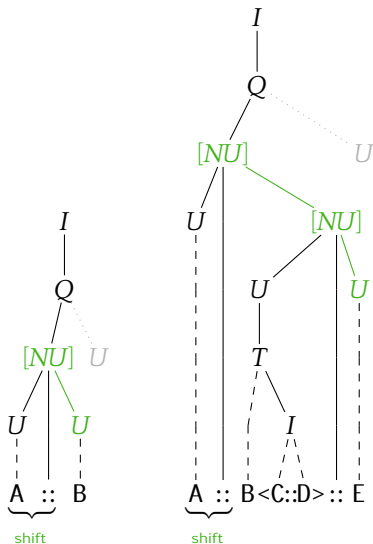
combined Grammar:

$$I \rightarrow U \mid Q$$

$$U \rightarrow i \mid T$$

$$Q \rightarrow [NU]$$

$$[NU] \rightarrow U :: [NU] \mid U :: U$$

$$T \rightarrow i \langle I \rangle$$


# GOING FORMAL: COMBININGS

Write a context-free grammar as  $G = \langle N, \Sigma, P, S \rangle$ , with vocabulary  $V = \Sigma \uplus N$ .

$G \hookrightarrow G'$  if  $\exists \mu: V'^* \rightarrow V^*$  a homomorphism s.t.

1.  $\mu(S') = S$ ,
2.  $\forall a \in \Sigma, \mu(a) = a$ ,
3.  $\mu(N') \subseteq N \cdot V^*$
4.  $\{A' \rightarrow \mu(\alpha') \mid A' \rightarrow \alpha' \in P'\}$   
 $= \{A' \rightarrow \alpha\delta \mid A' \in N', \mu(A') = A\delta, \text{ and } A \rightarrow \alpha \in P\}$

We call it a  **$k$ -combing** ( $\hookrightarrow_k$ ) if  $\mu(N') \subseteq N \cdot V^{\leq k}$ .



# GOING FORMAL: COMBININGS

## Example

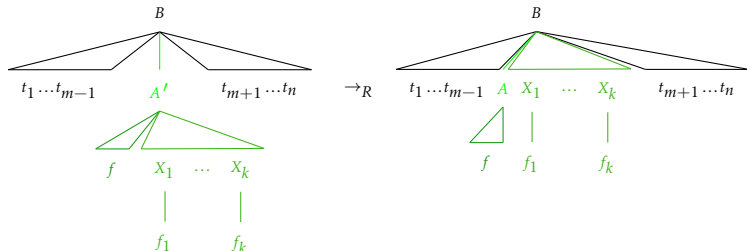
$$\begin{array}{ll}
 I \rightarrow U \mid Q & I \rightarrow U \mid Q \\
 U \rightarrow i \mid T & U \rightarrow i \mid T \\
 Q \rightarrow NU & Q \rightarrow [NU] \\
 N \rightarrow U :: N \mid U :: [NU] \rightarrow U :: [NU] \mid U :: U & \\
 T \rightarrow i \langle I \rangle & T \rightarrow i \langle I \rangle
 \end{array}$$

Define  $\mu([NU]) = NU$  and  $\mu(A) = A$  for the other nonterminals  $A$ .

# (UN-)COMBININGS AS TREE MAPPINGS

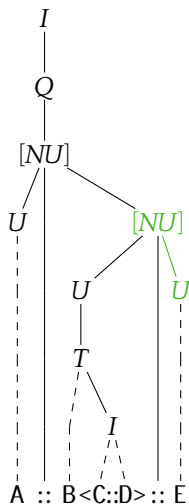
One rule per nonterminal  $A'$  s.t.

$$\mu(A') = AX_1 \cdots X_q:$$



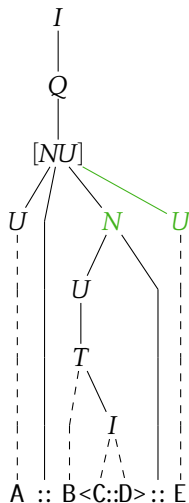
# (UN-)COMBININGS AS TREE MAPPINGS

## Example



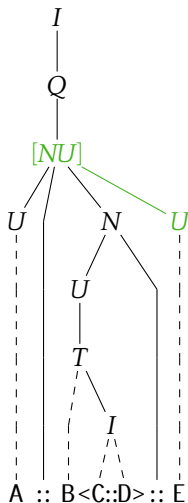
# (UN-)COMBININGS AS TREE MAPPINGS

## Example



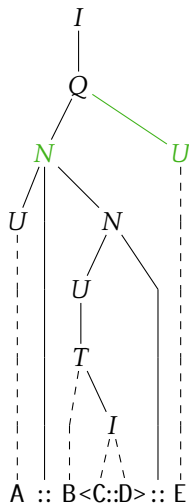
# (UN-)COMBININGS AS TREE MAPPINGS

## Example



# (UN-)COMBININGS AS TREE MAPPINGS

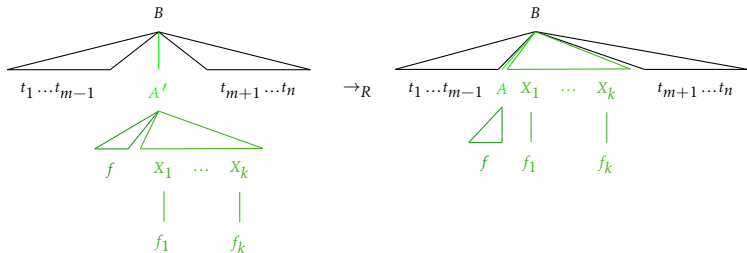
## Example



# (UN-)COMBININGS AS TREE MAPPINGS

One rule per nonterminal  $A'$  s.t.

$$\mu(A') = AX_1 \cdots X_q:$$

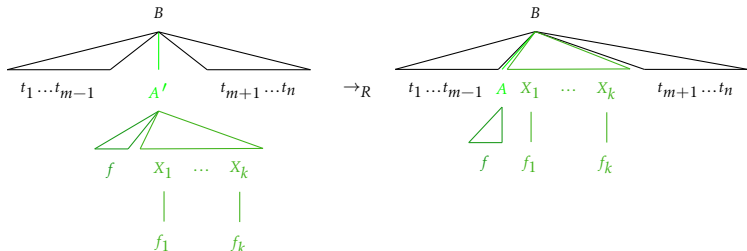


Formally, a **forest rewriting system**, which is confluent, terminating, and preserves terminal yields:

# (UN-)COMBININGS AS TREE MAPPINGS

One rule per nonterminal  $A'$  s.t.

$$\mu(A') = AX_1 \cdots X_q:$$



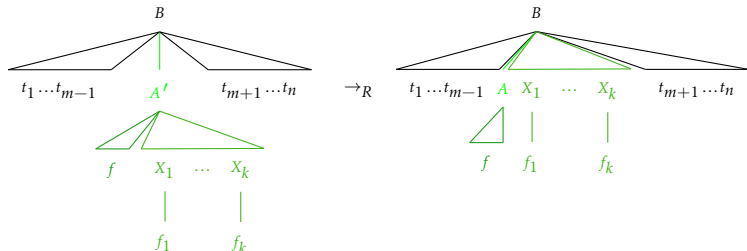
Formally, a forest rewriting system, which is **confluent**, **terminating**, and preserves terminal yields:



# (UN-)COMBININGS AS TREE MAPPINGS

One rule per nonterminal  $A'$  s.t.

$$\mu(A') = AX_1 \cdots X_q:$$



Formally, a forest rewriting system, which is confluent, terminating, and **preserves terminal yields**:

Theorem

If  $G \hookrightarrow G'$ , then  $L(G) = L(G')$ .

# SELECTIVE ML

## Definition

$G$  is **selML**( $k, m$ ) if  $\exists G', G \#^k \hookrightarrow_k G'$  and  $G'$  is LR( $m$ ).

## Combing Strategies

Given  $k$ :

- ▶ always comb “as much context as possible”?  
ML( $k, m$ ) grammars (Bertsch and Nederhof, 2007)
- ▶ always comb “as little as possible”?

# SELECTIVE ML

## Definition

$G$  is **selML**( $k, m$ ) if  $\exists G', G \#^k \hookrightarrow_k G'$  and  $G'$  is LR( $m$ ).

## Combing Strategies

Given  $k$ :

- ▶ always comb “as much context as possible”?  
ML( $k, m$ ) grammars (Bertsch and Nederhof, 2007)
- ▶ always comb “as little as possible”?

# SELECTIVE ML

## Definition

$G$  is **selML**( $k, m$ ) if  $\exists G', G \#^k \hookrightarrow_k G'$  and  $G'$  is LR( $m$ ).

## Combing Strategies

Given  $k$ :

- ▶ always comb “as much context as possible”?  
ML( $k, m$ ) grammars (Bertsch and Nederhof, 2007)
- ▶ always comb “as little as possible”?

# SELECTIVE ML

## Definition

$G$  is  $\text{selML}(k, m)$  if  $\exists G', G \#^k \hookrightarrow_k G'$  and  $G'$  is  $\text{LR}(m)$ .

## Combing Strategies

Given  $k$ :

- ▶ always comb “as much context as possible”?  
 $\text{ML}(k, m)$  grammars (Bertsch and Nederhof, 2007)
- ▶ always comb “as little as possible”?

# SEML( $k, m$ ) PARSERS

## PRINCIPLES

- ▶ LR-like construction
- ▶ add context only when required
- ▶ recompute states in case of conflicts:  
distinguished **conflict** and **deprecated** items.
- ▶ a combing can be extracted from the parser

# SEML( $k, m$ ) PARSERS

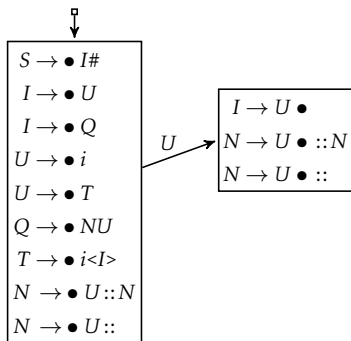
## EXAMPLE



$S \rightarrow \bullet I \#$
$I \rightarrow \bullet U$
$I \rightarrow \bullet Q$
$U \rightarrow \bullet i$
$U \rightarrow \bullet T$
$Q \rightarrow \bullet NU$
$T \rightarrow \bullet i \langle I \rangle$
$N \rightarrow \bullet U :: N$
$N \rightarrow \bullet U ::$

# SELMML( $k, m$ ) PARSERS

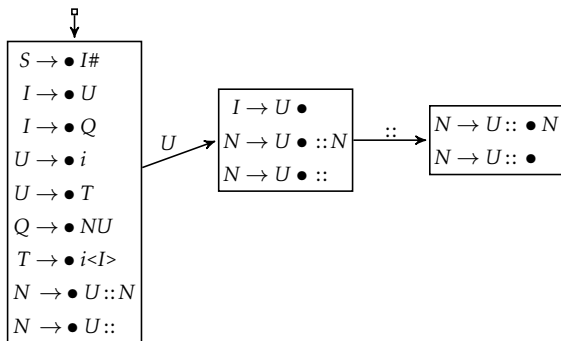
## EXAMPLE





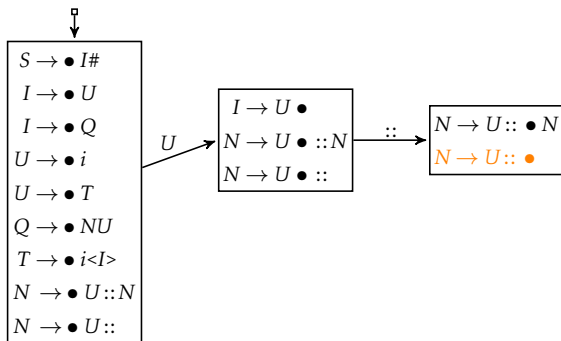
# SELMML( $k, m$ ) PARSERS

## EXAMPLE



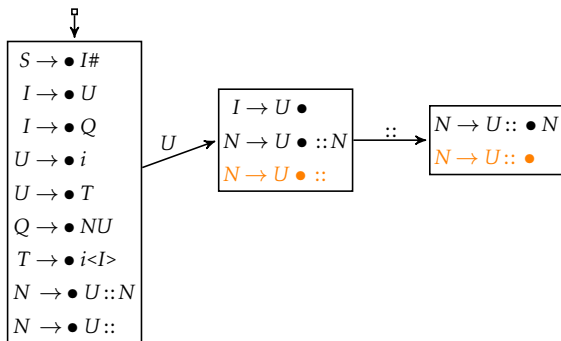
# SELMML( $k, m$ ) PARSERS

## EXAMPLE



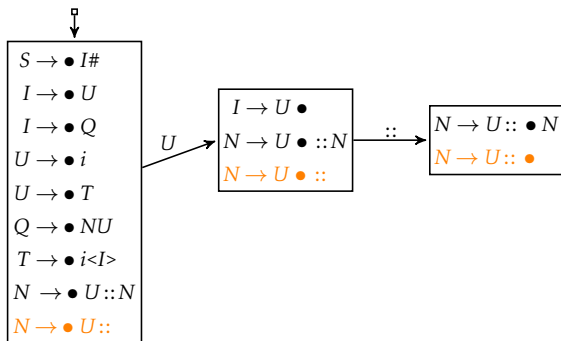
# SELMML( $k, m$ ) PARSERS

## EXAMPLE



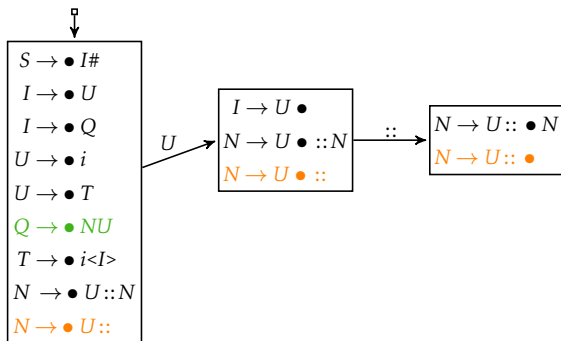
# SELMML( $k, m$ ) PARSERS

## EXAMPLE



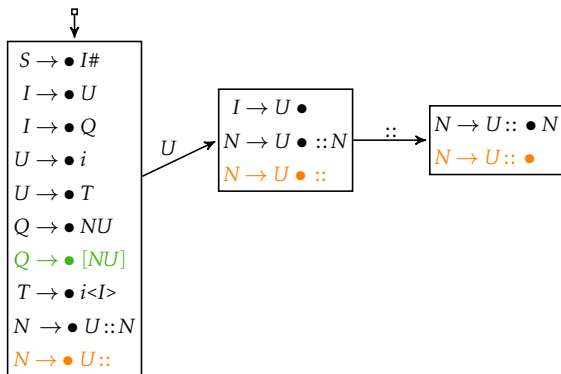
# SELMML( $k, m$ ) PARSERS

## EXAMPLE



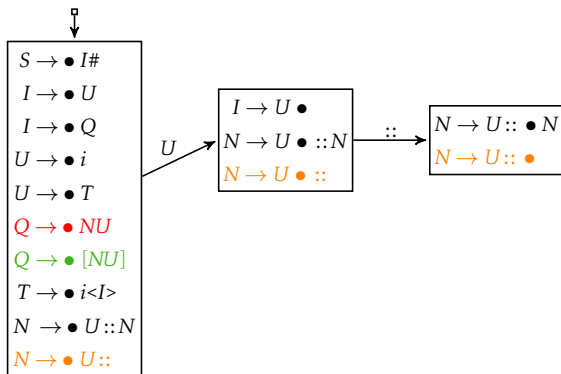
# SELMML( $k, m$ ) PARSERS

## EXAMPLE



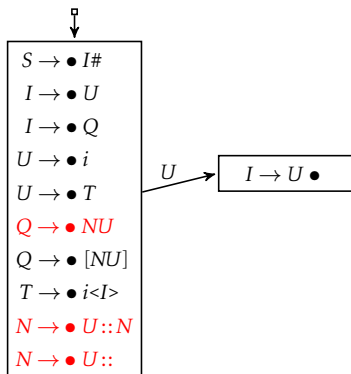
# SELMML( $k, m$ ) PARSERS

## EXAMPLE



# SELMML( $k, m$ ) PARSERS

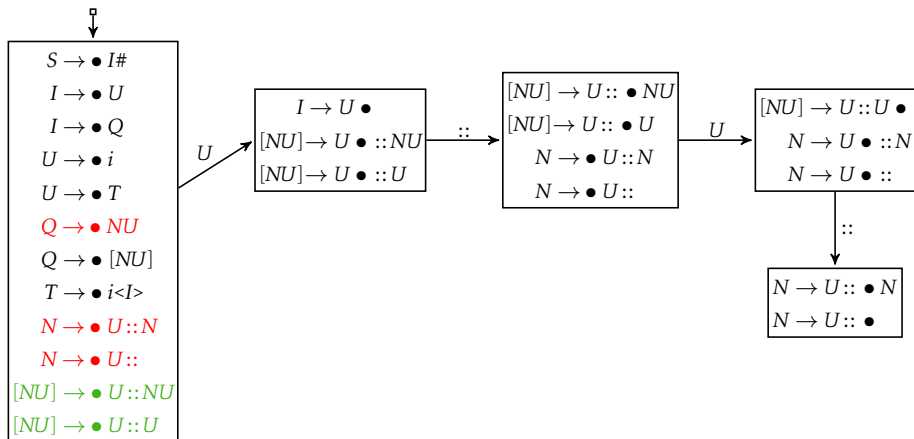
## EXAMPLE





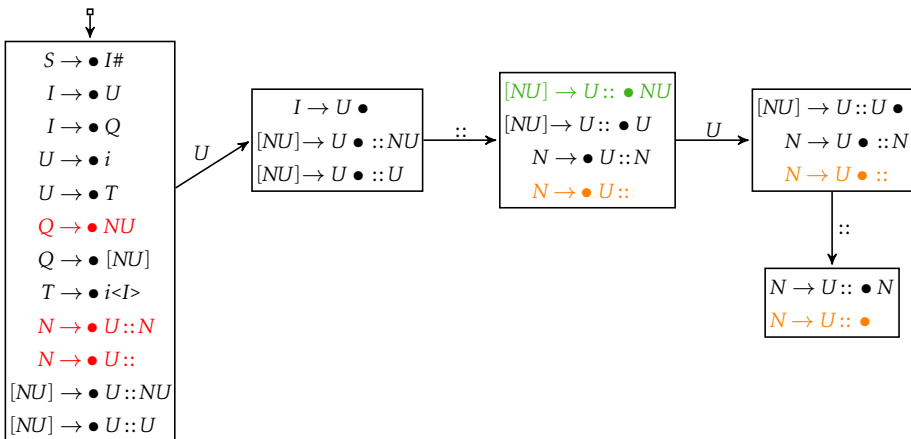
# SEMLL( $k, m$ ) PARSERS

## EXAMPLE



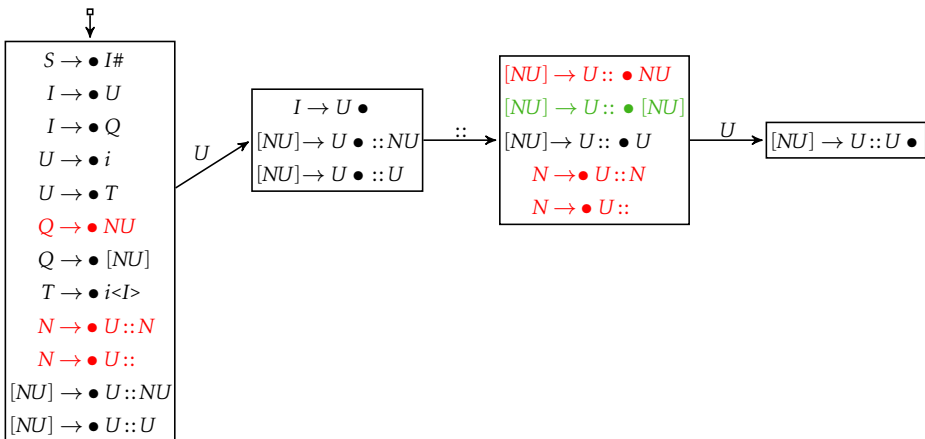
# SEML( $k, m$ ) PARSERS

## EXAMPLE



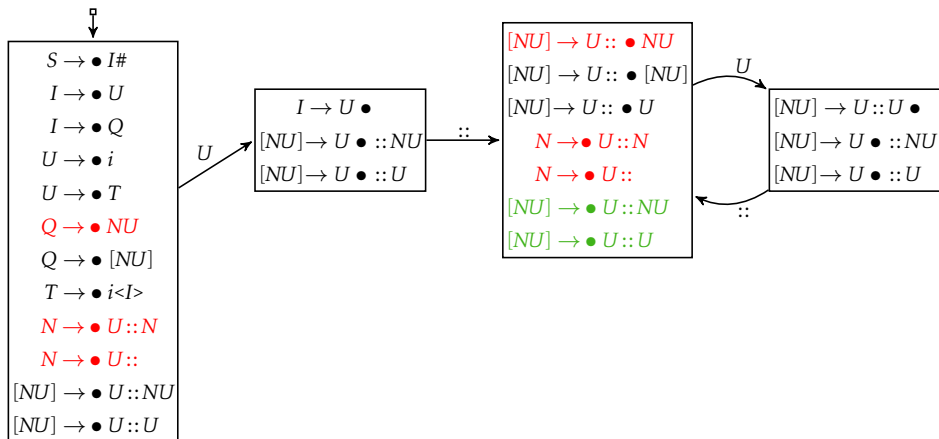
# SELMML( $k, m$ ) PARSERS

## EXAMPLE



# SEMLL( $k, m$ ) PARSERS

## EXAMPLE



# THEORETICAL RESULTS

For each  $k$  and  $m$ :

grammar class  $\text{selML}(k,m)$  includes  $\text{ML}(k,m)$ ,  
and thus  $\text{LR}(m)$

monotonicity  $\text{selML}(k,m)$  implies  $\text{selML}(k',m')$   
for all  $k' \geq k$  and  $m' \geq m$

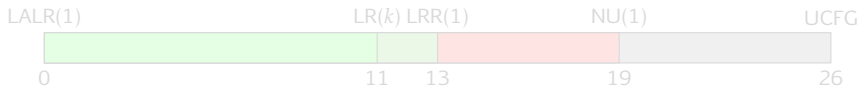
unambiguity of  $\text{selML}$  grammars

deterministic the language of a  $\text{selML}$   
grammar is a DCFL

undecidability of whether a grammar is  
 $\text{selML}(k,m)$  for some  $k$  and  $m$

# EXPERIMENTAL RESULTS

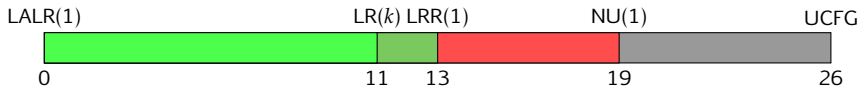
- ▶ **proof of concept** implementation available at <http://www.cs.st-andrews.ac.uk/~mjn/code/mlparsing/>
- ▶ grammar classes on 26 small grammars (Basten, 2008; Schmitz, 2010),  $k + m \leq 50$ :



- ▶ parser sizes:
  - ▶ often much smaller than ML
  - ▶ inconclusive wrt. LR

# EXPERIMENTAL RESULTS

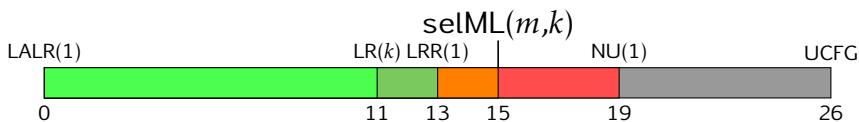
- ▶ **proof of concept** implementation available at <http://www.cs.st-andrews.ac.uk/~mjn/code/mlparsing/>
- ▶ grammar classes on 26 small grammars (Basten, 2008; Schmitz, 2010),  $k + m \leq 50$ :



- ▶ parser sizes:
  - ▶ often much smaller than ML
  - ▶ inconclusive wrt. LR

# EXPERIMENTAL RESULTS

- ▶ **proof of concept** implementation available at <http://www.cs.st-andrews.ac.uk/~mjn/code/mlparsing/>
- ▶ grammar classes on 26 small grammars (Basten, 2008; Schmitz, 2010),  $k + m \leq 50$ :

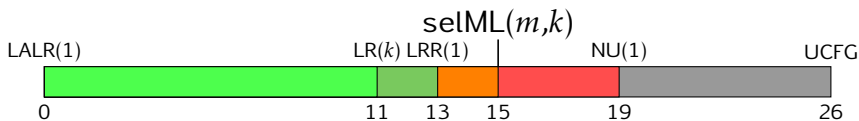


- ▶ parser sizes:
  - ▶ often much smaller than ML
  - ▶ inconclusive wrt. LR



# EXPERIMENTAL RESULTS

- ▶ **proof of concept** implementation available at <http://www.cs.st-andrews.ac.uk/~mjn/code/mlparsing/>
- ▶ grammar classes on 26 small grammars (Basten, 2008; Schmitz, 2010),  $k + m \leq 50$ :



- ▶ **parser sizes:**
  - ▶ often much smaller than ML
  - ▶ inconclusive wrt. LR

# CONCLUDING REMARKS

- ▶ avoid (some) LR conflicts
- ▶ deterministic parsing: linear-time complexity and exclusion of ambiguities
- ▶ grammar transformation view
- ▶ LR-like parser: can be used for generalized parsing.

# REFERENCES

- Basten, H.J.S., 2008. The usability of ambiguity detection methods for context-free grammars. In Vinju, J. and Johnstone, A., editors, *LDTA'08*, volume 238(5) of *Elec. Notes in Theor. Comput. Sci.*, pages 35–46. Elsevier. doi:10.1016/j.entcs.2009.09.039.
- Bertsch, E. and Nederhof, M.J., 2007. Some observations on LR-like parsing with delayed reduction. *Information Processing Letters*, 104(6):195–199. doi:10.1016/j.ipl.2007.07.003.
- Bouwers, E., Bravenboer, M., and Visser, E., 2008. Grammar engineering support for precedence rule recovery and compatibility checking. In *LDTA 2007*, volume 203(2) of *Elec. Notes in Theor. Comput. Sci.*, pages 85–101. Elsevier. doi:10.1016/j.entcs.2008.03.046.
- ISO Standard, 1998. *ISO/IEC 14882:1998: Programming Languages — C++*. International Organization for Standardization, Geneva, Switzerland.
- Lämmel, R., 2001. Grammar adaptation. In Oliveira, J.N. and Zave, P., editors, *FME 2001*, volume 2021 of *LNCS*, pages 550–570. Springer. doi:10.1007/3-540-45251-6\_32.
- Leermakers, R., 1992. Recursive ascent parsing: from Earley to Marcus. *Theor. Comput. Sci.*, 104(2):299–312. doi:10.1016/0304-3975(92)90127-2.
- Marcus, M.P., 1980. *A Theory of Syntactic Recognition for Natural Language*. Series in Artificial Intelligence. MIT Press.
- Schmitz, S., 2010. An experimental ambiguity detection tool. *Science of Computer Programming*, 75(1–2):71–84. doi:10.1016/j.scico.2009.07.002.