

Generic binarization for parsing and translation

Matthias Büchse

TU Dresden

Alexander Koller

Uni Potsdam

Heiko Vogler

TU Dresden

2013-08-05

Outline

Introduction

- CFG binarization for parsing
- SCFG binarization for translation
- Further formalisms

Generic binarization

- Common data structure: IRTG
- IRTG binarization

Results

- Algorithm
- Experiment

Outline

Introduction

- CFG binarization for parsing
- SCFG binarization for translation
- Further formalisms

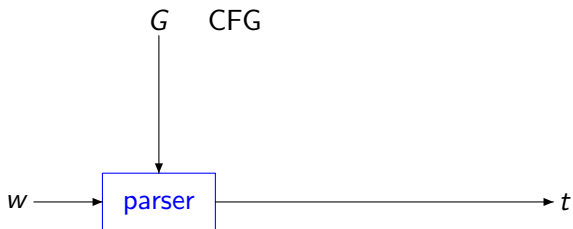
Generic binarization

- Common data structure: IRTG
- IRTG binarization

Results

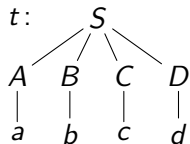
- Algorithm
- Experiment

CFG binarization for parsing

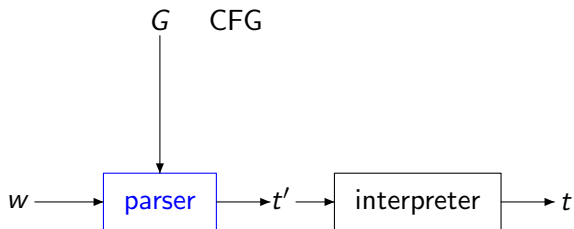


w :
 $abcd$

\rightsquigarrow

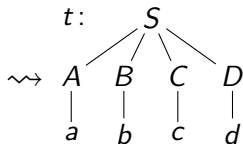


CFG binarization for parsing

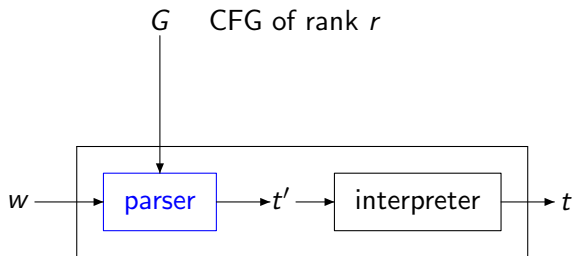


w:
abcd \rightsquigarrow

t':
 $S \rightarrow ABCD$
 $A \rightarrow a$ $B \rightarrow b$ $C \rightarrow c$ $D \rightarrow d$



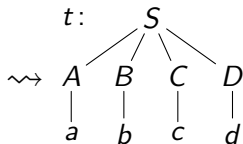
CFG binarization for parsing



$$O(|R| \cdot |w|^{r+1})$$

$w:$
 $abcd \rightsquigarrow$

$t':$ $S \rightarrow ABCD$
 $A \rightarrow a \quad B \rightarrow b \quad C \rightarrow c \quad D \rightarrow d$



CFG binarization for parsing

CFG G : (rank 4)

$S \rightarrow ABCD$

$A \rightarrow a$

$B \rightarrow b$

$C \rightarrow c$

$D \rightarrow d$

CFG binarization for parsing

CFG G : (rank 4)

$S \rightarrow ABCD$

$A \rightarrow a$

$B \rightarrow b$

$C \rightarrow c$

$D \rightarrow d$

CFG binarization for parsing

CFG G : (rank 4)

$S \rightarrow ABCD$

$A \rightarrow a$

$B \rightarrow b$

$C \rightarrow c$

$D \rightarrow d$

\rightsquigarrow

CFG G' : (rank 2)

$S \rightarrow A[BCD]$

$[BCD] \rightarrow B[CD]$

$[CD] \rightarrow CD$

\vdots

CFG binarization for parsing

CFG G : (rank 4)

$S \rightarrow ABCD$

$A \rightarrow a$

$B \rightarrow b$

$C \rightarrow c$

$D \rightarrow d$

\rightsquigarrow

CFG G' : (rank 2)

$S \rightarrow A[BCD]$

$[BCD] \rightarrow B[CD]$

$[CD] \rightarrow CD$

\vdots

rule-by-rule binarization problem *given* G :

for each rule of rank > 2

find equivalent collection of rules of rank ≤ 2

CFG binarization for parsing

CFG G : (rank 4)

$S \rightarrow ABCD$

$A \rightarrow a$

$B \rightarrow b$

$C \rightarrow c$

$D \rightarrow d$

\rightsquigarrow

CFG G' : (rank 2)

$S \rightarrow A[BCD]$

$[BCD] \rightarrow B[CD]$

$[CD] \rightarrow CD$

\vdots

rule-by-rule binarization problem *given* G :

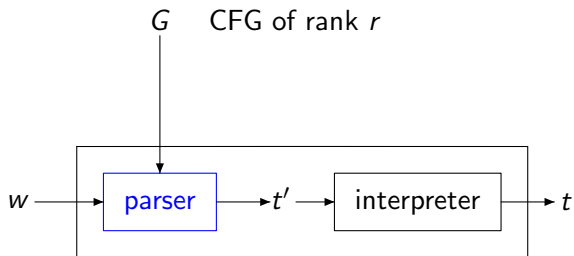
for each rule of rank > 2

find equivalent collection of rules of rank ≤ 2

binarization problem:

convert given CFG G into equivalent CFG G' of rank 2

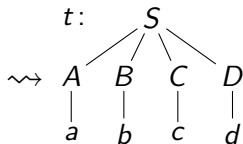
CFG binarization for parsing



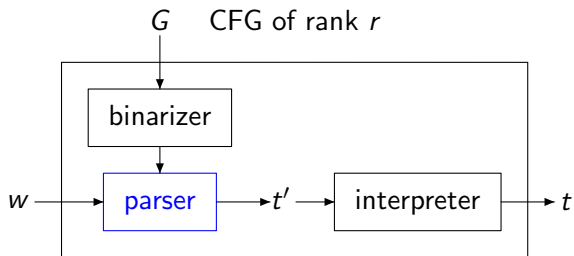
$$O(|R| \cdot |w|^{r+1})$$

$w:$
 $abcd \rightsquigarrow$

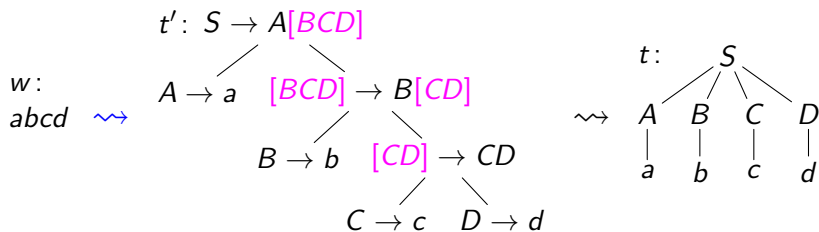
$t':$ $S \rightarrow ABCD$
 $A \rightarrow a \quad B \rightarrow b \quad C \rightarrow c \quad D \rightarrow d$



CFG binarization for parsing



$$O(r \cdot |R| \cdot |w|^3)$$



Outline

Introduction

- CFG binarization for parsing
- SCFG binarization for translation**
- Further formalisms

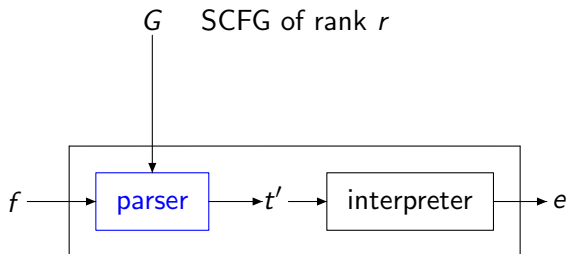
Generic binarization

- Common data structure: IRTG
- IRTG binarization

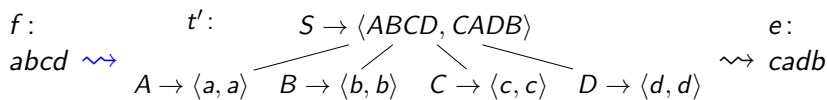
Results

- Algorithm
- Experiment

SCFG binarization for translation



$$O(|R| \cdot |w|^{r+1})$$



SCFG binarization for translation

SCFG G : (rank 4)

$$S \rightarrow \langle ABCD, CADB \rangle$$
$$A \rightarrow \langle a, a \rangle$$
$$B \rightarrow \langle b, b \rangle$$
$$C \rightarrow \langle c, c \rangle$$
$$D \rightarrow \langle d, d \rangle$$

rule-by-rule binarization problem *given* G :

for each rule of rank > 2

find equivalent collection of rules of rank ≤ 2

SCFG binarization for translation

SCFG G : (rank 4)

$S \rightarrow \langle ABCD, CADB \rangle$

$A \rightarrow \langle a, a \rangle$

$B \rightarrow \langle b, b \rangle$

$C \rightarrow \langle c, c \rangle$

$D \rightarrow \langle d, d \rangle$

rule-by-rule binarization problem *given* G :

for each rule of rank > 2

find equivalent collection of rules of rank ≤ 2

SCFG binarization for translation

SCFG G : (rank 4)

$S \rightarrow \langle ABCD, CADB \rangle$

$A \rightarrow \langle a, a \rangle$

$B \rightarrow \langle b, b \rangle$

$C \rightarrow \langle c, c \rangle$

$D \rightarrow \langle d, d \rangle$

SCFG G' : (rank 2)

$S \rightarrow \langle A[BCD], CADB \rangle$ ⚡

\rightsquigarrow

rule-by-rule binarization problem *given* G :

for each rule of rank > 2

find equivalent collection of rules of rank ≤ 2

SCFG binarization for translation

SCFG G : (rank 4)

$S \rightarrow \langle ABCD, CADB \rangle$

$A \rightarrow \langle a, a \rangle$

$B \rightarrow \langle b, b \rangle$

$C \rightarrow \langle c, c \rangle$

$D \rightarrow \langle d, d \rangle$

SCFG G' : (rank 2)

$S \rightarrow \langle A[BCD], CADB \rangle$ ⚡

$\rightsquigarrow S \rightarrow \langle [ABC]D, CADB \rangle$ ⚡

rule-by-rule binarization problem *given* G :

for each rule of rank > 2

find equivalent collection of rules of rank ≤ 2

SCFG binarization for translation

SCFG G : (rank 4)

$S \rightarrow \langle ABCD, CADB \rangle$

$A \rightarrow \langle a, a \rangle$

$B \rightarrow \langle b, b \rangle$

$C \rightarrow \langle c, c \rangle$

$D \rightarrow \langle d, d \rangle$

SCFG G' : (rank 2)

$S \rightarrow \langle A[BCD], CADB \rangle$ $\not\Leftarrow$

\rightsquigarrow $S \rightarrow \langle [ABC]D, CADB \rangle$ $\not\Leftarrow$

$S \rightarrow \langle [AB][CD], CADB \rangle$ $\not\Leftarrow$

rule-by-rule binarization problem *given* G :

for each rule of rank > 2

find equivalent collection of rules of rank ≤ 2

SCFG binarization for translation

SCFG G : (rank 4)

$S \rightarrow \langle ABCD, CADB \rangle$

$A \rightarrow \langle a, a \rangle$

$B \rightarrow \langle b, b \rangle$

$C \rightarrow \langle c, c \rangle$

$D \rightarrow \langle d, d \rangle$

SCFG G' : (rank 2)

$S \rightarrow \langle A[BCD], CADB \rangle$ $\not\Leftarrow$

\rightsquigarrow $S \rightarrow \langle [ABC]D, CADB \rangle$ $\not\Leftarrow$

$S \rightarrow \langle [AB][CD], CADB \rangle$ $\not\Leftarrow$

rule-by-rule binarization problem *given* G :

for each rule of rank > 2

find equivalent collection of rules of rank ≤ 2

solution need not exist

SCFG binarization for translation

SCFG G : (rank 4)

$S \rightarrow \langle ABCD, CADB \rangle$

$A \rightarrow \langle a, a \rangle$

$B \rightarrow \langle b, b \rangle$

$C \rightarrow \langle c, c \rangle$

$D \rightarrow \langle d, d \rangle$

SCFG G' : (rank 2)

$S \rightarrow \langle A[BCD], CADB \rangle$ ✘

$\rightsquigarrow S \rightarrow \langle [ABC]D, CADB \rangle$ ✘

$S \rightarrow \langle [AB][CD], CADB \rangle$ ✘

rule-by-rule binarization problem *given* G :

for each rule of rank > 2

find equivalent collection of rules of rank ≤ 2

solution need not exist

algorithm by Huang et al. (2009)

SCFG binarization for translation

SCFG G : (rank 4)

$S \rightarrow \langle ABCD, CADB \rangle$

$A \rightarrow \langle a, a \rangle$

$B \rightarrow \langle b, b \rangle$

$C \rightarrow \langle c, c \rangle$

$D \rightarrow \langle d, d \rangle$

SCFG G' : (rank 2)

$S \rightarrow \langle A[BCD], CADB \rangle$ ~~⚡~~

\rightsquigarrow $S \rightarrow \langle [ABC]D, CADB \rangle$ ~~⚡~~

$S \rightarrow \langle [AB][CD], CADB \rangle$ ~~⚡~~

rule-by-rule binarization problem *given* G :

for each rule of rank > 2

find equivalent collection of rules of rank ≤ 2

binarization problem:

convert given SCFG G into equivalent SCFG G' of rank 2

SCFG binarization for translation

SCFG G : (rank 4)

$S \rightarrow \langle ABCD, CADB \rangle$

$A \rightarrow \langle a, a \rangle$

$B \rightarrow \langle b, b \rangle$

$C \rightarrow \langle c, c \rangle$

$D \rightarrow \langle d, d \rangle$

SCFG G' : (rank 0)

$S \rightarrow \langle abcd, cadb \rangle$

\rightsquigarrow

rule-by-rule binarization problem *given* G :

for each rule of rank > 2

find equivalent collection of rules of rank ≤ 2

binarization problem:

convert given SCFG G into equivalent SCFG G' of rank 2

SCFG binarization for translation

SCFG G : (rank 4)

SCFG G' : (rank 0)

$S \rightarrow \langle ABCD, CADB \rangle$

$S \rightarrow \langle abcd, cadb \rangle$

$A \rightarrow \langle a, a \rangle$

\rightsquigarrow

$B \rightarrow \langle b, b \rangle$

$C \rightarrow \langle c, c \rangle$

$D \rightarrow \langle d, d \rangle$

rule-by-rule binarization problem *given* G :

for each rule of rank > 2

find equivalent collection of rules of rank ≤ 2

binarization problem:

convert given SCFG G into equivalent SCFG G' of rank 2

again: solution need not exist (Aho and Ullman, 1969)

Outline

Introduction

- CFG binarization for parsing
- SCFG binarization for translation
- Further formalisms**

Generic binarization

- Common data structure: IRTG
- IRTG binarization

Results

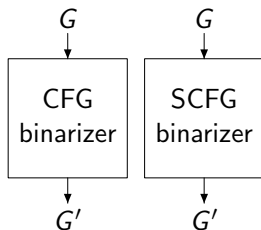
- Algorithm
- Experiment

Further formalisms

- ▶ CFG
- ▶ SCFG

(Chomsky)

(Lewis and Stearns, 1966)



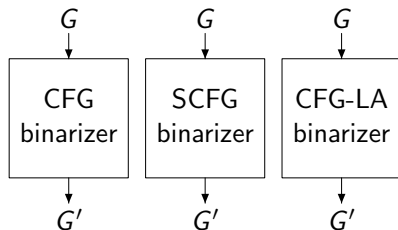
Further formalisms

- ▶ CFG
- ▶ SCFG
- ▶ CFG with latent annotations

(Chomsky)

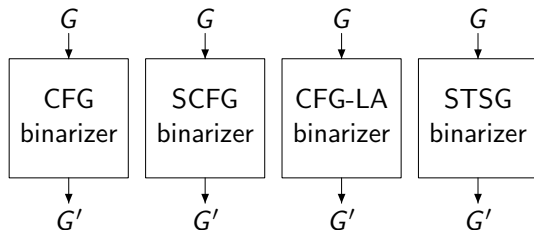
(Lewis and Stearns, 1966)

(Matsuzaki et al., 2005)



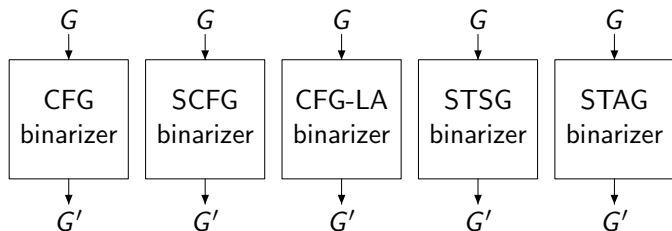
Further formalisms

- ▶ CFG (Chomsky)
- ▶ SCFG (Lewis and Stearns, 1966)
- ▶ CFG with latent annotations (Matsuzaki et al., 2005)
- ▶ synchronous tree-substitution grammars (Eisner, 2003)



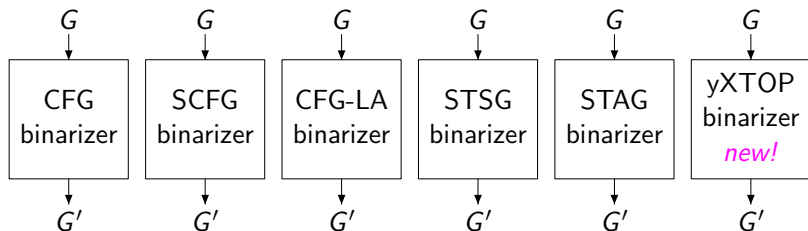
Further formalisms

- ▶ CFG (Chomsky)
- ▶ SCFG (Lewis and Stearns, 1966)
- ▶ CFG with latent annotations (Matsuzaki et al., 2005)
- ▶ synchronous tree-substitution grammars (Eisner, 2003)
- ▶ synchronous tree-adjoining grammars (Nesson et al., 2006)



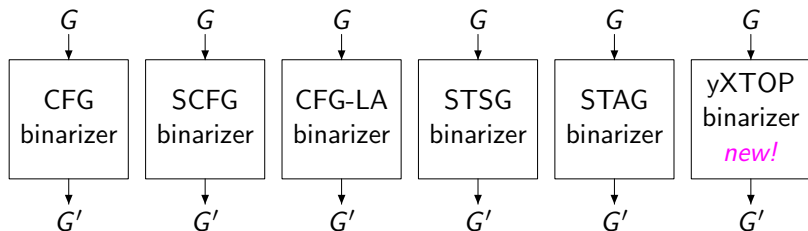
Further formalisms

- ▶ CFG (Chomsky)
- ▶ SCFG (Lewis and Stearns, 1966)
- ▶ CFG with latent annotations (Matsuzaki et al., 2005)
- ▶ synchronous tree-substitution grammars (Eisner, 2003)
- ▶ synchronous tree-adjoining grammars (Nesson et al., 2006)
- ▶ tree-to-string transducers (Graehl et al., 2008)



Further formalisms

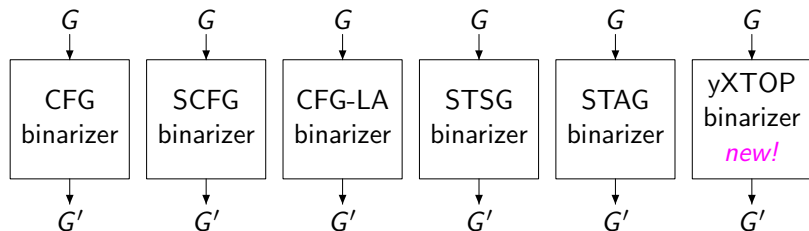
- ▶ CFG (Chomsky)
- ▶ SCFG (Lewis and Stearns, 1966)
- ▶ CFG with latent annotations (Matsuzaki et al., 2005)
- ▶ synchronous tree-substitution grammars (Eisner, 2003)
- ▶ synchronous tree-adjoining grammars (Nesson et al., 2006)
- ▶ tree-to-string transducers (Graehl et al., 2008)



six (potential) algorithms

Further formalisms

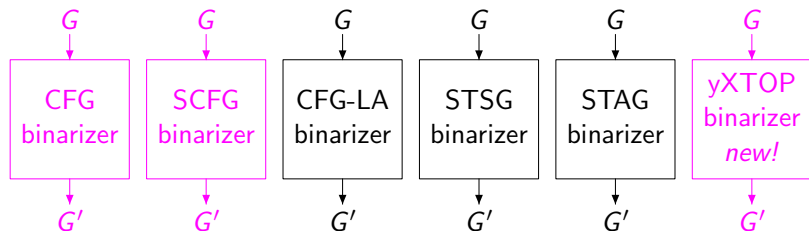
- ▶ CFG (Chomsky)
- ▶ SCFG (Lewis and Stearns, 1966)
- ▶ CFG with latent annotations (Matsuzaki et al., 2005)
- ▶ synchronous tree-substitution grammars (Eisner, 2003)
- ▶ synchronous tree-adjoining grammars (Nesson et al., 2006)
- ▶ tree-to-string transducers (Graehl et al., 2008)



six (potential) algorithms \rightsquigarrow one shared algorithm!

Further formalisms

- ▶ CFG (Chomsky)
- ▶ SCFG (Lewis and Stearns, 1966)
- ▶ CFG with latent annotations (Matsuzaki et al., 2005)
- ▶ synchronous tree-substitution grammars (Eisner, 2003)
- ▶ synchronous tree-adjoining grammars (Nesson et al., 2006)
- ▶ tree-to-string transducers (Graehl et al., 2008)



six (potential) algorithms \rightsquigarrow one shared algorithm!

Outline

Introduction

- CFG binarization for parsing
- SCFG binarization for translation
- Further formalisms

Generic binarization

- Common data structure: IRTG
- IRTG binarization

Results

- Algorithm
- Experiment

Every grammar formalism
mentioned so far
is a class of IRTGs.

IRTG ... interpreted regular tree grammar

Example

synchronous context-free grammar:

$$A \rightarrow \langle BCD, DaBC \rangle$$

$$B \rightarrow \langle b, b \rangle$$

$$C \rightarrow \langle c, c \rangle$$

$$D \rightarrow \langle d, d \rangle$$

Example

synchronous context-free grammar:

$$A \rightarrow \langle BCD, DaBC \rangle$$

$$B \rightarrow \langle b, b \rangle$$

$$C \rightarrow \langle c, c \rangle$$

$$D \rightarrow \langle d, d \rangle$$

Example

synchronous context-free grammar:

$A \rightarrow \langle BCD, DaBC \rangle$ ▶ derive pairs for B , C , and D

$B \rightarrow \langle b, b \rangle$

$C \rightarrow \langle c, c \rangle$

$D \rightarrow \langle d, d \rangle$

Example

synchronous context-free grammar:

- $A \rightarrow \langle BCD, DaBC \rangle$ ▶ derive pairs for B , C , and D
- $B \rightarrow \langle b, b \rangle$ ▶ construct
- $C \rightarrow \langle c, c \rangle$ component 1: $x_1 x_2 x_3$
- $D \rightarrow \langle d, d \rangle$

Example

synchronous context-free grammar:

$A \rightarrow \langle BCD, DaBC \rangle$ ▶ derive pairs for B , C , and D

$B \rightarrow \langle b, b \rangle$

▶ construct

$C \rightarrow \langle c, c \rangle$

component 1: $x_1 x_2 x_3$

$D \rightarrow \langle d, d \rangle$

component 2: $x_3 a x_1 x_2$

Example

synchronous context-free grammar:

- $A \rightarrow \langle BCD, DaBC \rangle$ ▶ derive pairs for B , C , and D
- $B \rightarrow \langle b, b \rangle$ ▶ construct
- $C \rightarrow \langle c, c \rangle$ component 1: $x_1 x_2 x_3$ } α
- $D \rightarrow \langle d, d \rangle$ component 2: $x_3 a x_1 x_2$ }

Example

synchronous context-free grammar:

$$\begin{array}{ll} A \rightarrow \langle BCD, DaBC \rangle & \blacktriangleright \text{derive pairs for } B, C, \text{ and } D \\ B \rightarrow \langle b, b \rangle & \blacktriangleright \text{construct} \\ C \rightarrow \langle c, c \rangle & \text{component 1: } x_1 x_2 x_3 \\ D \rightarrow \langle d, d \rangle & \text{component 2: } x_3 a x_1 x_2 \end{array} \left. \vphantom{\begin{array}{l} B \\ C \\ D \end{array}} \right\} \alpha$$

interpreted regular tree grammar:

$$\begin{array}{ll} A \rightarrow \alpha(B, C, D) & \text{con}^3(x_1, x_2, x_3) \xleftarrow{h_1} \alpha \xrightarrow{h_2} \text{con}^4(x_3, a, x_1, x_2) \\ B \rightarrow \alpha_1 & b \xleftarrow{\alpha_1} \xrightarrow{\alpha_1} b \\ C \rightarrow \alpha_2 & c \xleftarrow{\alpha_2} \xrightarrow{\alpha_2} c \\ D \rightarrow \alpha_3 & d \xleftarrow{\alpha_3} \xrightarrow{\alpha_3} d \end{array}$$

Example

synchronous context-free grammar:

$$\begin{array}{ll} A \rightarrow \langle BCD, DaBC \rangle & \blacktriangleright \text{derive pairs for } B, C, \text{ and } D \\ B \rightarrow \langle b, b \rangle & \blacktriangleright \text{construct} \\ C \rightarrow \langle c, c \rangle & \text{component 1: } x_1 x_2 x_3 \\ D \rightarrow \langle d, d \rangle & \text{component 2: } x_3 a x_1 x_2 \end{array} \left. \vphantom{\begin{array}{l} B \\ C \\ D \end{array}} \right\} \alpha$$

interpreted regular tree grammar:

$$\begin{array}{ll} A \rightarrow \alpha(B, C, D) & \text{con}^3(x_1, x_2, x_3) \xleftarrow{h_1} \alpha \xrightarrow{h_2} \text{con}^4(x_3, a, x_1, x_2) \\ B \rightarrow \alpha_1 & b \xleftarrow{\alpha_1} \xrightarrow{\alpha_1} b \\ C \rightarrow \alpha_2 & c \xleftarrow{\alpha_2} \xrightarrow{\alpha_2} c \\ D \rightarrow \alpha_3 & d \xleftarrow{\alpha_3} \xrightarrow{\alpha_3} d \end{array}$$

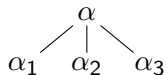
Semantics

$A \rightarrow \alpha(B, C, D)$

$B \rightarrow \alpha_1$

$C \rightarrow \alpha_2$

$D \rightarrow \alpha_3$



derivation
tree

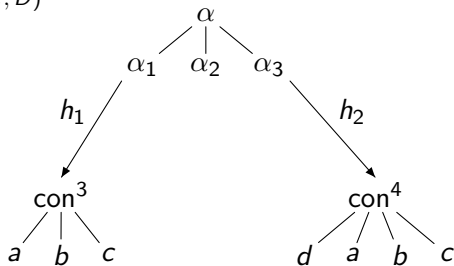
Semantics

$A \rightarrow \alpha(B, C, D)$

$B \rightarrow \alpha_1$

$C \rightarrow \alpha_2$

$D \rightarrow \alpha_3$



derivation
tree

semantic
term

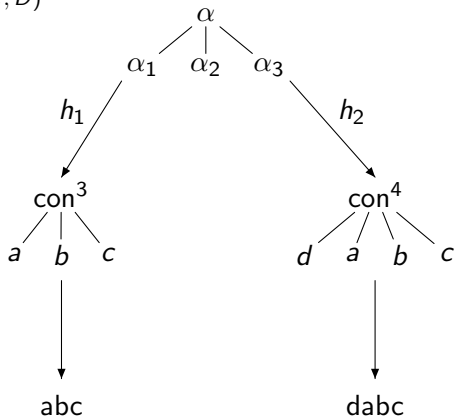
Semantics

$A \rightarrow \alpha(B, C, D)$

$B \rightarrow \alpha_1$

$C \rightarrow \alpha_2$

$D \rightarrow \alpha_3$



derivation
tree

semantic
term

semantic
object

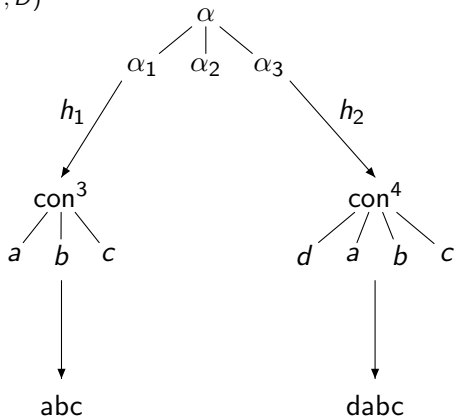
Semantics

$A \rightarrow \alpha(B, C, D)$

$B \rightarrow \alpha_1$

$C \rightarrow \alpha_2$

$D \rightarrow \alpha_3$



derivation
tree

semantic
term

semantic
object

interpretation 1

interpretation 2

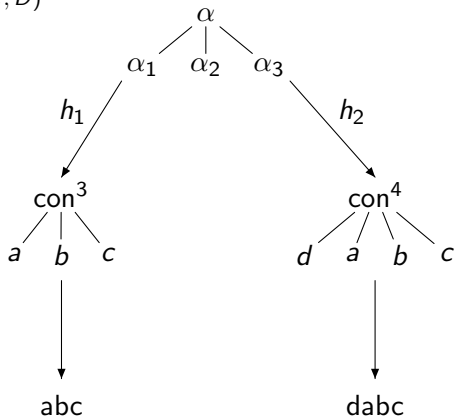
Semantics

$A \rightarrow \alpha(B, C, D)$

$B \rightarrow \alpha_1$

$C \rightarrow \alpha_2$

$D \rightarrow \alpha_3$



derivation
tree

semantic
term

semantic
object

interpretation 1

interpretation 2

$n \dots$ number of interpretations

Why the trouble?

Why the trouble?

Every grammar formalism
mentioned so far
is a class of IRTGs.

Outline

Introduction

- CFG binarization for parsing
- SCFG binarization for translation
- Further formalisms

Generic binarization

- Common data structure: IRTG
- IRTG binarization**

Results

- Algorithm
- Experiment

Binarization problem

rule-by-rule binarization problem *given* G :

for each rule of rank > 2

find equivalent collection of rules of rank ≤ 2

Binarization problem

rule-by-rule binarization problem *given* G :

for each rule of rank > 2

find equivalent collection of rules of rank ≤ 2

binarization problem:

convert given IRTG G into equivalent IRTG G' of rank 2

Binarization problem

rule-by-rule binarization problem *given* G :

for each rule of rank > 2

find **equivalent** collection of rules of rank ≤ 2

binarization problem:

convert given IRTG G into **equivalent** IRTG G' of rank 2

Binarization problem

rule-by-rule binarization problem *given* G :

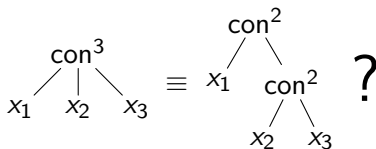
for each rule of rank > 2

find equivalent collection of rules of rank ≤ 2

binarization problem:

convert given IRTG G into equivalent IRTG G' of rank 2

equivalence problem: is it true that



Binarization problem

rule-by-rule binarization problem *given* G :

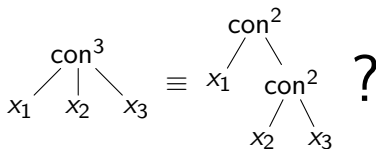
for each rule of rank > 2

find equivalent collection of rules of rank ≤ 2

binarization problem:

convert given IRTG G into equivalent IRTG G' of rank 2

equivalence problem: is it true that



undecidable in general

Binarization problem

rule-by-rule binarization problem *given* G :

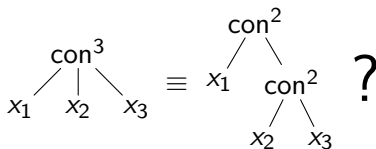
for each rule of rank > 2

find equivalent collection of rules of rank ≤ 2

binarization problem:

convert given IRTG G into equivalent IRTG G' of rank 2

equivalence problem: is it true that

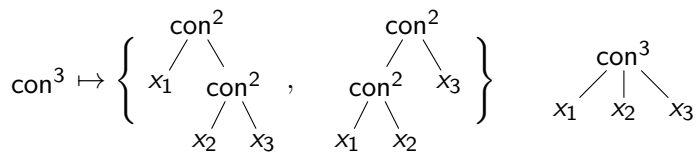


undecidable in general \rightsquigarrow adapt problem specification!

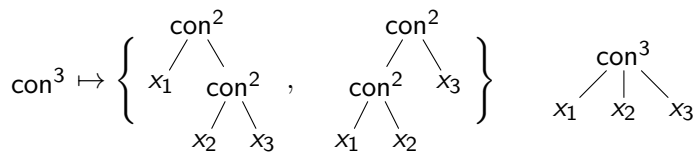
Binarization rule \flat : operation symbols \rightarrow binary terms

$$\text{con}^3 \mapsto \left\{ \begin{array}{c} \text{con}^2 \\ / \quad \backslash \\ x_1 \quad \text{con}^2 \\ \quad / \quad \backslash \\ \quad x_2 \quad x_3 \end{array} , \begin{array}{c} \text{con}^2 \\ / \quad \backslash \\ \text{con}^2 \quad x_3 \\ / \quad \backslash \\ x_1 \quad x_2 \end{array} \right\}$$

Binarization rule \flat : operation symbols \rightarrow binary terms

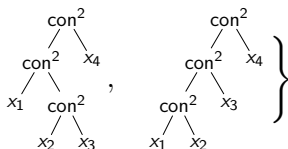
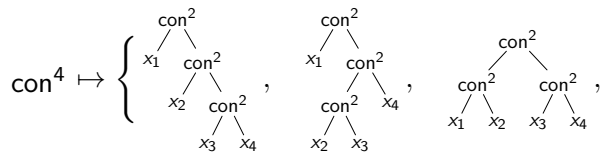
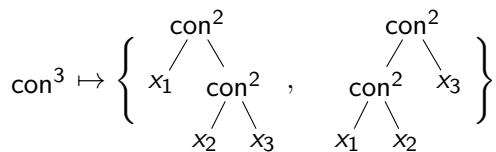


Binarization rule \flat : operation symbols \rightarrow binary terms



$$x_1 \cdot (x_2 \cdot x_3) = (x_1 \cdot x_2) \cdot x_3 = x_1 \cdot x_2 \cdot x_3$$

Binarization rule \flat : operation symbols \rightarrow binary terms



⋮

Refined binarization problem

rule-by-rule binarization problem *given* G
with respect to b_1, \dots, b_n :

for each rule of rank > 2

find collection of rules of rank ≤ 2

equivalent with respect to b_1, \dots, b_n

Refined binarization problem

rule-by-rule binarization problem *given* G
with respect to b_1, \dots, b_n :

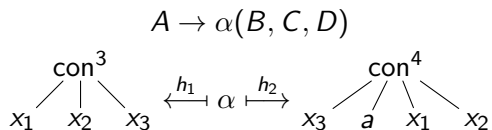
for each rule of rank > 2

find collection of rules of rank ≤ 2

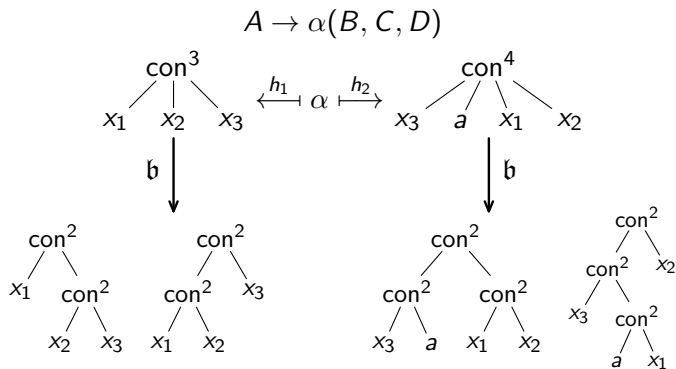
equivalent with respect to b_1, \dots, b_n

remember: solution need not exist

Using the binarization rule \flat constructively

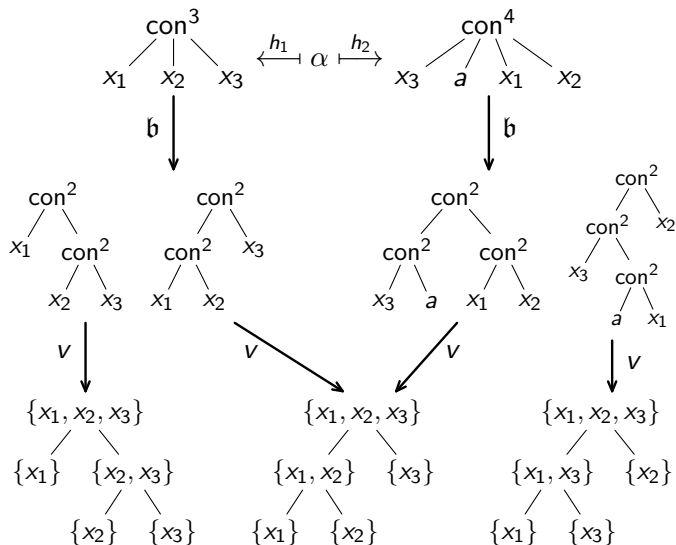


Using the binarization rule \mathfrak{b} constructively

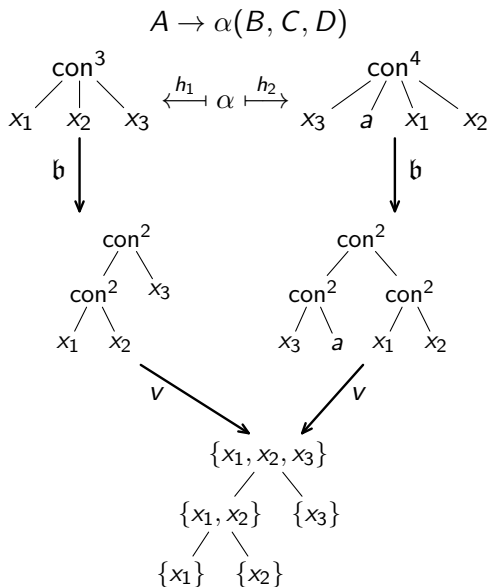


Using the binarization rule \mathfrak{b} constructively

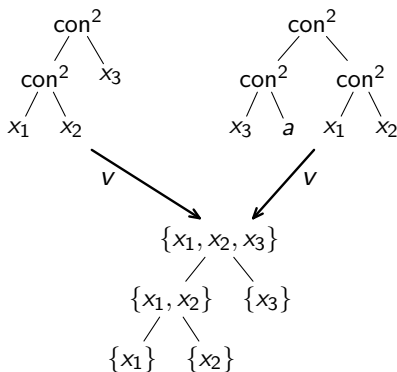
$$A \rightarrow \alpha(B, C, D)$$



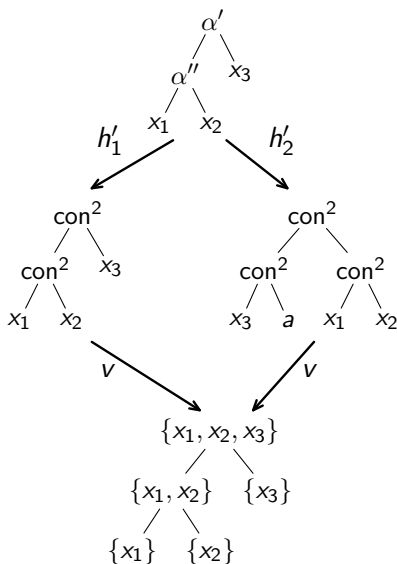
Using the binarization rule \mathfrak{b} constructively



Using the binarization rule \flat constructively

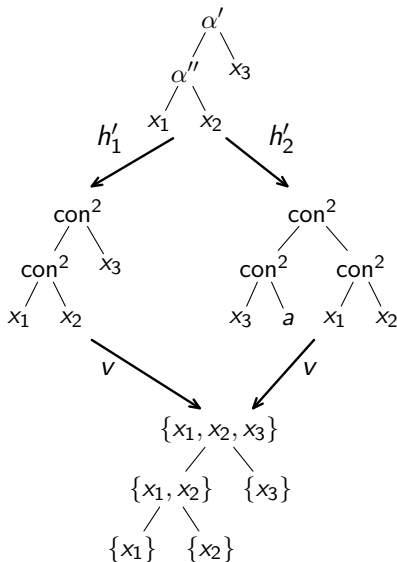


Using the binarization rule \mathfrak{b} constructively



Using the binarization rule \flat constructively

$$A \rightarrow \alpha'([BC], D) \quad [BC] \rightarrow \alpha''(B, C)$$



Outline

Introduction

- CFG binarization for parsing
- SCFG binarization for translation
- Further formalisms

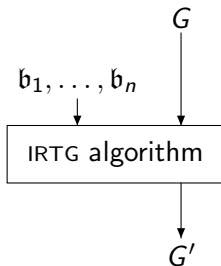
Generic binarization

- Common data structure: IRTG
- IRTG binarization

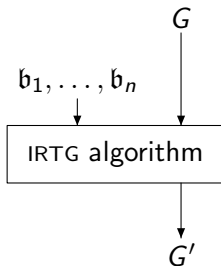
Results

- Algorithm**
- Experiment

Algorithm

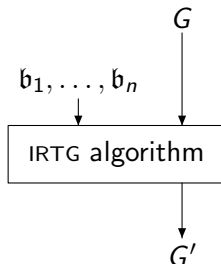


Algorithm



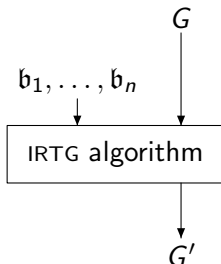
- ▶ G' is equivalent to G

Algorithm



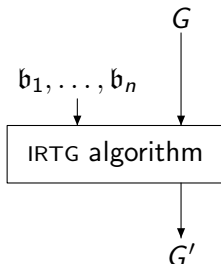
- ▶ G' is equivalent to G
- ▶ G' has rank 2
if the rule-by-rule binarization problem *given* G
with respect to b_1, \dots, b_n has a solution

Algorithm



- ▶ G' is equivalent to G
- ▶ G' has rank 2
 - if the rule-by-rule binarization problem *given* G with respect to b_1, \dots, b_n has a solution
 - otherwise: best effort

Algorithm

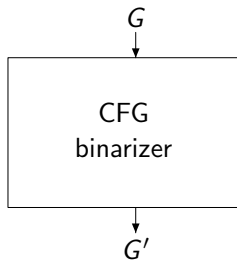


- ▶ G' is equivalent to G
- ▶ G' has rank 2
 - if the rule-by-rule binarization problem given G with respect to b_1, \dots, b_n has a solution
 - otherwise: best effort
- ▶ runtime $O(|R| \cdot c^n)$

Applications

Solve rule-by-rule binarization problem (whenever possible)

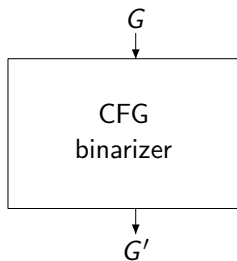
for every CFG:



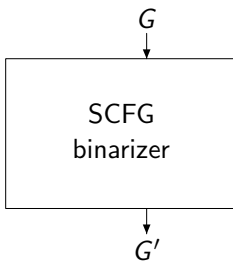
Applications

Solve rule-by-rule binarization problem (whenever possible)

for every CFG:



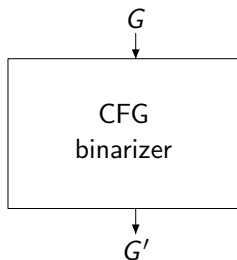
for every SCFG:



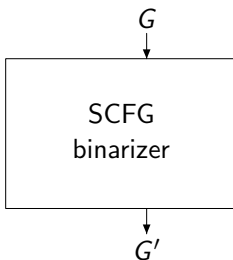
Applications

Solve rule-by-rule binarization problem (whenever possible)

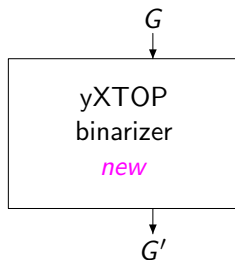
for every CFG:



for every SCFG:



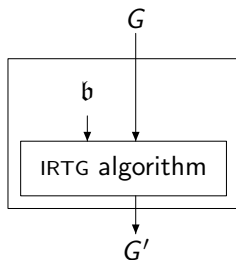
for every
tree-to-string
transducer:



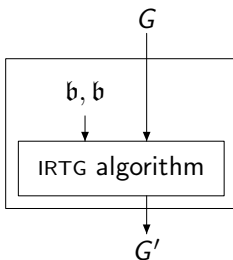
Applications

Solve rule-by-rule binarization problem (whenever possible)

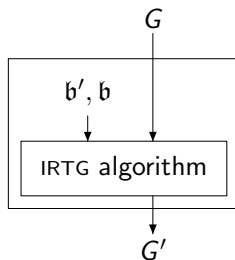
for every CFG:



for every SCFG:



for every
tree-to-string
transducer:

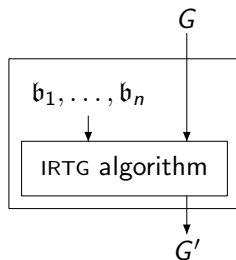


... with merely two binarization rules b, b' !

More generally

Solve rule-by-rule binarization problem (whenever possible)

for every $G \in \mathcal{C}$:



provided that b_1, \dots, b_n are **complete** for every $G \in \mathcal{C}$

Outline

Introduction

- CFG binarization for parsing
- SCFG binarization for translation
- Further formalisms

Generic binarization

- Common data structure: IRTG
- IRTG binarization

Results

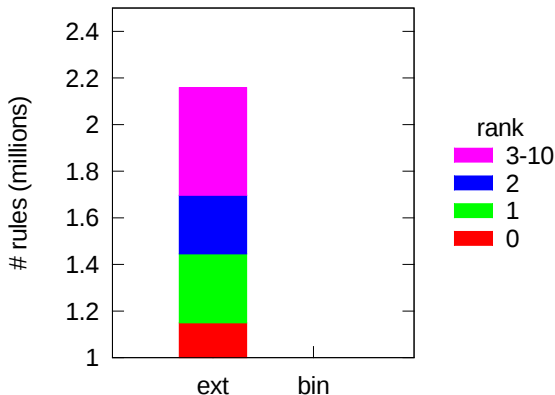
- Algorithm
- Experiment**

Experiment: tree-to-string transducers

- ▶ extracted from 1 M parallel sentences (Europarl Eng-Ger)

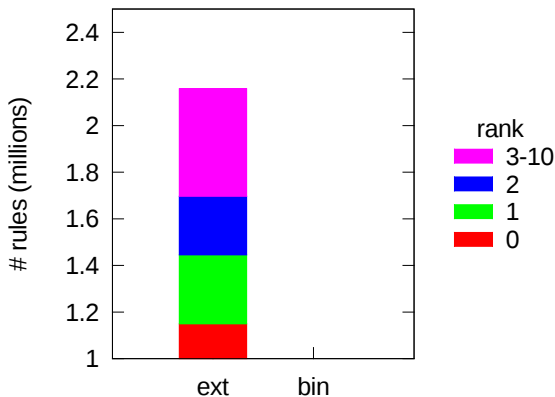
Experiment: tree-to-string transducers

- ▶ extracted from 1 M parallel sentences (Europarl Eng-Ger)
- ▶ 2.15 M rules up to rank 10



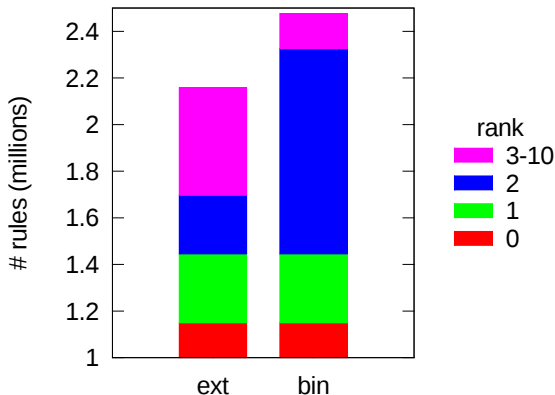
Experiment: tree-to-string transducers

- ▶ extracted from 1 M parallel sentences (Europarl Eng-Ger)
- ▶ 2.15 M rules up to rank 10
- ▶ implementation in Haskell



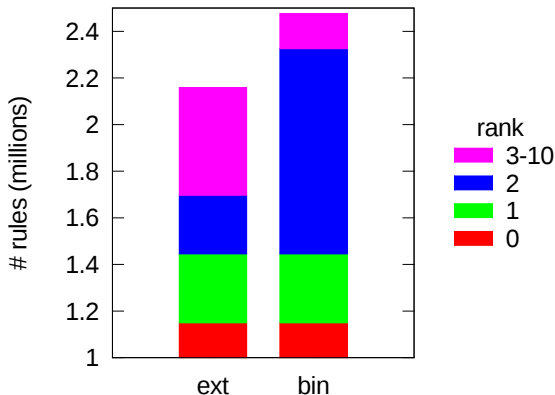
Experiment: tree-to-string transducers

- ▶ extracted from 1 M parallel sentences (Europarl Eng-Ger)
- ▶ 2.15 M rules up to rank 10
- ▶ implementation in Haskell
- ▶ 67% rules could be binarized



Experiment: tree-to-string transducers

- ▶ extracted from 1 M parallel sentences (Europarl Eng-Ger)
- ▶ 2.15 M rules up to rank 10
- ▶ implementation in Haskell
- ▶ 67% rules could be binarized
- ▶ took 4.4 min on a single core of an i5 2520M



Need to binarize?
The algorithm is there.
Supply binarization rules.

Thank you!