

# Lecture 7: Introduction to syntax-based MT

Andreas Maletti

Statistical Machine Translation

Stuttgart — December 16, 2011

# Lecture 7

## Goals

- Overview
- Tree substitution grammars (tree automata)
- Synchronous grammars (tree transducers)

# Contents

- 1 Overview
- 2 Tree representations
- 3 Bar-Hillel Construction
- 4 Variants
- 5 Tree Transducers

## Word-based system (FST)

And then the matter was decided , and everything was put in place

ف كان ان تم الحسم و وضعت الأمور في نصاب ها  
*f kAn An tm AlHsm w wDEt Al>mwr fy nSAb hA*

### Derivation

Input:

And then the matter was decided , and everything was put in place

Output:

## Word-based system (FST)

And then the matter was decided , and everything was put in place

ف	كان	ان	تم	المسألة	و	وضعت	الأمر	في	نصاب	ها
<i>f</i>	<i>kAn</i>	<i>An</i>	<i>tm</i>	<i>AlHsm</i>	<i>w</i>	<i>wDEt</i>	<i>Al&gt;mwr</i>	<i>fy</i>	<i>nSAb</i>	<i>hA</i>

### Derivation

Input:

then *the matter was decided , and everything was put in place*

Output:

## Word-based system (FST)

And then the matter was decided , and everything was put in place

ها نصاب في الأمور وضعت و الحسم تم ان كان ف  
hA nSAb fy Al>mwr wDEt w AlHsm tm An kAn f

### Derivation

Input:

the matter was decided , and everything was put in place

Output:

# Word-based system (FST)

And then the matter was decided , and everything was put in place

ها نصاب في الأمور وضعت و تم الحسم ان كان ف  
hA nSAb fy Al>mwr wDEt w wDEt Al>mwr fy nSAb hA  
f kAn An tm AlHsm w wDEt Al>mwr fy nSAb hA

## Derivation

Input:

the matter was decided , and everything was put in place

Output:

f

## Word-based system (FST)

And then the matter was decided , and everything was put in place

ف	كان	ان	تم	الحسم	و	وضعت	الأمر	في	نصاب	ها
<i>f</i>	<i>kAn</i>	<i>An</i>	<i>tm</i>	<i>AlHsm</i>	<i>w</i>	<i>wDEt</i>	<i>Al&gt;mwr</i>	<i>fy</i>	<i>nSAb</i>	<i>hA</i>

### Derivation

Input:

the matter was decided , and everything was put in place

Output:

*f kAn*



# Word-based system (FST)

And then the matter was decided , and everything was put in place

ها نصاب في الأمور وضعت و الحسم تم ان كان ف  
hA nSAb fy Al>mwr wDEt w و AlHsm tm An kAn f

## Derivation

Input:

*the* matter *was decided , and everything was put in place*

Output:

*f kAn*

# Word-based system (FST)

And then the matter was decided , and everything was put in place

ف كان ان تم الحسم و وضعت الأمور في نصاب ها  
f kAn An tm AlHsm w wDEt Al>mwr fy nSAb hA

## Derivation

Input:

*the matter* was decided , and everything was put in place

Output:

f kAn

# Word-based system (FST)

And then the matter was decided , and everything was put in place

ف كان ان تم الحسم و وضعت الأمور في نصاب ها  
*f kAn An tm AlHsm w wDEt Al>mwr fy nSAb hA*

## Derivation

Input:

*the matter was **decided** , and everything was put in place*

Output:

*f kAn*

# Word-based system (FST)

And then the matter was decided , and everything was put in place

ف كان ان تم الحسم و وضعت الأمور في نصاب ها  
f kAn An tm AlHsm w wDEt Al>mwr fy nSAb hA

## Derivation

Input:

*the matter* , *and everything was put in place*

Output:

*f kAn An tm AlHsm*

# Word-based system (FST)

And then the matter was decided , and everything was put in place

ف كان ان تم الجسم و وضعت الأمور في نصابها  
*f kAn An tm AlHsm w wDEt Al>mwr fy nSAb hA*

## Derivation

Input:

*the matter* and *everything was put in place*

Output:

*f kAn An tm AlHsm*

# Word-based system (FST)

And then the matter was decided , and everything was put in place

ها نصاب في الأمور وضعت و تم الحسم كان ف  
hA nSAb fy Al>mwr wDEt w tm AlHsm kAn f

## Derivation

Input:

*the matter* everything *was put in place*

Output:

*f kAn An tm AlHsm w*

# Word-based system (FST)

And then the matter was decided , and everything was put in place

ها نصاب في الأمور وضعت و تم الحسم ان كان ف  
hA nSAb fy Al>mwr wDEt w ALHsm tm An kAn f

## Derivation

Input:

*the matter* was *put in place*

Output:

*f kAn An tm ALHsm w*

## Word-based system (FST)

And then the matter was decided , and everything was put in place

ها نصاب في الأمور وضعت و الحسم تم ان كان ف  
hA nSAb fy Al>mwr wDEt w ALHsm tm An kAn f

### Derivation

Input:

*the matter was **put** in place*

Output:

*f kAn An tm ALHsm w*



# Word-based system (FST)

And then the matter was decided , and everything was put in place

ها نصاب في الأمور وضعت و تم الحسم ان كان ف  
hA nSAb fy Al>mwr wDEt w wDEt Al>mwr fy nSAb hA  
f kAn An tm AlHsm w wDEt Al>mwr fy nSAb hA

## Derivation

Input:

*the matter* in *place*

Output:

*f kAn An tm AlHsm w wDEt*

# Word-based system (FST)

And then the matter was decided , and everything was put in place

ف	كان	ان	تم	الحسم	و	وضعت	الأمر	في	نصاب	ها
<i>f</i>	<i>kAn</i>	<i>An</i>	<i>tm</i>	<i>AlHsm</i>	<i>w</i>	<i>wDEt</i>	<i>Al&gt;mwr</i>	<i>fy</i>	<i>nSAb</i>	<i>hA</i>

## Derivation

Input:

in place

Output:

*f kAn An tm AlHsm w wDEt Al>mwr*

# Word-based system (FST)

And then the matter was decided , and everything was put in place

ف كان ان تم الحسم و وضعت الأمور في نصاب ها  
*f kAn An tm AlHsm w wDEt Al>mwr fy nSAb hA*

## Derivation

Input:

*place*

Output:

*f kAn An tm AlHsm w wDEt Al>mwr fy*

# Word-based system (FST)

And then the matter was decided , and everything was put in place

ف	كان	ان	تم	الحسم	و	وضعت	الأمر	في	نصاب	ها
<i>f</i>	<i>kAn</i>	<i>An</i>	<i>tm</i>	<i>AlHsm</i>	<i>w</i>	<i>wDEt</i>	<i>Al&gt;mwr</i>	<i>fy</i>	<i>nSAb</i>	<i>hA</i>

## Derivation

Input:

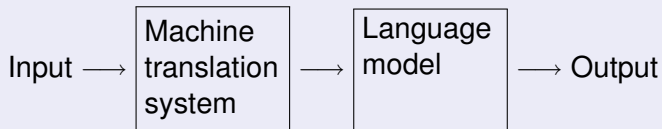


Output:

*f kAn An tm AlHsm w wDEt Al>mwr fy nSAb hA*

# Phrase-based machine translation

## Schema

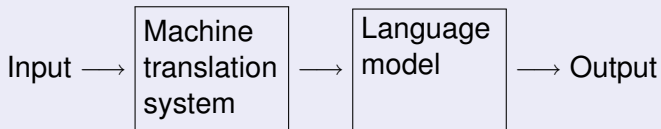


## Phrase-based systems

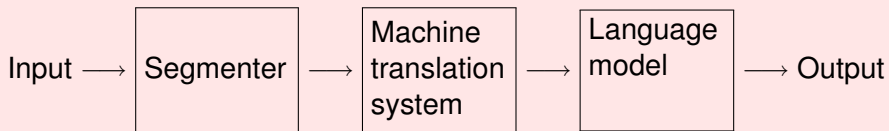


# Phrase-based machine translation

## Schema



## Phrase-based systems



## Phrase-based system (FST+Perm)

And then the matter was decided , and everything was put in place

ف	كان	ان	تم	الحسم	و	وضعت	الأمر	في	نصاب	ها
<i>f</i>	<i>kAn</i>	<i>An</i>	<i>tm</i>	<i>AlHsm</i>	<i>w</i>	<i>wDEt</i>	<i>Al&gt;mwr</i>	<i>fy</i>	<i>nSAb</i>	<i>hA</i>

### Derivation

Input:

*And then the matter was decided , and everything was put in place*

Output:

# Phrase-based system (FST+Perm)

And then the matter was decided , and everything was put in place

ف كان ان تم الحسم و وضعت الأمور في نصاب ها  
*f kAn An tm AlHsm w wDEt Al>mwr fy nSAb hA*

## Derivation

Input:

And then<sub>1</sub> the matter<sub>5</sub> was decided<sub>2</sub> , and everything<sub>3</sub> was put<sub>4</sub> in place<sub>6</sub>

Output:



# Phrase-based system (FST+Perm)

And then the matter was decided , and everything was put in place

ف كان ان تم الحسم و وضعت الامور في نصاب ها  
*f kAn An tm AlHsm w wDEt Al>mwr fy nSAb hA*

## Derivation

Input:

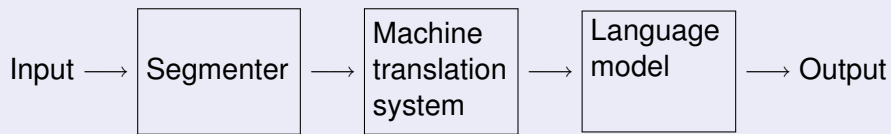
*And then*<sub>1</sub> *the matter*<sub>5</sub> *was decided*<sub>2</sub> *, and everything*<sub>3</sub> *was put*<sub>4</sub> *in place*<sub>6</sub>

Output:

*f kAn*<sub>1</sub> *An tm AlHsm*<sub>2</sub> *w*<sub>3</sub> *wDEt*<sub>4</sub> *Almwr*<sub>5</sub> *fy nSAb hA*<sub>6</sub>

# Machine translation (cont'd)

## Phrase-based systems

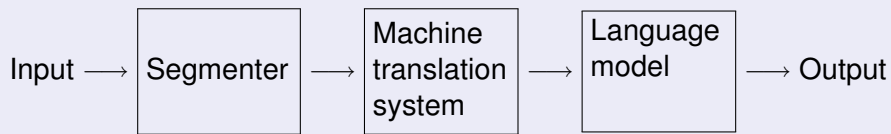


## Syntax-based systems

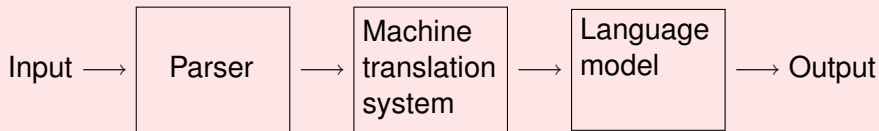


# Machine translation (cont'd)

## Phrase-based systems

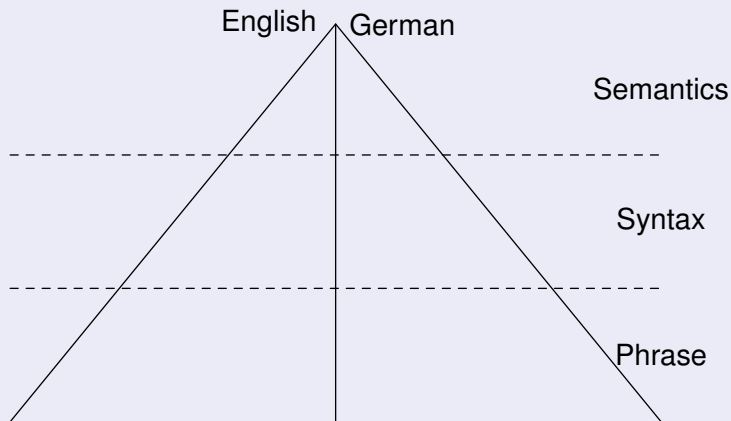


## Syntax-based systems



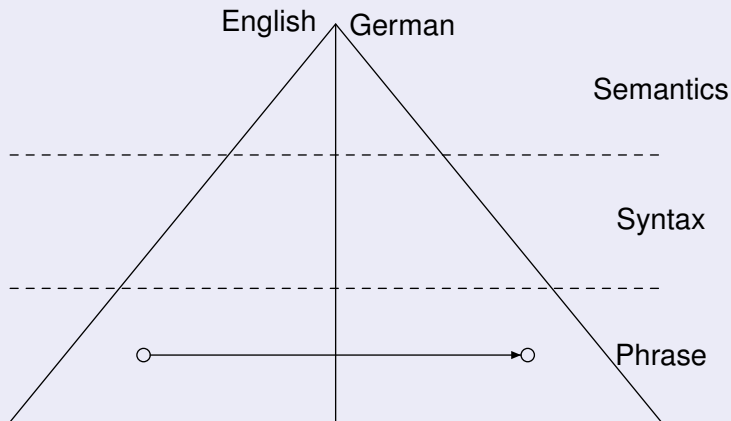
# Syntax-based Approach

## Overview



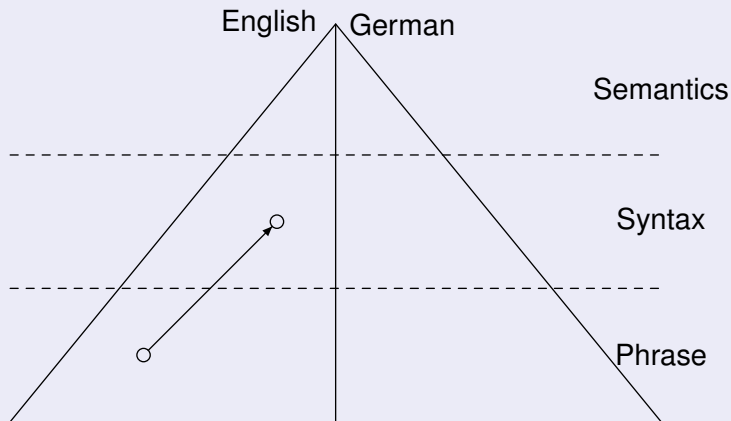
# Syntax-based Approach

## Overview



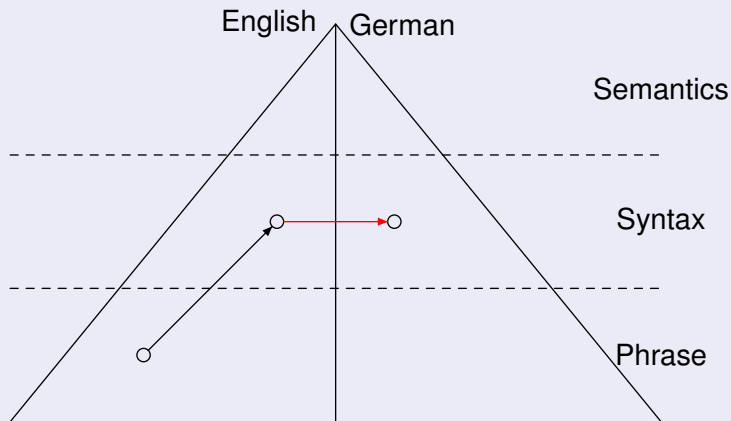
# Syntax-based Approach

## Overview



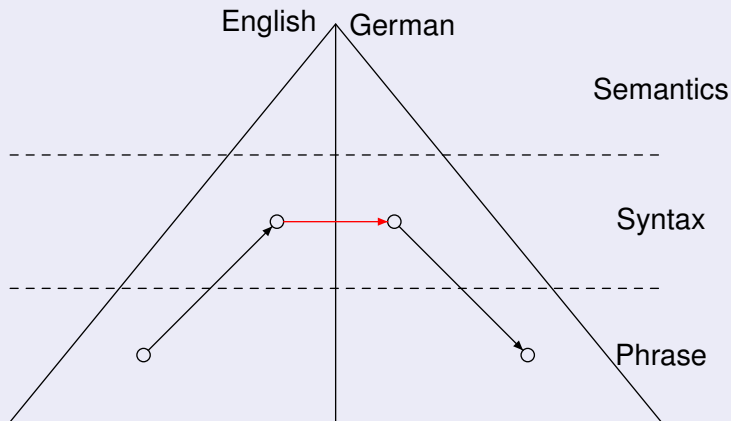
# Syntax-based Approach

## Overview



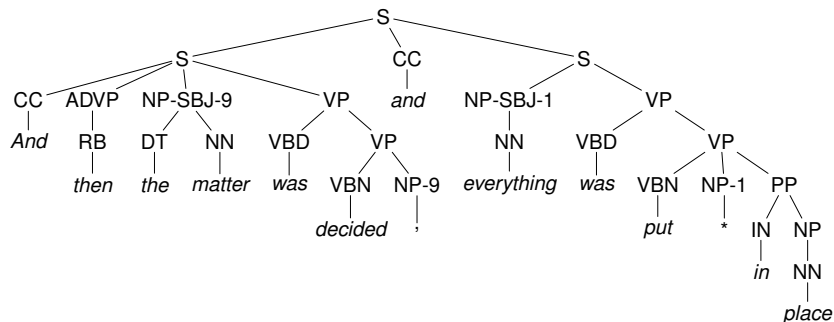
# Syntax-based Approach

## Overview





# Parser

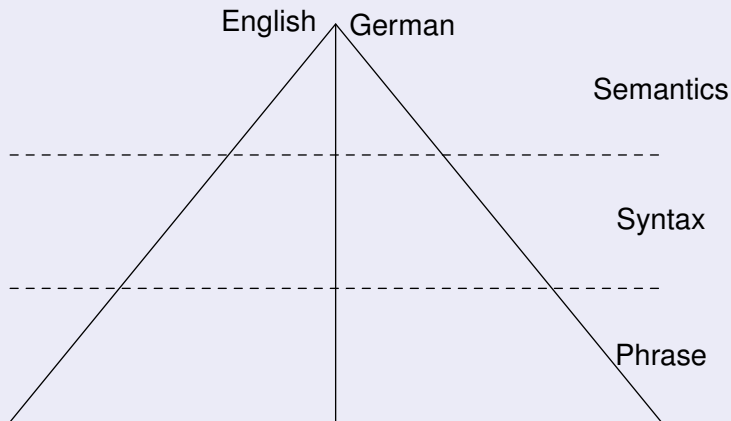


*And then the matter was decided , and everything was put in place*

(thanks to [KEVIN KNIGHT](#) for the data)

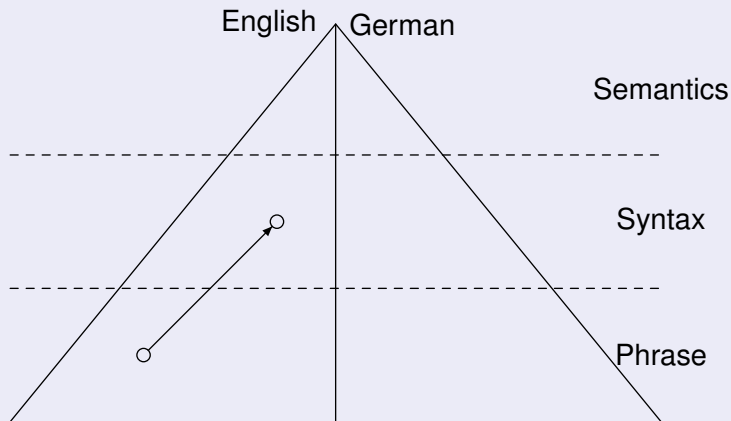
# Semantics-based Approach

## Overview



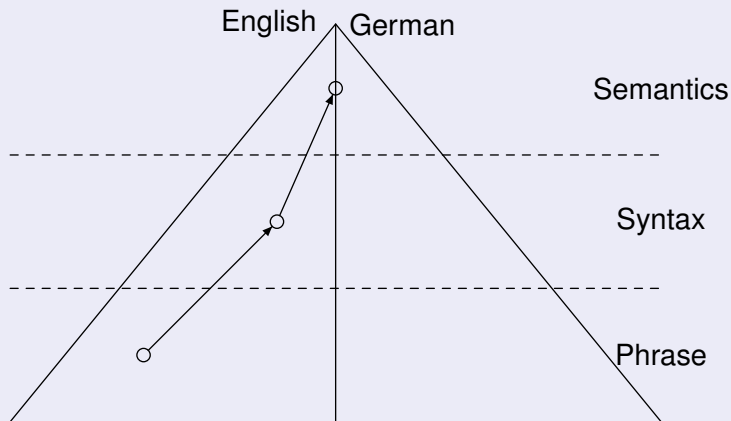
# Semantics-based Approach

## Overview



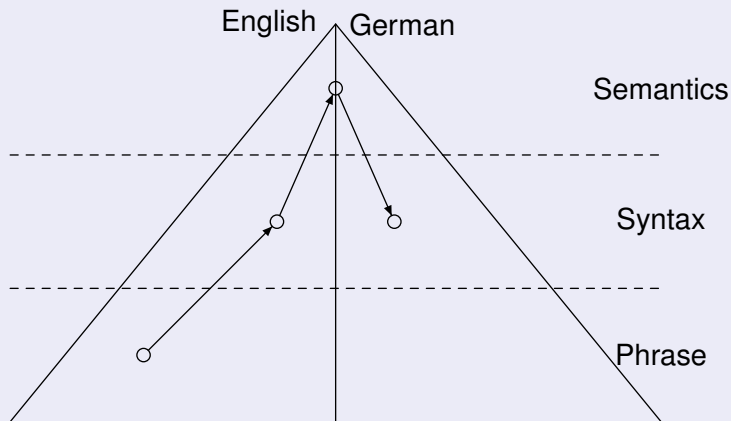
# Semantics-based Approach

## Overview



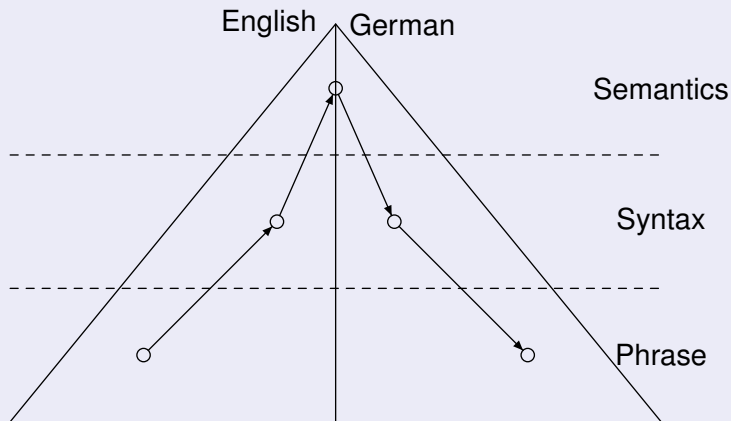
# Semantics-based Approach

## Overview



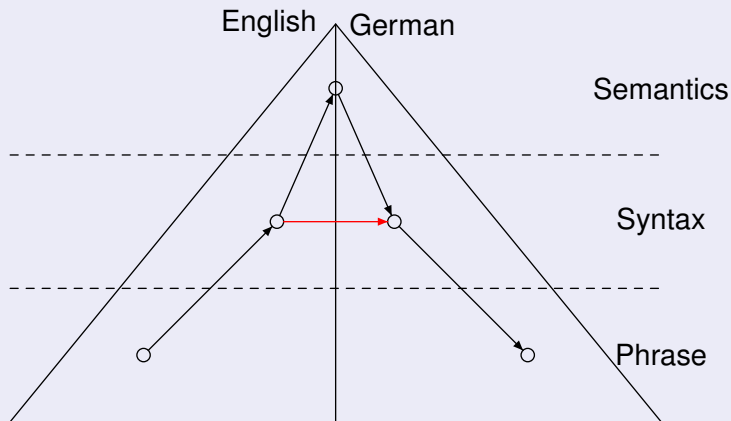
# Semantics-based Approach

## Overview



# Semantics-based Approach

## Overview



# Contents

- 1 Overview
- 2 Tree representations**
- 3 Bar-Hillel Construction
- 4 Variants
- 5 Tree Transducers



# Parsing and CFG

## Example (Context-free grammar)

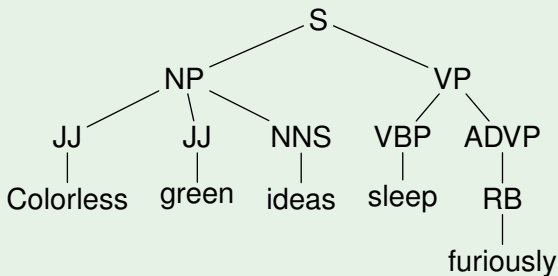
$S \rightarrow NP VP$	$NP \rightarrow JJ JJ NNS$
$VP \rightarrow VBP ADVP$	$ADVP \rightarrow RB$
$JJ \rightarrow \textit{Colorless}$	$JJ \rightarrow \textit{green}$
$NNS \rightarrow \textit{ideas}$	$VBP \rightarrow \textit{sleep}$
$RB \rightarrow \textit{furiously}$	

## Derivation

$S \rightarrow^*$  Colorless green ideas sleep furiously

# Parse tree

## Example

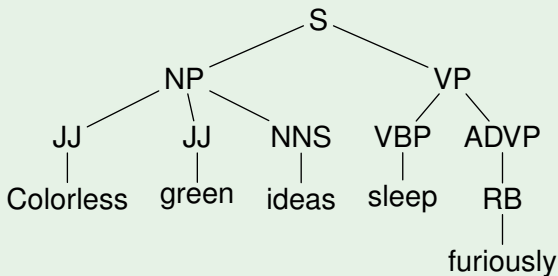


## Remark

We are interested in the parse tree, not just whether  $S \rightarrow^* w!$

# Parse tree

## Example

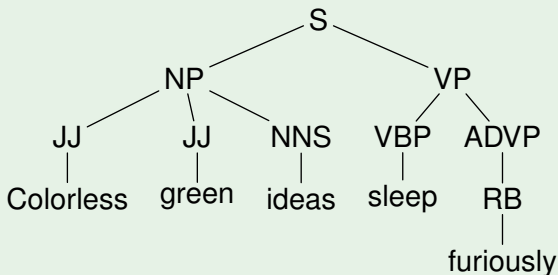


## Remark

We are interested in the parse tree, not just whether  $S \rightarrow^* w!$

# Parse tree

## Example



## Remark

We are interested in the parse tree, not just whether  $S \rightarrow^* w!$

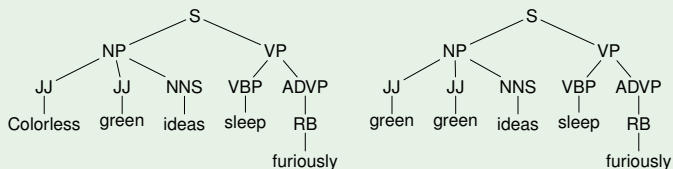
But there can be exponentially many parse trees for a sentence.

# Packed tree language

## Remark

A tree language is often called **forest** in NLP.

## Example

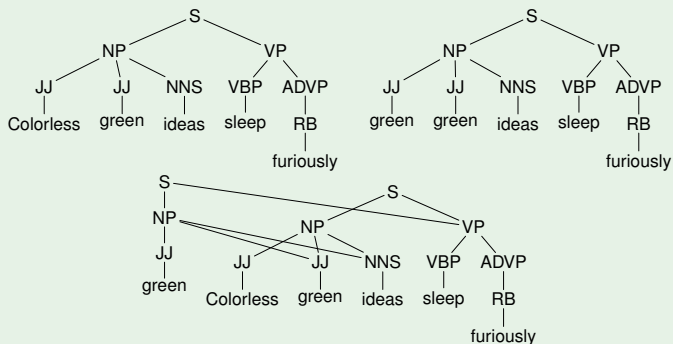


# Packed tree language

## Remark

A tree language is often called **forest** in NLP.

## Example



# Local tree language

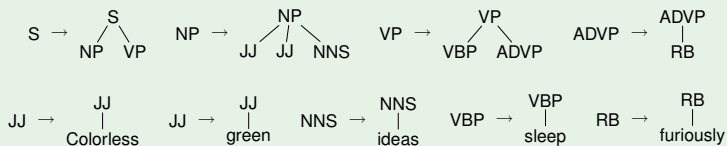
## Definition

A **local tree grammar** is a grammar with rules of the form

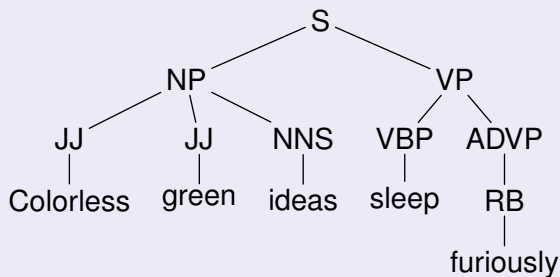
$$S \rightarrow \begin{array}{c} S \\ / \quad | \quad \backslash \\ N_1 \quad \dots \quad N_k \end{array}$$

# Local tree language

## Example



## Derivation





# Local tree language

## Definition

The tree languages generated by local tree grammars are the **local tree languages**.

## Theorem

*The set of derivations of a context-free grammar forms a local tree language.*

# Local tree language

## Definition

The tree languages generated by local tree grammars are the **local tree languages**.

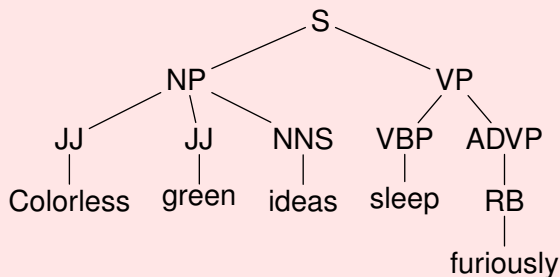
## Theorem

*The set of derivations of a context-free grammar forms a local tree language.*

# Local tree language

## Question

Is the tree language consisting of only

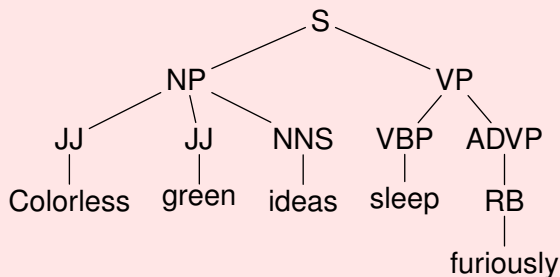


a local tree language?

# Local tree language

## Question

Is the tree language consisting of only



a local tree language?

## Answer

**NO!**

# Local tree language

## Notes

Local tree languages have undesirable properties:

- not closed under union
- cannot represent all finite languages
- ...

# Regular tree language

## Definition

A **regular tree grammar** is a grammar with rules of the form

$$q \rightarrow \begin{array}{c} S \\ / \quad | \quad \backslash \\ q_1 \quad \dots \quad q_k \end{array}$$

The such generated languages are the **regular tree languages**.

## Remark

Regular tree grammars are local tree grammars with hidden states.

# Regular tree language

## Definition

A **regular tree grammar** is a grammar with rules of the form

$$q \rightarrow \begin{array}{c} S \\ / \quad | \quad \backslash \\ q_1 \quad \dots \quad q_k \end{array}$$

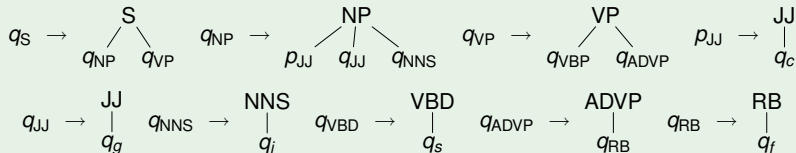
The such generated languages are the **regular tree languages**.

## Remark

Regular tree grammars are local tree grammars with hidden states.

# Regular tree language

## Example



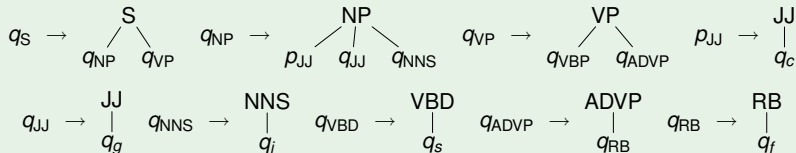
## Derivation

$q_S$



# Regular tree language

## Example

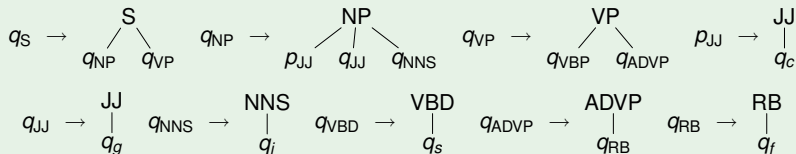


## Derivation

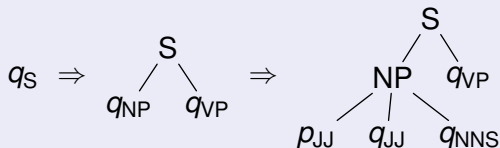


# Regular tree language

## Example



## Derivation



# Regular tree languages

## Principal properties

- Finite languages are regular
- Closed under all Boolean operations
- Closed under relabelings, linear homomorphisms, inverse homomorphisms
- Can be determinized and minimized (**bottom-up**)

## Summary

They are basically the tree version of finite-state automata with the same nice properties.

# Trees and their yield

## Definition

The **yield** of a tree is the string of its leaves (in natural order).

## Theorem

- *The yield language of a regular tree language is context-free.*
- *Each context-free language is the yield of a regular tree language.*

# Questions

## Question

Is  $\{\sigma(t, t) \mid t \text{ arbitrary tree}\}$  regular?

# Questions

## Question

Is  $\{\sigma(t, t) \mid t \text{ arbitrary tree}\}$  regular?

## Answer

NO!

# Questions

## Question

Is  $\{\sigma(t, t) \mid t \text{ arbitrary tree}\}$  regular?

## Answer

NO!

## Question

Is  $\{\sigma(\gamma^n(\alpha), \gamma^n(\alpha)) \mid n \in \mathbb{N}\}$  regular?

# Questions

## Question

Is  $\{\sigma(t, t) \mid t \text{ arbitrary tree}\}$  regular?

## Answer

NO!

## Question

Is  $\{\sigma(\gamma^n(\alpha), \gamma^n(\alpha)) \mid n \in \mathbb{N}\}$  regular?

## Answer

NO!

## Remark

Not every tree language with context-free yield language is regular!



# Back to parsing

## Observation

Most CFG-parsers are regular tree grammars (+ control) because

- they are based on a CFG ( $\rightarrow$  local tree grammar) and
- have hidden states (or features)

## Alternative

The features can be made explicit in the parse tree structure.

# Back to parsing

## Observation

Most CFG-parsers are regular tree grammars (+ control) because

- they are based on a CFG ( $\rightarrow$  local tree grammar) and
- have hidden states (or features)

## Alternative

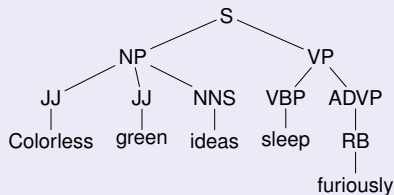
The features can be made explicit in the parse tree structure.

# Contents

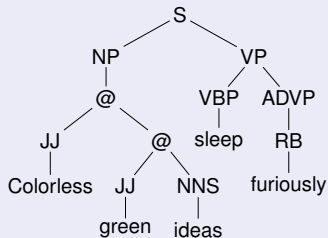
- 1 Overview
- 2 Tree representations
- 3 Bar-Hillel Construction**
- 4 Variants
- 5 Tree Transducers

# Binarization

## Input tree



## Binarized tree

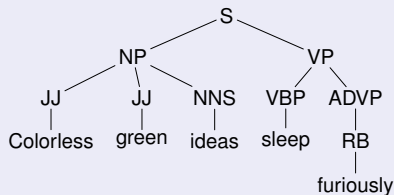


## Theorem

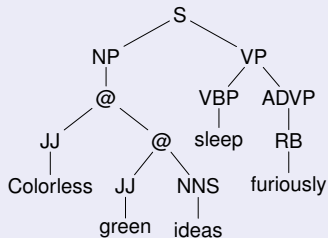
A tree language is regular if and only if its binarization is regular

# Binarization

## Input tree



## Binarized tree



## Theorem

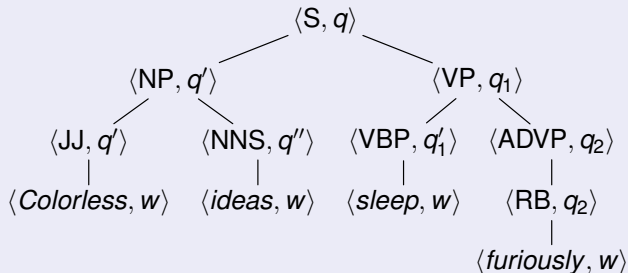
A tree language is regular if and only if its binarization is regular

# Individual runs

## Run on the yield

$\langle p \rangle$  Colorless  $\langle p_1 \rangle$  ideas  $\langle p_2 \rangle$  sleep  $\langle p_3 \rangle$  furiously  $\langle p' \rangle$

## Run on the input tree

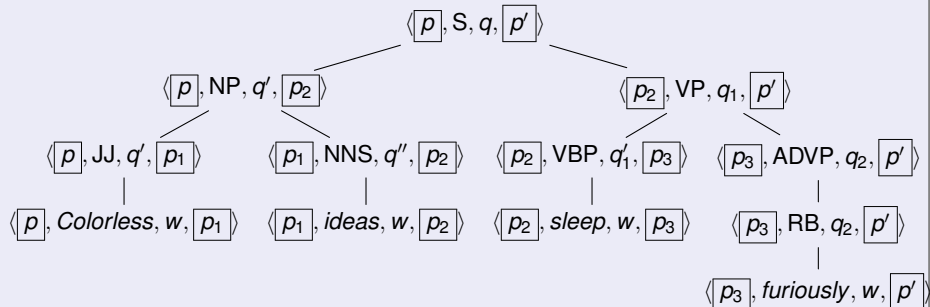


# Bar-Hillel construction

## Run on the yield

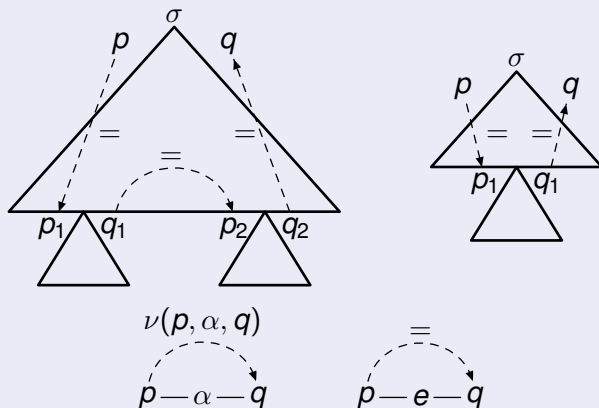
( $p$ ) Colorless ( $p_1$ ) ideas ( $p_2$ ) sleep ( $p_3$ ) furiously ( $p'$ )

## Composite run



# Bar-Hillel construction (cont'd)

## Illustration





## Bar-Hillel construction (cont'd)

### Theorem

*The regular restriction of a regular tree language is regular*

### Remark

Complexity:  $O(mn^3)$

- $m$ : size of the regular tree grammar
- $n$ : size of the regular grammar (or input string)

### Conclusion

We can parse with regular tree grammars in  $O(mn^3)$ .

# Contents

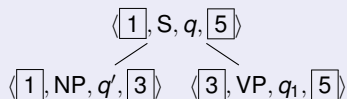
- 1 Overview
- 2 Tree representations
- 3 Bar-Hillel Construction
- 4 Variants**
- 5 Tree Transducers

# Top-down Bar-Hillel construction

## Run on the yield

(1) Colorless (2) ideas (3) sleep (4) furiously (5)

## Composite run

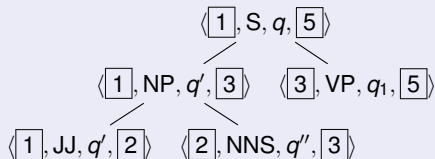


# Top-down Bar-Hillel construction

## Run on the yield

(1) Colorless (2) ideas (3) sleep (4) furiously (5)

## Composite run

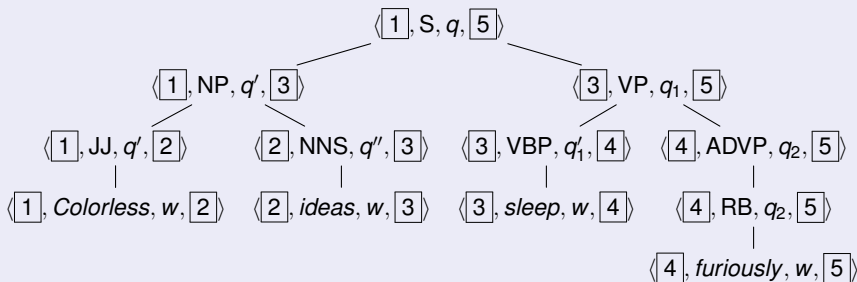


# Top-down Bar-Hillel construction

## Run on the yield

(1) Colorless (2) ideas (3) sleep (4) furiously (5)

## Composite run



# Bottom-up Bar-Hillel construction

## Run on the yield

(1) Colorless (2) ideas (3) sleep (4) furiously (5)

## Composite run

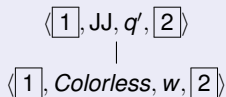
$\langle \boxed{1}, \text{Colorless}, w, \boxed{2} \rangle$

# Bottom-up Bar-Hillel construction

## Run on the yield

(1) Colorless (2) ideas (3) sleep (4) furiously (5)

## Composite run



# Bottom-up Bar-Hillel construction

## Run on the yield

(1) Colorless (2) ideas (3) sleep (4) furiously (5)

## Composite run



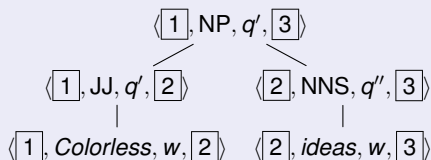


# Bottom-up Bar-Hillel construction

## Run on the yield

(1) Colorless (2) ideas (3) sleep (4) furiously (5)

## Composite run

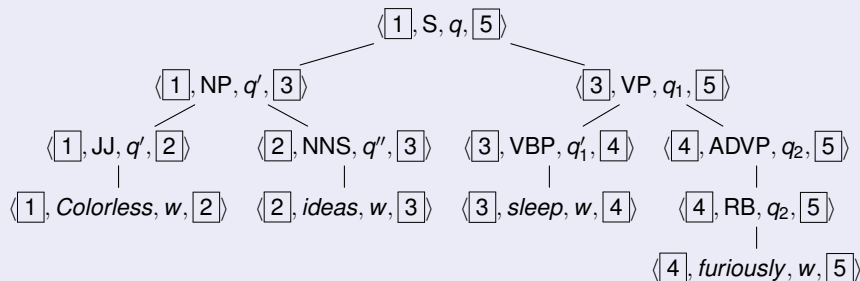


# Bottom-up Bar-Hillel construction

## Run on the yield

(1) Colorless (2) ideas (3) sleep (4) furiously (5)

## Composite run



# Summary

## Key points

- regular tree grammar as efficient tree data structure
- context-free behavior
- basis for most syntax-based translation models

## Further models

- weaker models are generally inadequate
- tree adjoining grammars (more expressive, but worse computational properties)
- Automata on directed acyclic graphs (see Daniel's lecture)

# Contents

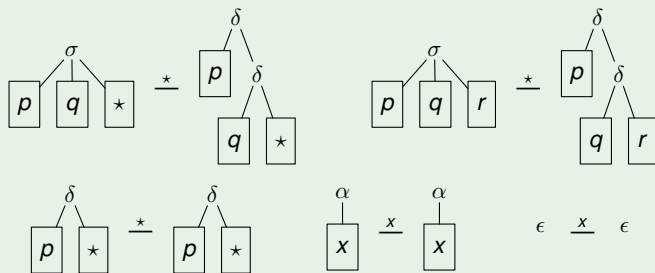
- 1 Overview
- 2 Tree representations
- 3 Bar-Hillel Construction
- 4 Variants
- 5 Tree Transducers**

# Extended Top-down Tree Transducer

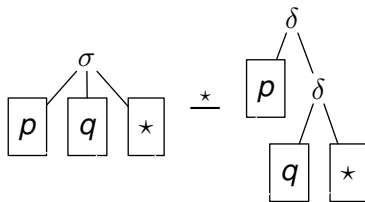
## Definition

Each rule now has an input and an output side, which are both full trees (not just a single symbol followed by states)

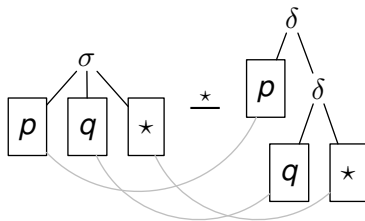
## Example



# Link Structure

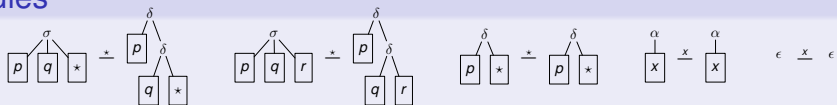


# Link Structure



# Derivation

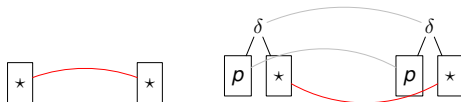
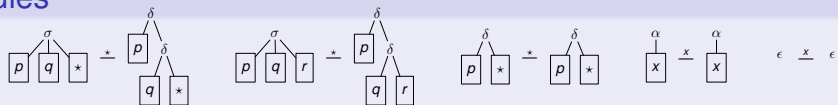
## Rules





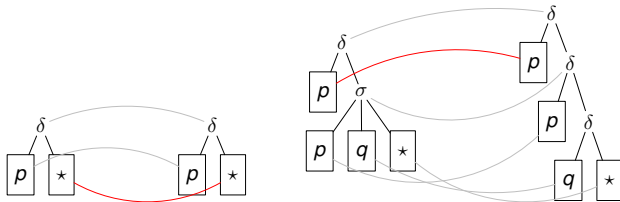
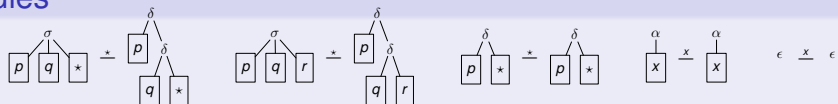
# Derivation

## Rules



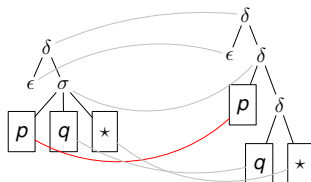
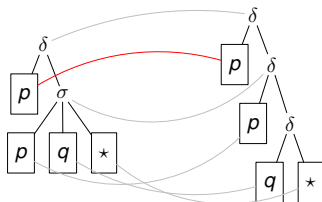
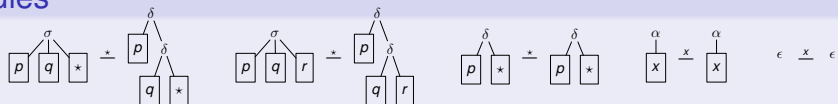
# Derivation

## Rules



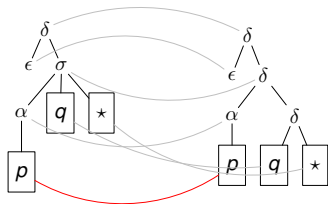
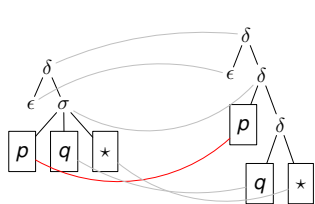
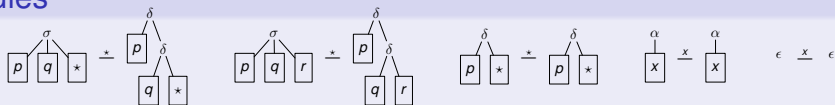
# Derivation

## Rules



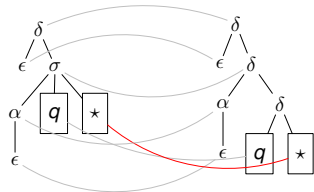
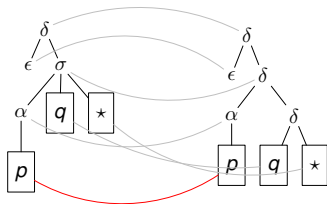
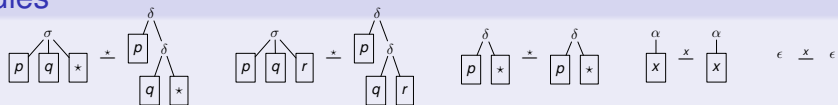
# Derivation

## Rules



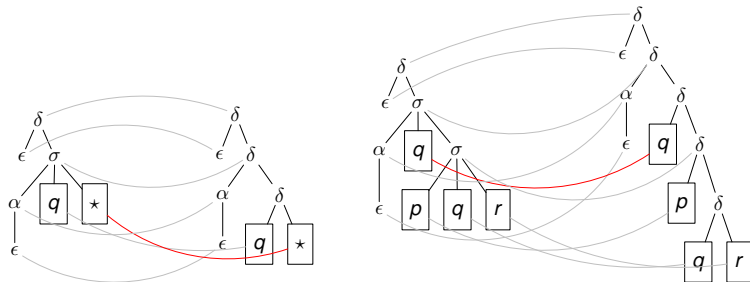
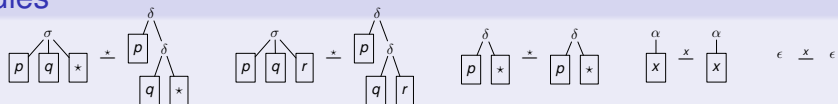
# Derivation

## Rules



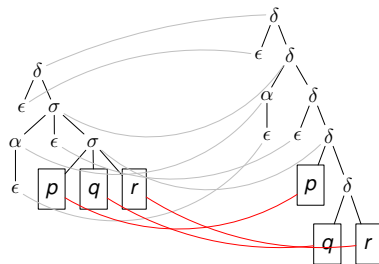
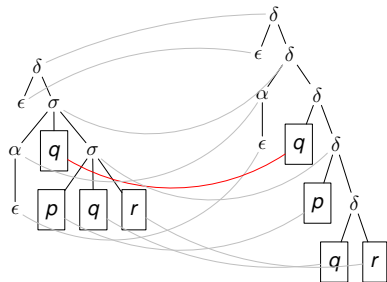
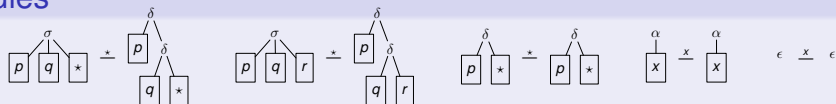
# Derivation

## Rules



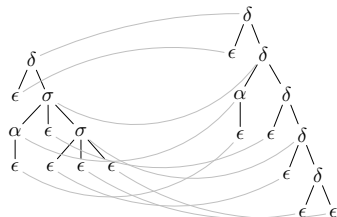
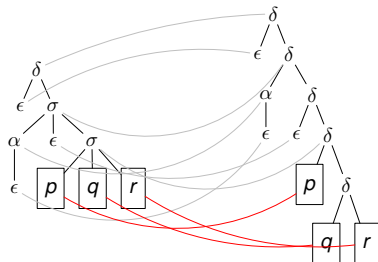
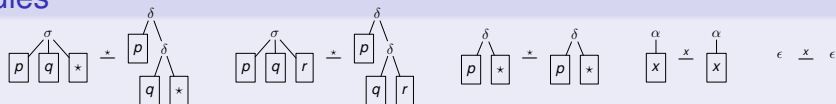
# Derivation

## Rules



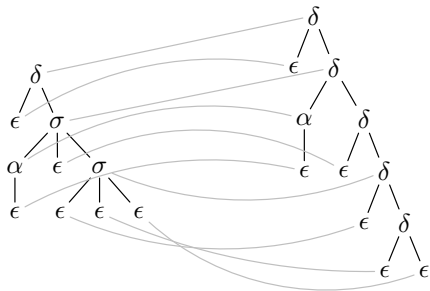
# Derivation

## Rules

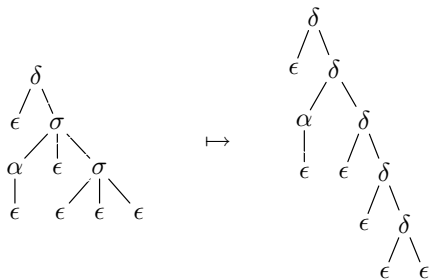




# Semantics of XTOP



# Semantics of XTOP



# References

- [Bar-Hillel, Perles, Shamir](#): *On formal properties of simple phrase structure grammars*. Language and Information, 1964
- [Berstel, Reutenauer](#): Recognizable formal power series on trees. *Theoret. Comput. Sci.* 18, p. 115–148, 1982
- [Borchardt](#): The theory of recognizable tree series. Ph.D. thesis TU Dresden, 2005
- [Gécseg, Steinby](#): *Tree Automata*. Akadémiai Kiadó, Budapest 1984
- [Nederhof, Satta](#): Probabilistic parsing as intersection. In IWPT 2003

Thank you for your attention!