

Weighted Automata in Statistical Machine Translation

Andreas Maletti

Institute for Natural Language Processing
University of Stuttgart

WATA — April 26, 2016

Review translation [by Google Translate]

- 1 The room it is not narrowly was a simple, bathtub was also attached.
- 2 Wi-fi, TV and I was available.
- 3 Church looked When morning awake open the curtain.
- 4 When looking at often, wives, went out and is invited to try to go [...].
- 5 But was a little cold, morning walks was good.

Review translation [by Google Translate]

- 1 The room it is not narrowly was a simple, bathtub was also attached.
- 2 Wi-fi, TV and I was available.
- 3 Church looked When morning awake open the curtain.
- 4 When looking at often, wives, went out and is invited to try to go [...].
- 5 But was a little cold, morning walks was good.

Original [Japanese — © tripadvisor]

- 1 部屋もシンプルでしたが狭くなく、バスタブもついていました。
- 2 Wi-fi、テレビも利用出来ました。
- 3 朝起きてカーテンを開けると教会が見えました。
- 4 しばし眺めていると、妻たちは、[...]るから行こうとさそわれ出かけました。
- 5 ちょっと寒かったけれど、朝の散策はグッドでしたよ。

Sample translation [by phrase-based Moses]

- 1 I think Danish is a hard language, though it looks like German.
- 2 Fortunately talking almost all Danes English, especially the young.
- 3 The boys come too late, but the girls come on time.

Danish-to-English Translation

Sample translation [by phrase-based Moses]

- 1 I think Danish is a hard language, though it looks like German.
- 2 Fortunately talking almost all Danes English, especially the young.
- 3 The boys come too late, but the girls come on time.

Original Danish

- 1 Jeg synes at dansk er et svært sprog, selvom det ligner tysk.
- 2 Heldigvis snakker næsten alle danskere engelsk, især de unge.
- 3 Drengene kom for sent, men pigerne kom til tiden.

Timeline

1960 ● Dark age

- rule-based systems (e.g.,  SYSTRAN)
- Chomskyan approach (perfect translation, poor coverage)

1991 ● Reformation

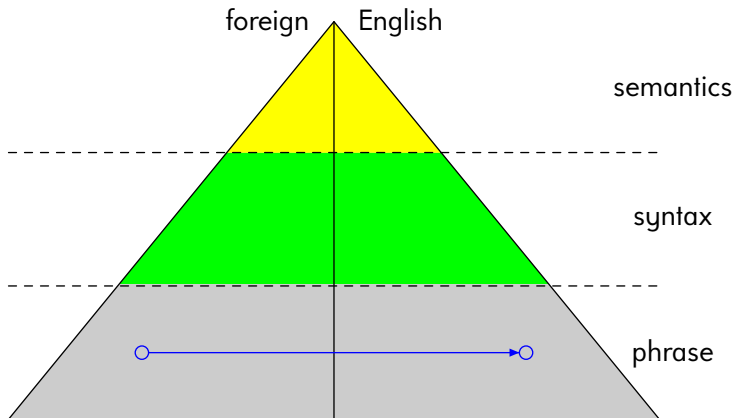
- phrase-based and syntax-based systems
- statistical approach (cheap, automatically trained)

2016 ● Potential future

- semantics-based systems (e.g., FrameNet-based)
- semi-supervised, statistical approach
- basic understanding of (translated) text

Machine Translation

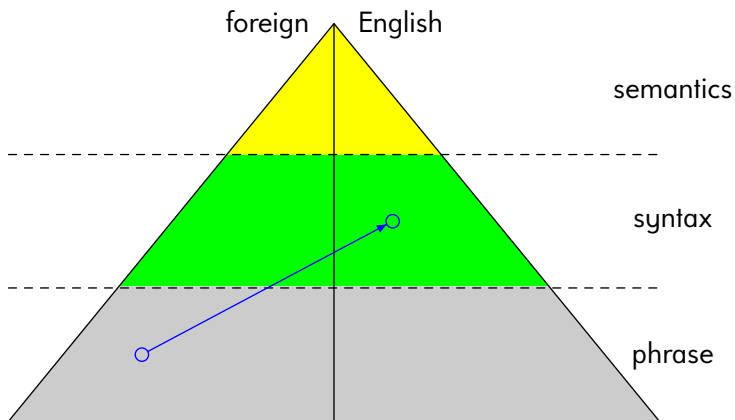
Vauquois triangle:



Translation model: **string-to-string**

Machine Translation

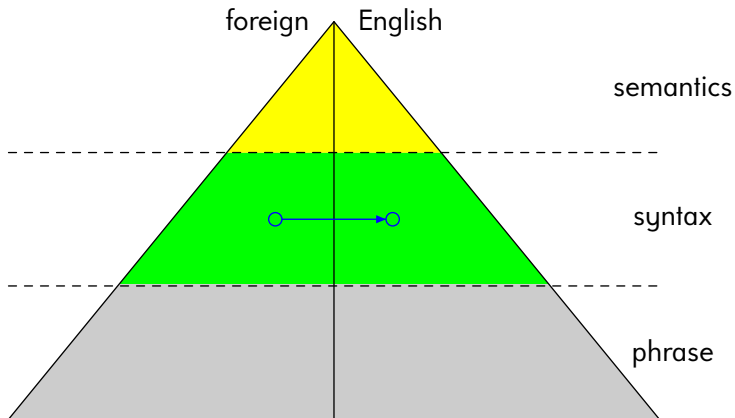
Vauquois triangle:



Translation model: [string-to-tree](#)

Machine Translation

Vauquois triangle:



Translation model: [tree-to-tree](#)

Training data

- parallel corpus
- word alignments
- parse trees

(for syntax-based systems)

Training data

- parallel corpus
- word alignments
- parse trees

(for syntax-based systems)

Parallel corpus

linguistic resource containing (sentence-by-sentence) example translations

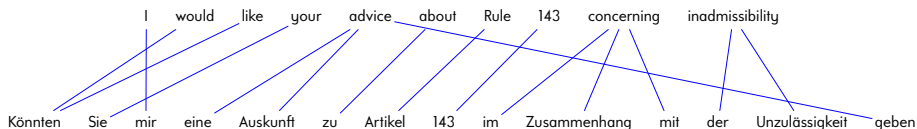
parallel corpus, word alignments, parse tree

I would like your advice about Rule 143 concerning inadmissibility

Könnten Sie mir eine Auskunft zu Artikel 143 im Zusammenhang mit der Unzulässigkeit geben

Machine Translation

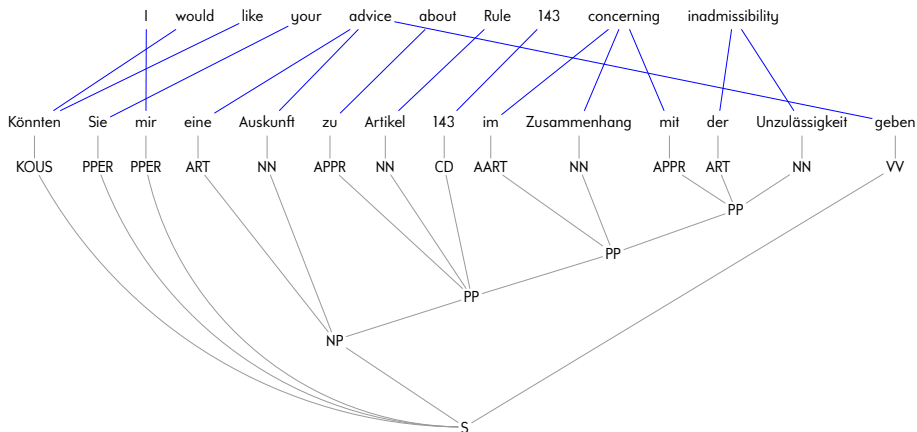
parallel corpus, **word alignments**, parse tree



via GIZA++ [Och, Ney: A systematic comparison of various statistical alignment models. *Computational Linguistics* 29(1), 2003]

Machine Translation

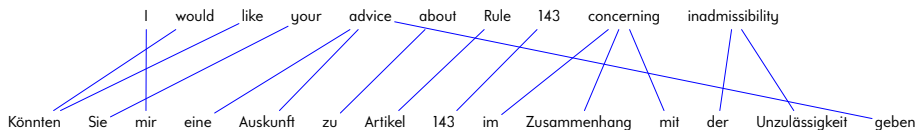
parallel corpus, word alignments, **parse tree**



via Berkeley parser [Petrov, Barrett, Thibaux, Klein: Learning accurate, compact, and interpretable tree annotation. *Proc. ACL*, 2006]

Phrase-based Model

Training example:

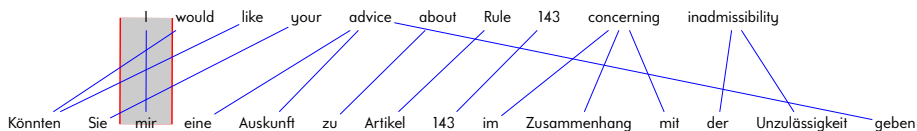


Extracted rules:

I	—	mir	would like	—	Könnten
your	—	Sie	about	—	zu
Rule	—	Artikel	143	—	143
concerning	—	im Zusammenhang mit	about Rule	—	zu Artikel
inadmissibility	—	der Unzulässigkeit			

Phrase-based Model

Training example:

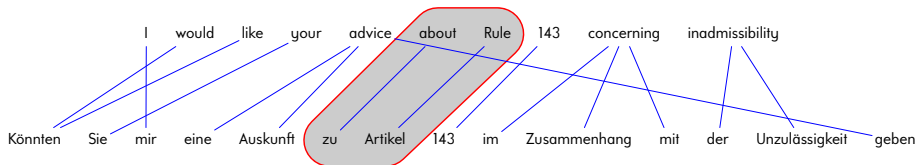


Extracted rules:

I	—	mir	would like	—	Könnten
your	—	Sie	about	—	zu
Rule	—	Artikel	143	—	143
concerning	—	im Zusammenhang mit	about Rule	—	zu Artikel
inadmissibility	—	der Unzulässigkeit			

Phrase-based Model

Training example:



Extracted rules:

I	—	mir	would like	—	Könnten
your	—	Sie	about	—	zu
Rule	—	Artikel	143	—	143
concerning	—	im Zusammenhang mit	about Rule	—	zu Artikel
inadmissibility	—	der Unzulässigkeit			

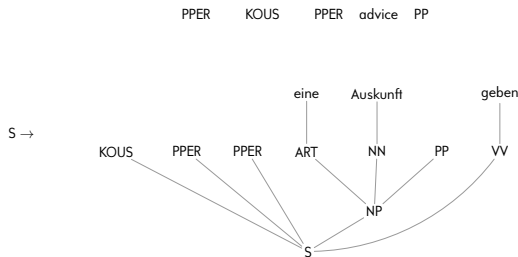
Notes

- essentially weighted finite-state transducer
- weights estimated using maximum likelihood

Weighted Synchronous Grammars

Synchronous tree substitution grammar: productions $N \rightarrow (r, r_1)$

- nonterminal N
- right-hand side r of context-free grammar production
- right-hand side r_1 of tree substitution grammar production

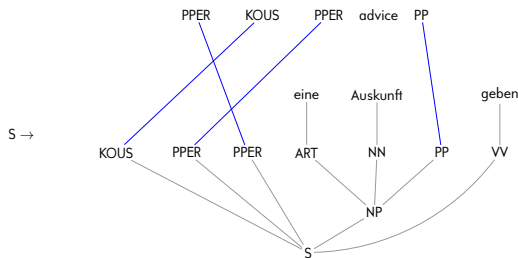


variant of [M., Graehl, Hopkins, Knight: The power of extended top-down tree transducers. *SIAM Journal on Computing* 39(2), 2009]

Weighted Synchronous Grammars

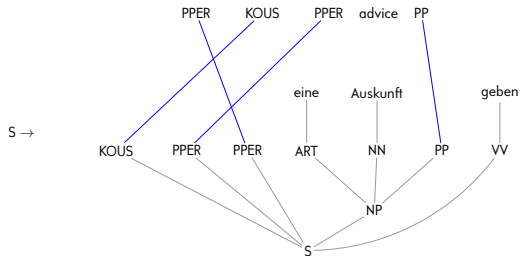
Synchronous tree substitution grammar: productions $N \rightarrow (r, r_1)$

- nonterminal N
- right-hand side r of context-free grammar production
- right-hand side r_1 of tree substitution grammar production
- (bijective) synchronization of nonterminals



variant of [M., Graehl, Hopkins, Knight: The power of extended top-down tree transducers. *SIAM Journal on Computing* 39(2), 2009]

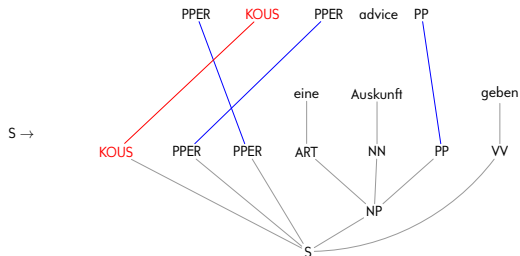
Synchronous Grammars



Production application

- 1 Selection of synchronous nonterminals

Synchronous Grammars



Production application

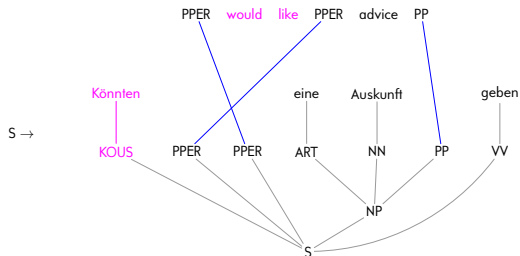
- 1 Selection of synchronous nonterminals
- 2 Selection of suitable production

KOUS →

would like

Könnten
KOUS

Synchronous Grammars



Production application

- 1 Selection of synchronous nonterminals
- 2 Selection of suitable production
- 3 Replacement on both sides

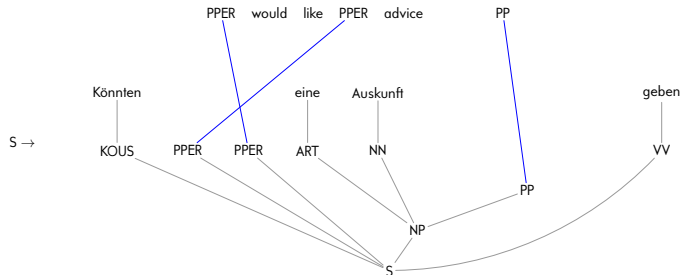
KOU5 →

would like

Könnten

KOU5

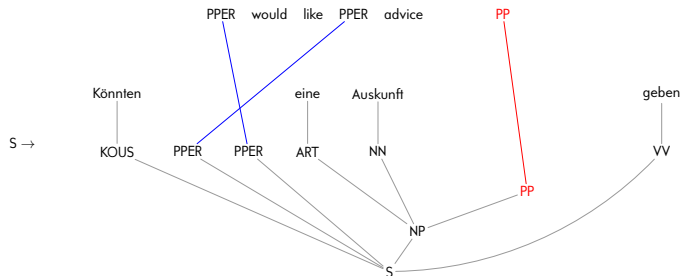
Synchronous Grammars



Production application

- 1 synchronous nonterminals

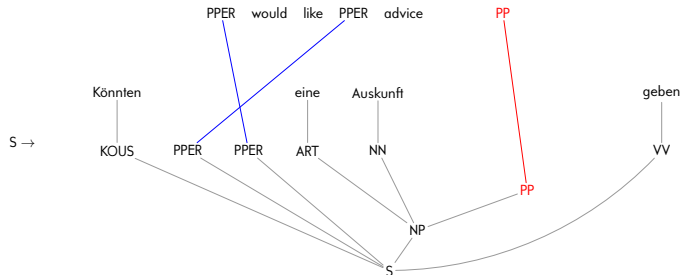
Synchronous Grammars



Production application

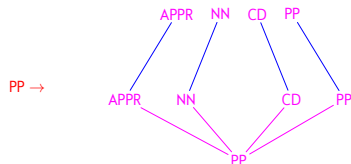
- 1 synchronous nonterminals

Synchronous Grammars

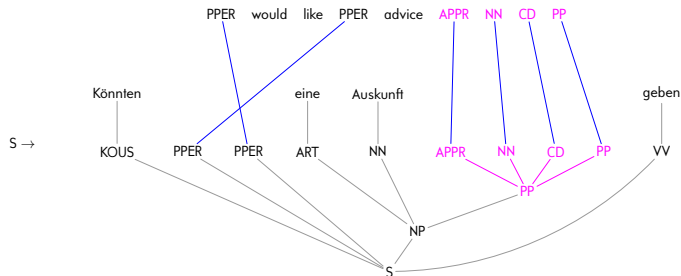


Production application

- 1 synchronous nonterminals
- 2 suitable production

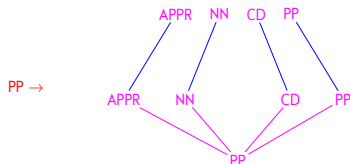


Synchronous Grammars



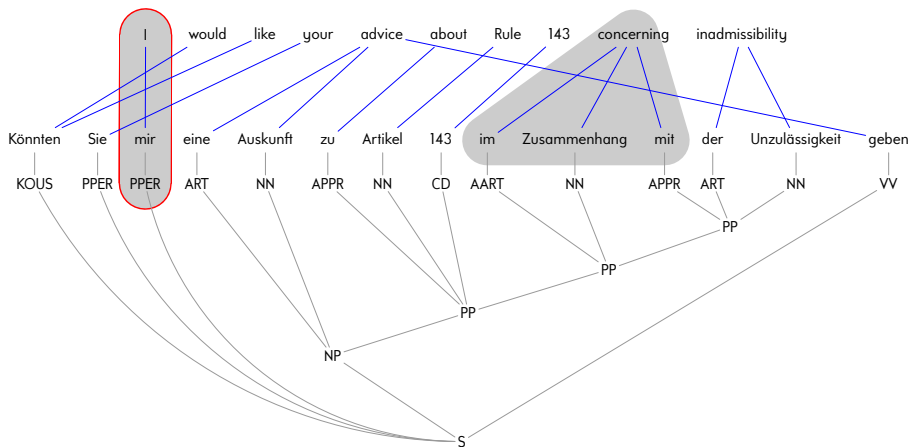
Production application

- 1 synchronous nonterminals
- 2 suitable production
- 3 replacement



Production Extraction

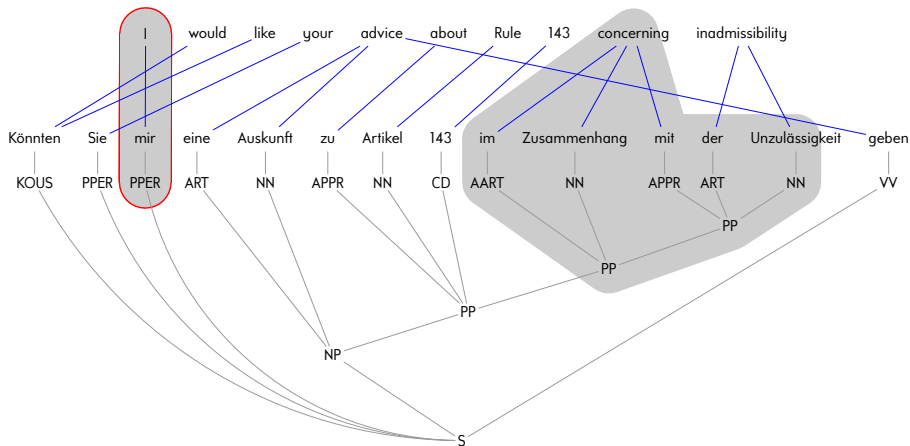
(extractable productions marked in red)



following [Galley, Hopkins, Knight, Marcu: What's in a translation rule? *Proc. NAACL*, 2004]

Production Extraction

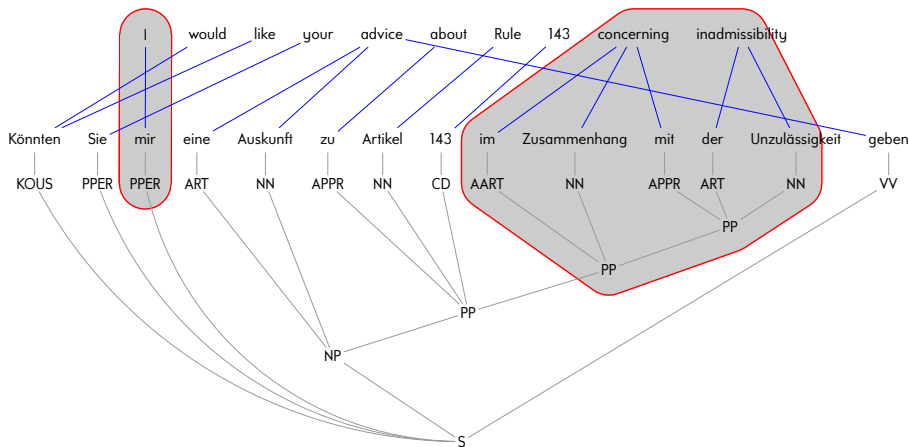
(extractable productions marked in red)



following [Galley, Hopkins, Knight, Marcu: What's in a translation rule? *Proc. NAACL*, 2004]

Production Extraction

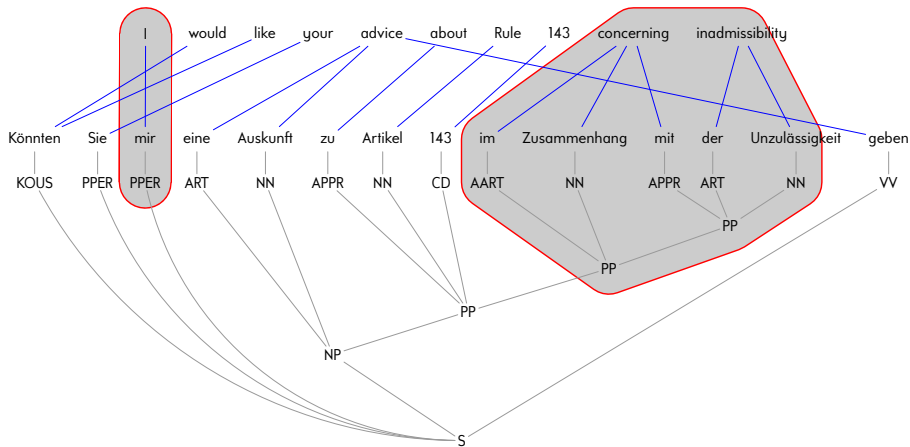
(extractable productions marked in red)



following [Galley, Hopkins, Knight, Marcu: What's in a translation rule? *Proc. NAACL*, 2004]

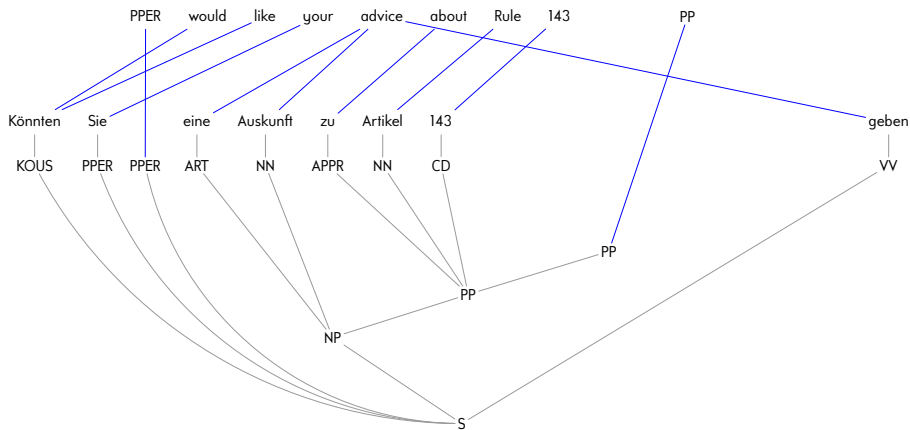
Production Extraction

Removal of extractable production:



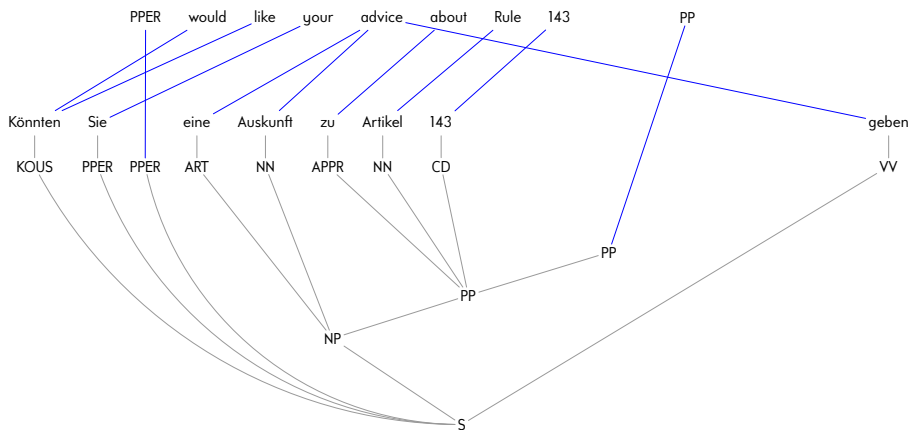
Production Extraction

Removal of extractable production:



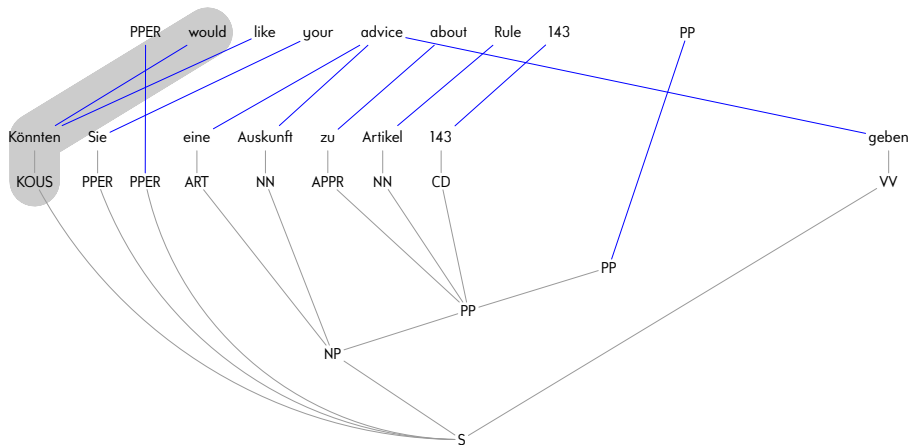
Production Extraction

Repeated production extraction:



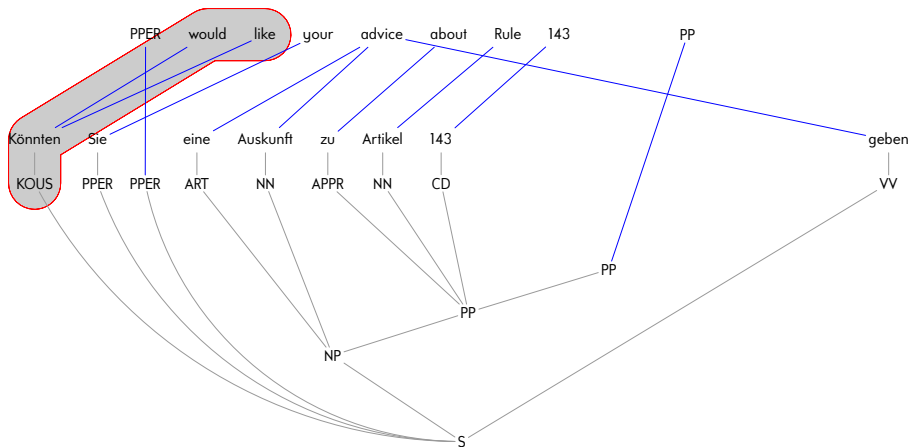
Production Extraction

Repeated production extraction:



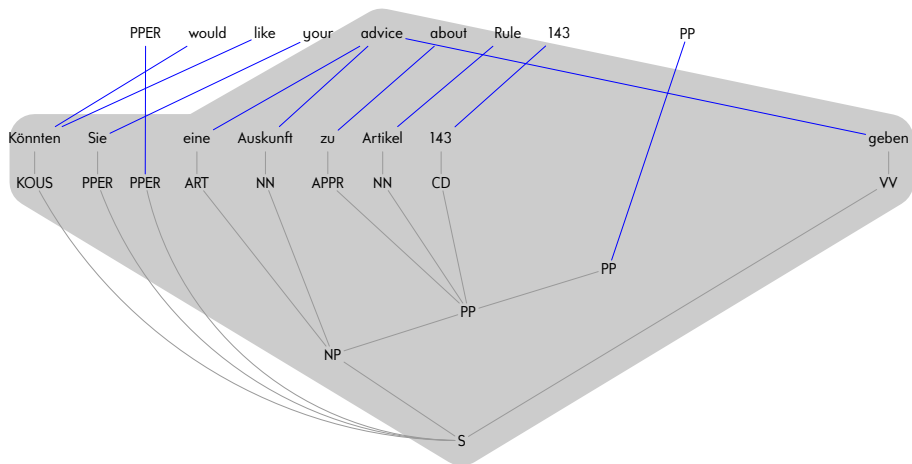
Production Extraction

Repeated production extraction: (extractable productions marked in red)



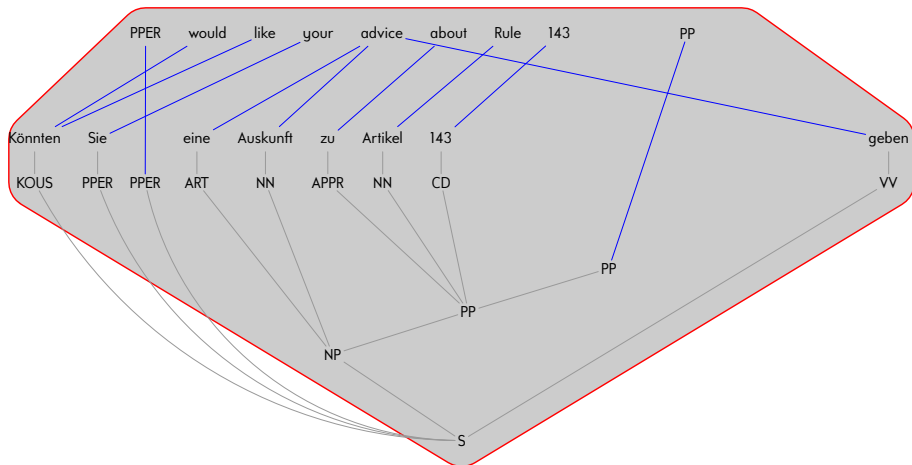
Production Extraction

Repeated production extraction: (extractable productions marked in red)



Production Extraction

Repeated production extraction: (extractable productions marked in red)



Advantages

- very simple
- implemented in framework ‘Moses’
[[Koehn](#) et al.: *Moses* — Open source toolkit for statistical machine translation. *Proc. ACL*, 2007]
- “context-free”

Advantages

- very simple
- implemented in framework ‘Moses’
[Koehn et al.: Moses — Open source toolkit for statistical machine translation. *Proc. ACL*, 2007]
- “context-free”

Disadvantages

- problems with discontinuities
- composition and binarization not possible
[M., Graehl, Hopkins, Knight: The power of extended top-down tree transducers. *SIAM Journal on Computing* 39(2), 2009]
[Zhang, Huang, Gildea, Knight: Synchronous Binarization for Machine Translation. *Proc. NAACL*, 2006]
- “context-free”

English → German translation task:

(higher BLEU is better)

Type	System	BLEU		
		vanilla	WMT 2013	WMT 2015
string-to-string	phrase-based	16.7	20.3	23.3
	hierarchical	17.0	—	—
string-to-tree	STSG	15.2	19.4	24.5
tree-to-tree	STSG	14.5	—	15.3

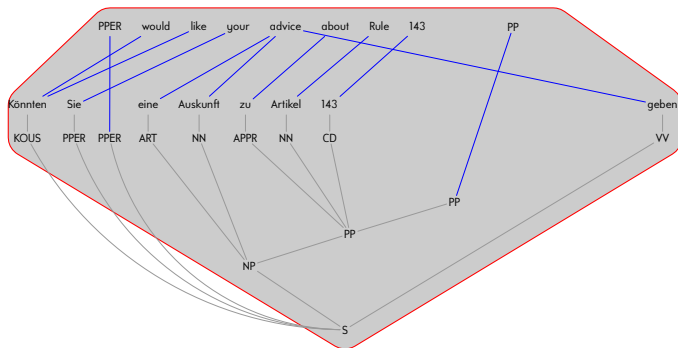
from [Seemann, Braune, M.: A systematic evaluation of MBOT in statistical machine translation. *Proc. MT-Summit*, 2015]
and [Bojar et al.: Findings of the 2013 workshop on statistical machine translation. *Proc. WMT*, 2013]
and [Bojar et al.: Findings of the 2015 workshop on statistical machine translation. *Proc. WMT*, 2015]

Observation

- syntax-based systems competitive with manual adjustments
- much less so for vanilla systems
- very unfortunate situation [more supervision yields lower scores]

- 1 Background
- 2 Extending the Expressive Power
- 3 Investigating their Expressive Power

Production Extraction



- very specific production
- every production for 'advice' contains sentence structure
(syntax "in the way")

Synchronous Grammars

Synchronous multi tree substitution grammar: $N \rightarrow (r, \langle r_1, \dots, r_n \rangle)$

variant of [M: Why synchronous tree substitution grammars?. *Proc. NAACL*, 2010]

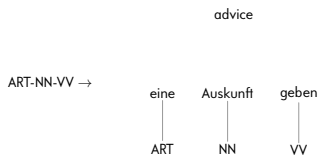
- nonterminal N
- right-hand side r of context-free grammar production
- right-hand **sides** r_1, \dots, r_n of regular tree grammar production

Synchronous Grammars

Synchronous multi tree substitution grammar: $N \rightarrow (r, \langle r_1, \dots, r_n \rangle)$

variant of [M.: Why synchronous tree substitution grammars?. *Proc. NAACL*, 2010]

- nonterminal N
- right-hand side r of context-free grammar production
- right-hand **sides** r_1, \dots, r_n of regular tree grammar production

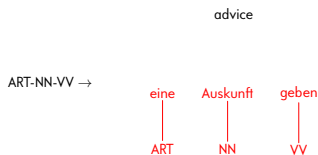


Synchronous Grammars

Synchronous multi tree substitution grammar: $N \rightarrow (r, \langle r_1, \dots, r_n \rangle)$

variant of [M.: Why synchronous tree substitution grammars?. *Proc. NAACL*, 2010]

- nonterminal N
- right-hand side r of context-free grammar production
- right-hand **sides** r_1, \dots, r_n of regular tree grammar production



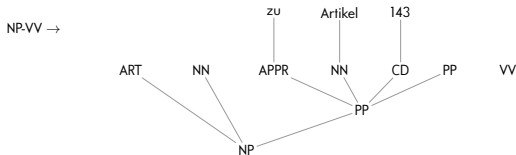
Synchronous Grammars

Synchronous multi tree substitution grammar: $N \rightarrow (r, \langle r_1, \dots, r_n \rangle)$

variant of [M: Why synchronous tree substitution grammars?. *Proc. NAACL*, 2010]

- nonterminal N
- right-hand side r of context-free grammar production
- right-hand **sides** r_1, \dots, r_n of regular tree grammar production

ART-NN-VV about Rule 143 PP

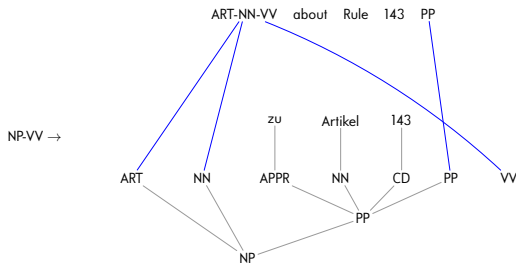


Synchronous Grammars

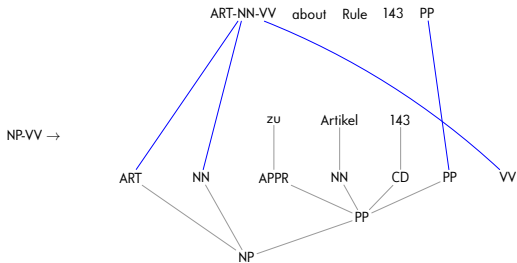
Synchronous multi tree substitution grammar: $N \rightarrow (r, \langle r_1, \dots, r_n \rangle)$

variant of [M.: Why synchronous tree substitution grammars?. *Proc. NAACL*, 2010]

- nonterminal N
- right-hand side r of context-free grammar production
- right-hand **sides** r_1, \dots, r_n of regular tree grammar production
- synchronization via map NT r_1, \dots, r_n to NT r



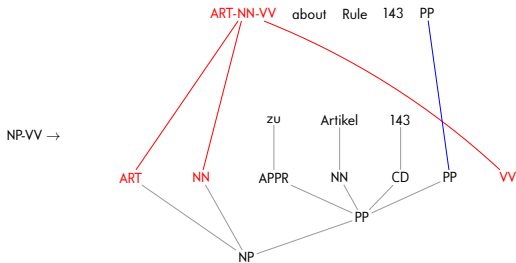
Synchronous Grammars



Production application

- 1 synchronous nonterminals

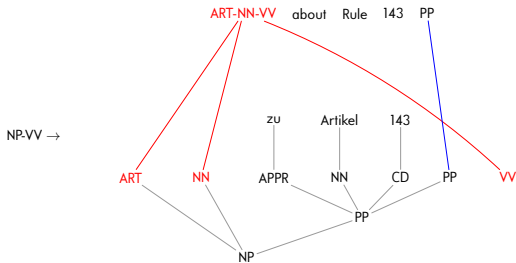
Synchronous Grammars



Production application

- 1 synchronous nonterminals

Synchronous Grammars



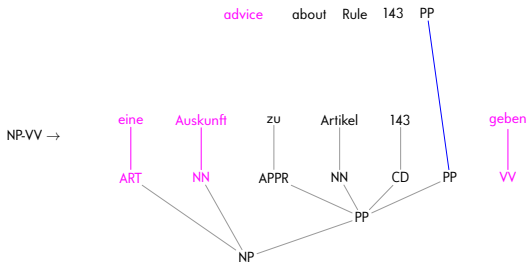
Production application

- 1 synchronous nonterminals
- 2 suitable production

ART-NN-VV →

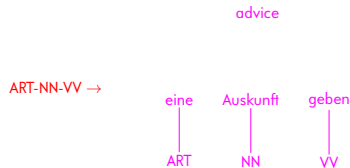


Synchronous Grammars



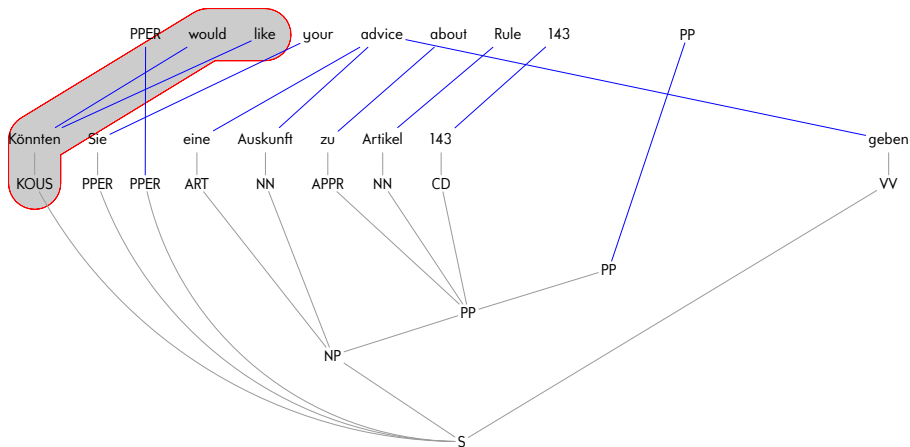
Production application

- 1 synchronous nonterminals
- 2 suitable production
- 3 replacement



Production Extraction

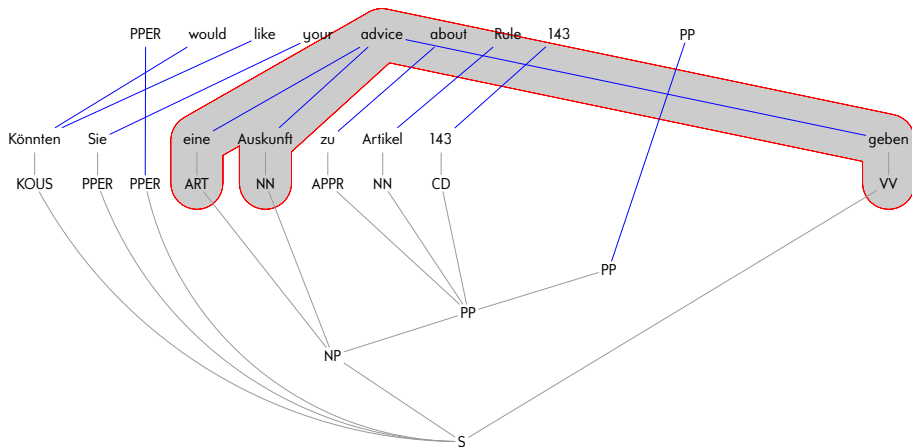
(extractable productions marked in red)



variant of [M.: How to train your multi bottom-up tree transducer. *Proc. ACL*, 2011]

Production Extraction

(extractable productions marked in red)



variant of [M.: How to train your multi bottom-up tree transducer. *Proc. ACL*, 2011]

Advantages

- complicated discontinuities
- implemented in framework 'Moses'
- binarizable, composable

[[Braune, Seemann, Quernheim, M.](#): Shallow local multi bottom-up tree transducers in SMT. *Proc. ACL*, 2013]

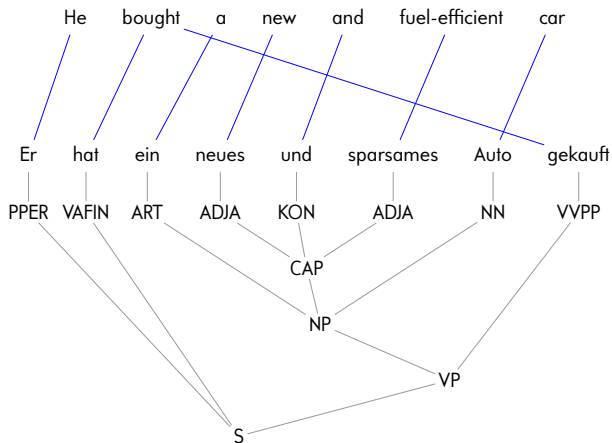
Advantages

- complicated discontinuities
- implemented in framework 'Moses'
[Braune, Seemann, Quernheim, M.: Shallow local multi bottom-up tree transducers in SMT. *Proc. ACL*, 2013]
- binarizable, composable

Disadvantages

- output non-regular (trees) or non-context-free (strings)
- not symmetric (input context-free; output not)

Discontinuity



System	Number of productions		
	E.-to-German	E.-to-Arabic	E.-to-Chinese
t-to-t STSG	7M	24M	8M
t-to-t SMTSG	41M	151M	84M
s-to-t STSG	14M	55M	17M
s-to-t SMTSG	144M	491M	162M
phrase-based	406M	842M	209M
s-to-s SMTSG	1,084M	2,208M	683M

from [Seemann, Braune, M.: A systematic evaluation of MBOT in statistical machine translation. *Proc. MT-Summit*, 2015]

String-to-tree systems vs. phrase-based:

Task	BLEU		
	STSG	SMTSG	phrase-based
English → German	15.0	*15.5	16.8
English → Arabic	48.2	*49.1	51.9
English → Chinese	17.7	*18.4	18.1
English → Polish	21.3	*23.4	24.4
English → Russian	24.7	*26.1	27.9

from [Seemann, Braune, M.: A systematic evaluation of MBOT in statistical machine translation. *Proc. MT-Summit*, 2015]
and [Seemann, M.: Discontinuous statistical machine translation with target-side dependency syntax. *Proc. WMT*, 2015]

Conclusions

- consistent improvements
- 1 magnitude more productions
- SMTSG alleviate some of the problems of syntax-based systems

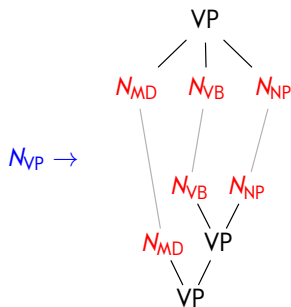
- 1 Background
- 2 Extending the Expressive Power
- 3 Investigating their Expressive Power

Synchronous Grammars

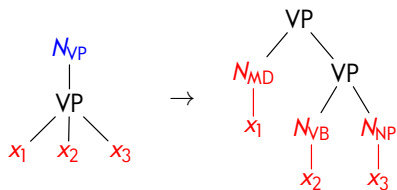
Notes

- tree-to-tree models easier for theoretical investigation
- strongly related to tree transducers
- we disallow trivial input sides of just a nonterminal (ϵ -free)

Synchronous grammar:



Tree transducer:



Synchronous Grammars

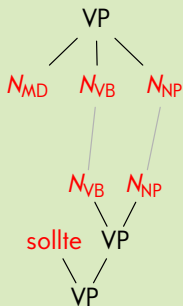
Major linear tree transducers:

synchronization	bijjective	injective
input sides		(output \rightarrow input)
shallow	nondeleting top-down ...	top-down ...
general	nondeleting extended ...	extended ...

Further distinction

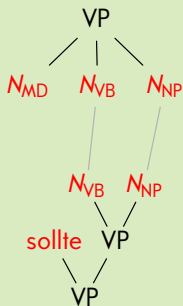
- allow productions on disconnected input nonterminals
→ regular look-ahead
- allow arbitrary trees for disconnected input nonterminals
→ no look-ahead

Illustration



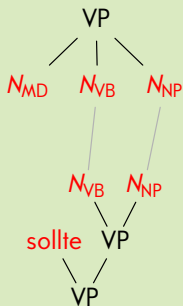
- **no look-ahead:** can plug any (terminal) tree for N_{MD} [e.g., NP(DT(the), NN(tower))]

Illustration



- **no look-ahead:** can plug any (terminal) tree for N_{MD}
[e.g., NP(DT(the), NN(tower))]
- **regular look-ahead:** use special “no-output”-productions $N \rightarrow (r)$
[e.g., $N_{MD} \rightarrow (MD(\text{should}))$]

Illustration



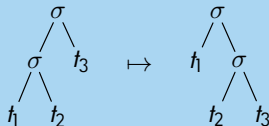
- **no look-ahead:** can plug any (terminal) tree for N_{MD}
[e.g., NP(DT(the), NN(tower))]
- **regular look-ahead:** use special “no-output”-productions $N \rightarrow (r)$
[e.g., $N_{MD} \rightarrow (MD(\text{should}))$]
- SMTSG always have regular look-ahead (any number of components includes 0)

Evaluation criteria



rotations implementable?

(for arbitrary t_1, t_2, t_3)



symmetric?



domain regular?



range regular?



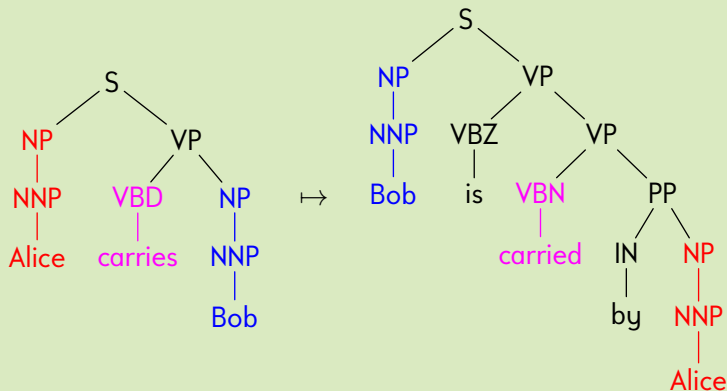
closed under composition?

following [Knight: Capturing practical natural language transformations. *Machine Translation* 21(2), 2007]
and [May, Knight, Vogler: Efficient inference through cascades of weighted tree transducers. *Proc. ACL*, 2010]

Icons by interactivemania (<http://www.interactivemania.com/>) and UN Office for the Coordination of Humanitarian Affairs

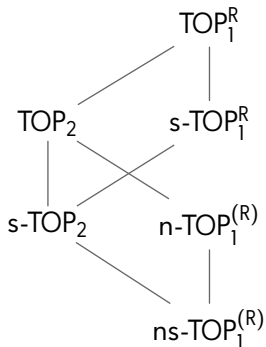
Synchronous Grammars

Illustration of rotations








Top-down Tree Transducer

Hasse diagram with composition closure indicated in subscript:

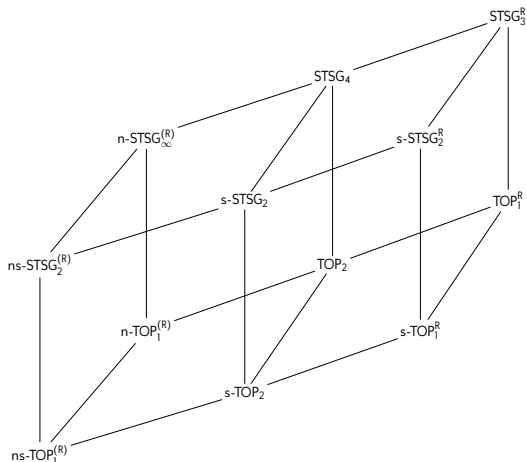


Top-down Tree Transducer

Model \ Criterion					
ns-TOP	X	X	✓	✓	✓
n-TOP	X	X	✓	✓	✓
s-TOP	X	X	✓	✓	X ₂
s-TOP ^R	X	X	✓	✓	✓
TOP	X	X	✓	✓	X ₂
TOP ^R	X	X	✓	✓	✓

Synchronous Tree Substitution Grammars






Hasse diagram with the composition closure indicated in subscript:



composition closures by

[Engelfriet, Fülöp, M.: Composition closure of linear extended top-down tree transducers. *Theory of Computing Systems*, to appear 2016]

Synchronous Tree Substitution Grammars

Model \ Criterion					
n-TOP	X	X	✓	✓	✓
TOP	X	X	✓	✓	X ₂
TOP ^R	X	X	✓	✓	✓
ns-STSG	✓	✓	✓	✓	X ₂
n-STSG	✓	X	✓	✓	X _∞
s-STSG ^(R)	✓	X	✓	✓	X ₂
STSG	✓	X	✓	✓	X ₄
STSG ^R	✓	X	✓	✓	X ₃

Advantages of SMTSG

- always have regular look-ahead
- can always be made nondeleting & shallow
- closed under composition

[Engelfriet, Liliu, M.: Extended multi bottom-up tree transducers — composition and decomposition. *Acta Informatica* 46(8), 2009]

Advantages of SMTSG

- always have regular look-ahead
- can always be made nondeleting & shallow
- closed under composition

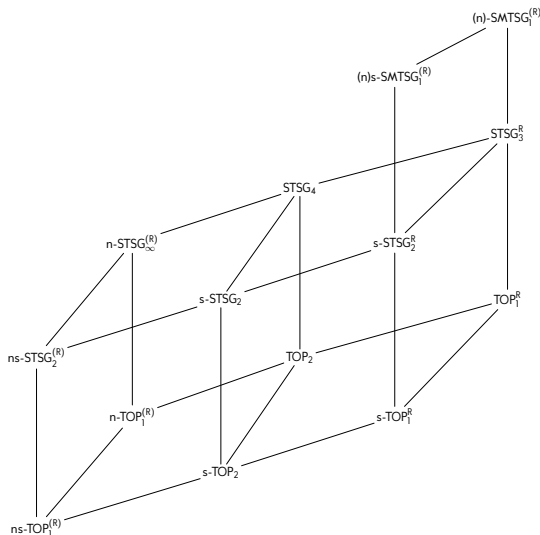
Disadvantages of SMTSG

- non-regular range






[Engelfriet, Liliu, M.: Extended multi bottom-up tree transducers — composition and decomposition. *Acta Informatica* 46(8), 2009]

Synchronous Multi Tree Substitution Grammars

Hasse diagram with the composition closure indicated in subscript:



Synchronous Multi Tree Substitution Grammars

Model \ Criterion					
n-TOP	X	X	✓	✓	✓
TOP	X	X	✓	✓	X ₂
TOP ^R	X	X	✓	✓	✓
ns-STSG	✓	✓	✓	✓	X ₂
n-STSG	✓	X	✓	✓	X _∞
s-STSG ^(R)	✓	X	✓	✓	X ₂
STSG	✓	X	✓	✓	X ₄
STSG ^R	✓	X	✓	✓	X ₃
(n)s-SMTSG ^(R)	✓	X	✓	X	✓
(n)-SMTSG ^(R)	✓	X	✓	X	✓
reg.-range SMTSG	✓	X	✓	✓	✓
symmetric SMTSG	✓	✓	✓	✓	✓

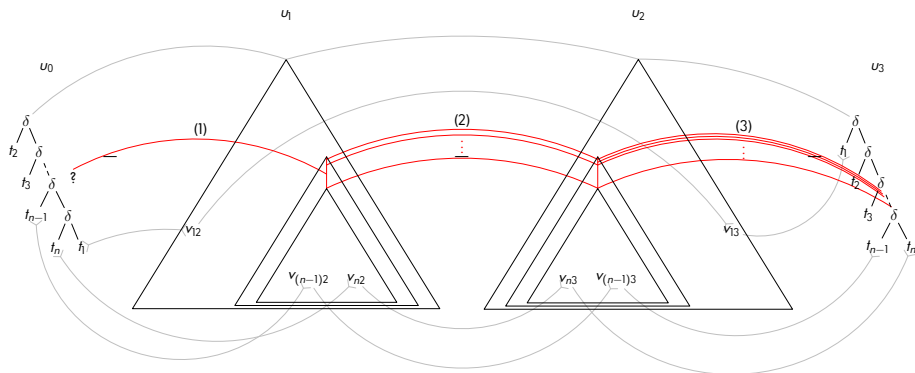
(string-level) range characterization by

[Gildea: On the string translations produced by multi bottom-up tree transducers. *Computational Linguistics* 38(3), 2012]

Synchronous Multi Tree Substitution Grammars

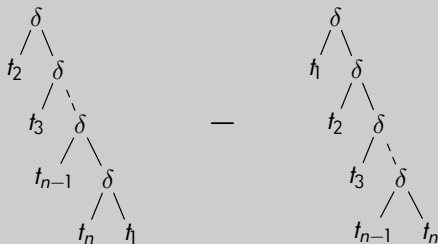
Theorem

$$(\text{STSG}^R)^3 \subsetneq \text{reg.-range SMTSG}$$



[M: The power of weighted regularity-preserving multi bottom-up tree transducers. *Int. J. Found. Comput. Sci.* 26(7), 2015]

Counterexample relation



- abstracts a well-known linguistic transformation called **topicalization**
- implementable by SMTSG, but not by any composition of STSG

Illustration of topicalization

- It rained **yesterday night**.

Topicalized: **Yesterday night**, it rained.

Illustration of topicalization

- It rained **yesterday night**.

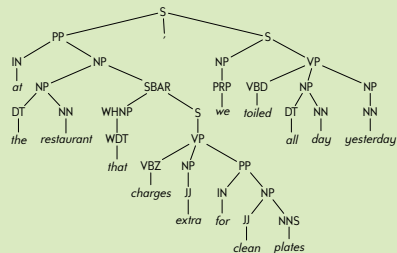
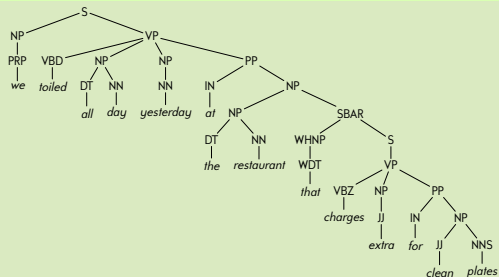
Topicalized: **Yesterday night**, it rained.

- We toiled all day yesterday
at the restaurant that charges extra for clean plates.

Topicalized: **At the restaurant that charges extra for clean plates**,
we toiled all day yesterday.

Synchronous Multi Tree Substitution Grammars

On the tree level



Contributions

- SMTSG implementation and evaluation

[Braune, Seemann, Quernheim, M.: Shallow local multi bottom-up tree transducers in SMT. *Proc. ACL*, 2013]

[Seemann, Braune, M.: String-to-tree multi bottom-up tree transducers. *Proc. ACL*, 2015]

[Seemann, Braune, M.: A systematic evaluation of MBOT in statistical machine translation. *Proc. MT-Summit*, 2015]

- characterization of expressive power of STSG and SMTSG

[Engelfriet, Lilin, M.: Extended multi bottom-up tree transducers — composition and decomposition. *Acta Informatica* 46(8), 2009]

[Engelfriet, Fülöp, M.: Composition closure of linear extended top-down tree transducers. *Theory of Computing Systems*, 2015]

[M.: The power of weighted regularity-preserving multi bottom-up tree transducers. *Int. J. Found. Comput. Sci.*, 2015]

- new proof technique (based on synchronization links)

[Fülöp, M.: Linking theorems for tree transducers. Manuscript, 2014]

similar ideas used in

[Bojanczyk: Transducers with origin information. *Proc. ICALP*, 2014]

[Filiot, Maneth, Reynier, Talbot: Decision problems of tree transducers with origin. *Proc. ICALP*, 2015]

Open Questions

- better production extraction? [better algorithms]
- additional expressive power necessary? [new models]
- further improvements possible? [tweaks]

Open Questions

- better production extraction? [better algorithms]
- additional expressive power necessary? [new models]
- further improvements possible? [tweaks]

Thank you for the attention.