

Feature-Rich Translation by Quasi-Synchronous Lattice Parsing

Kevin Gimpel and Noah A. Smith



Carnegie Mellon

Introduction

- Two trends in machine translation research
 - Many approaches to decoding
 - Phrase-based
 - Hierarchical phrase-based
 - Tree-to-string
 - String-to-tree
 - Tree-to-tree
 - Regardless of decoding approach, addition of richer features can improve translation quality



Carnegie Mellon

Introduction

- Two trends in machine translation research
 - Many approaches to decoding
 - Phrase-based
 - Hierarchical phrase-based
 - Tree-to-string
 - String-to-tree
 - Tree-to-tree
 - Regardless of decoding approach, addition of richer features can improve translation quality
- Decoding algorithms are strongly tied to features permitted



Carnegie Mellon

Phrase-Based Decoding

konnten sie es übersetzen ?
| / \ /
could you translate it ?



Carnegie Mellon

Phrase-Based Decoding

konnten sie es übersetzen ?
| / \ /
could you translate it ?

Phrase Table

- 1 konnten → could
- 2 konnten sie → could you
- 3 es übersetzen → translate it
- 4 sie es übersetzen → you translate it
- 5 es → it
- 6 ? → ?
- ...



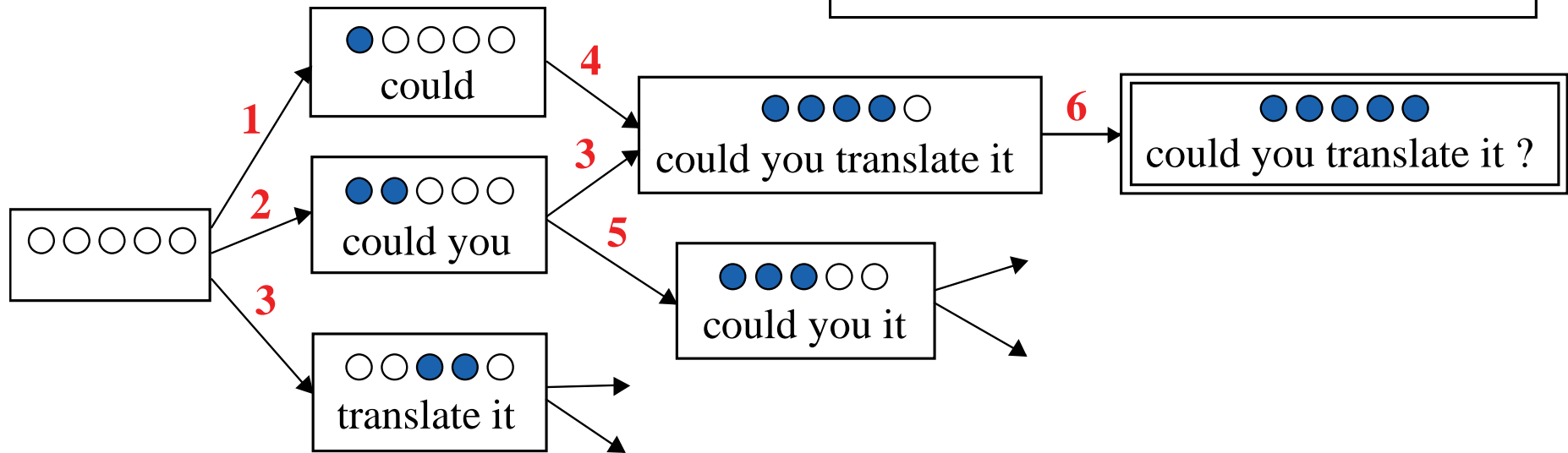
Carnegie Mellon

Phrase-Based Decoding

konnten sie es übersetzen ?
 | / \ /
 could you translate it ?

Phrase Table

- 1** konnten → could
- 2** konnten sie → could you
- 3** es übersetzen → translate it
- 4** sie es übersetzen → you translate it
- 5** es → it
- 6** ? → ?
- ...

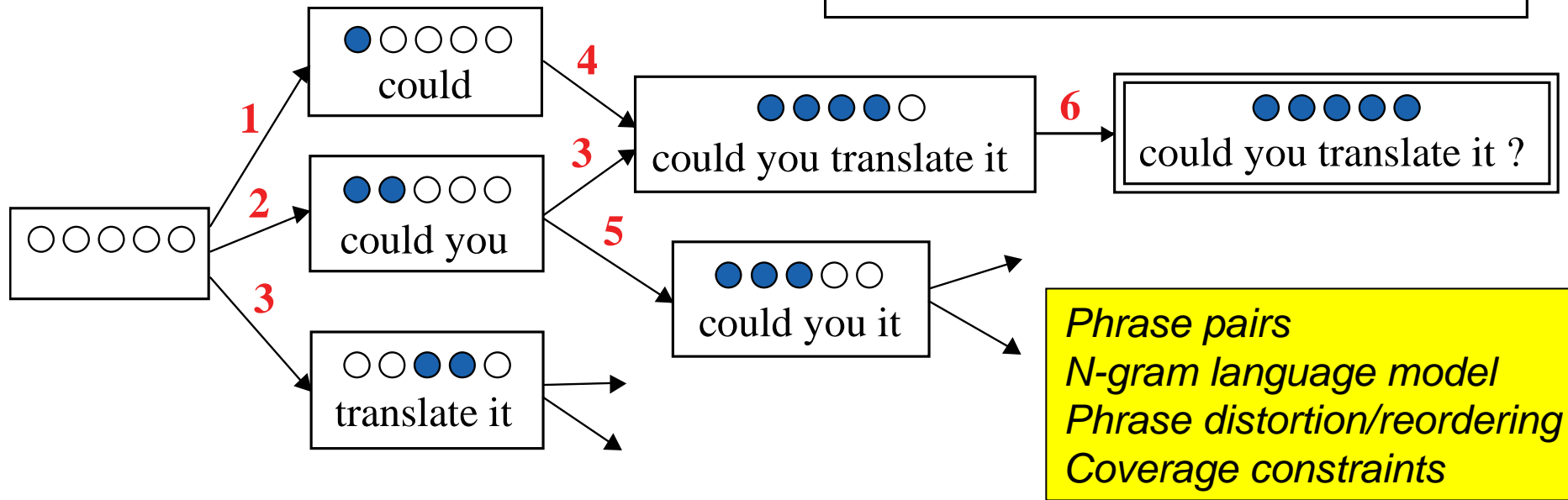


Carnegie Mellon

Phrase-Based Decoding

konnten sie es übersetzen ?
 | / \ /
 could you translate it ?

Phrase Table	
1	konnten → could
2	konnten sie → could you
3	es übersetzen → translate it
4	sie es übersetzen → you translate it
5	es → it
6	? → ?
...	



Hierarchical Phrase-Based Decoding

0 konnten₁ sie₂ es₃ übersetzen₄ ?₅

could you translate it ?

SCFG Rules

- 1 $X \rightarrow \text{es übersetzen / translate it}$
- 2 $X \rightarrow \text{es / it}$
- 3 $X \rightarrow \text{übersetzen / translate}$
- 4 $X \rightarrow \text{konnten sie } X ? / \text{could you } X ?$
- 5 $X \rightarrow \text{konnten sie } X_1 X_2 ? / \text{could you } X_2 X_1 ?$
- ...



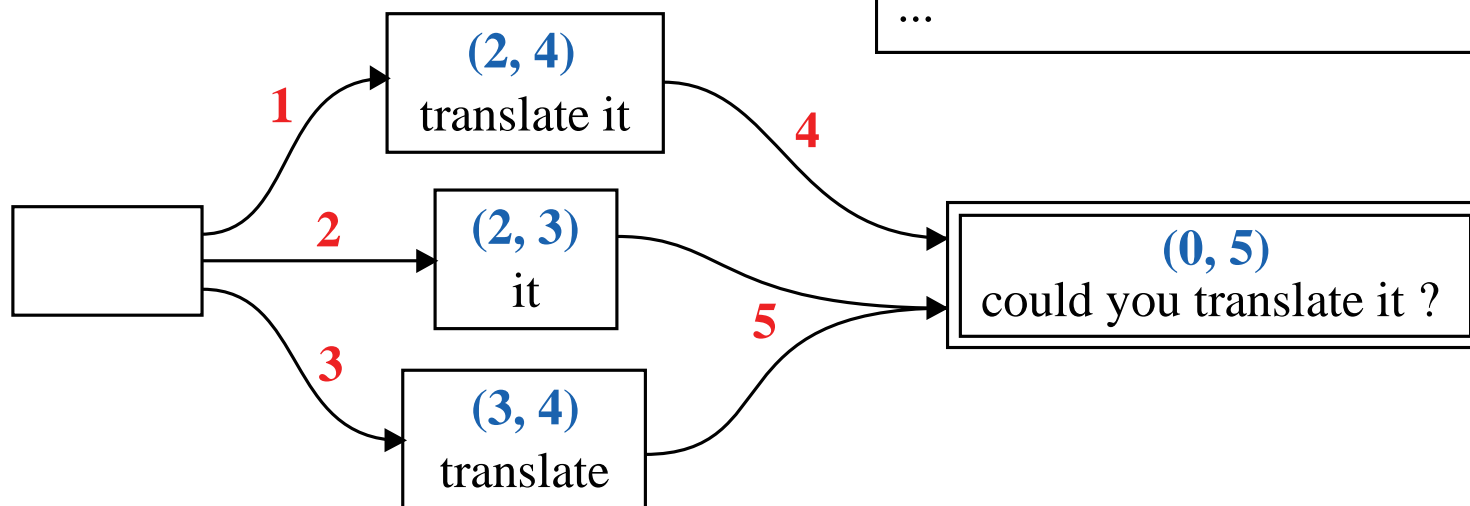
Carnegie Mellon

Hierarchical Phrase-Based Decoding

0 konnten₁ sie₂ es₃ übersetzen₄ ?₅

could you translate it ?

SCFG Rules	
1	X → es übersetzen / translate it
2	X → es / it
3	X → übersetzen / translate
4	X → konnten sie X ? / could you X ?
5	X → konnten sie X ₁ X ₂ ? / could you X ₂ X ₁ ?
...	...



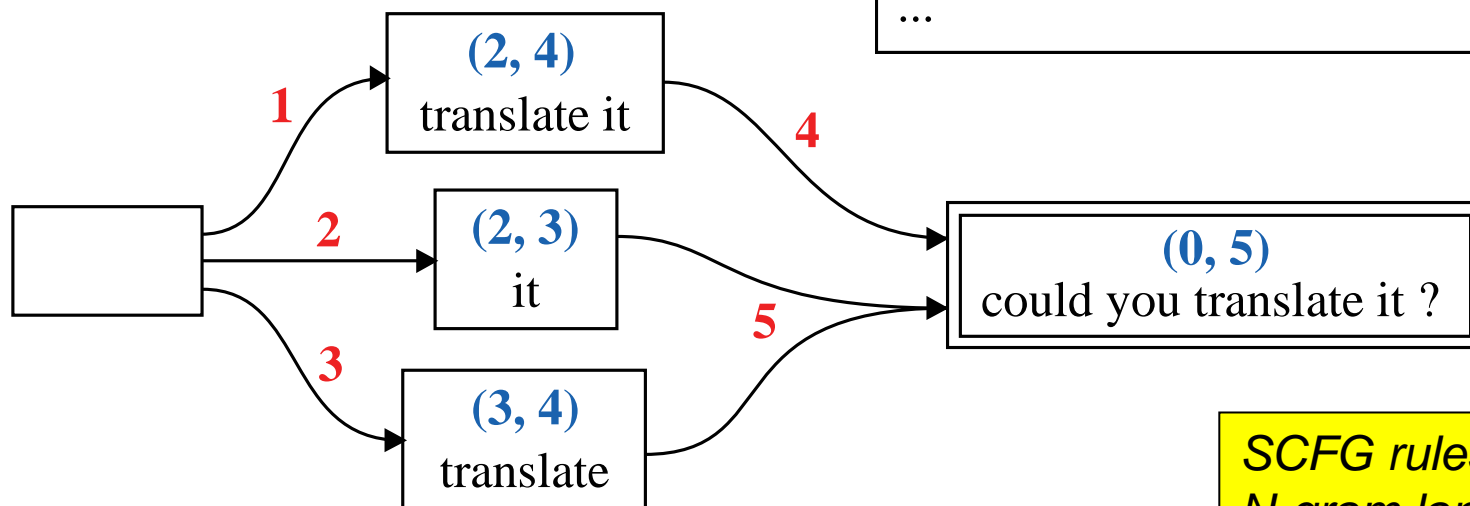
Carnegie Mellon

Hierarchical Phrase-Based Decoding

0 konnten₁ sie₂ es₃ übersetzen₄ ?₅

could you translate it ?

SCFG Rules	
1	X → es übersetzen / translate it
2	X → es / it
3	X → übersetzen / translate
4	X → konnten sie X ? / could you X ?
5	X → konnten sie X ₁ X ₂ ? / could you X ₂ X ₁ ?
...	...



SCFG rules
N-gram language model
Coverage constraints

Our goal:

An MT framework that allows as many features as possible without committing to any particular decoding approach



Overview

- Initial step towards a “universal decoder” that can permit any feature of source and target words/trees/alignments
- Experimental platform for comparison of formalisms, feature sets, and training methods
- Building blocks:
 - Quasi-synchronous grammar (Smith & Eisner 2006)
 - Generic approximate inference methods for non-local features (Chiang 2007; Gimpel & Smith 2009)



Carnegie Mellon

Outline

- Introduction
- **Model**
- Quasi-Synchronous Grammar
- Training and Decoding
- Experiments
- Conclusions and Future Work



Carnegie Mellon

$$\langle \mathbf{t}^*, \tau_{\mathbf{t}}^*, \mathbf{a}^* \rangle = \operatorname{argmax}_{\langle \mathbf{t}, \tau_{\mathbf{t}}, \mathbf{a} \rangle} p(\mathbf{t}, \tau_{\mathbf{t}}, \mathbf{a} \mid \mathbf{s}, \tau_{\mathbf{s}})$$

target words target tree alignment of target tree nodes to source tree nodes source words source tree



Carnegie Mellon

Parameterization

$$\langle \mathbf{t}^*, \tau_{\mathbf{t}}^*, \mathbf{a}^* \rangle = \operatorname{argmax}_{\langle \mathbf{t}, \tau_{\mathbf{t}}, \mathbf{a} \rangle} p(\mathbf{t}, \tau_{\mathbf{t}}, \mathbf{a} \mid \mathbf{s}, \tau_{\mathbf{s}})$$

- We use a single globally-normalized log-linear model:

$$p(\mathbf{t}, \tau_{\mathbf{t}}, \mathbf{a} \mid \mathbf{s}, \tau_{\mathbf{s}}) = \frac{\exp\{\boldsymbol{\theta}^\top \mathbf{g}(\mathbf{s}, \tau_{\mathbf{s}}, \mathbf{a}, \mathbf{t}, \tau_{\mathbf{t}})\}}{\sum_{\mathbf{a}', \mathbf{t}', \tau_{\mathbf{t}'}} \exp\{\boldsymbol{\theta}^\top \mathbf{g}(\mathbf{s}, \tau_{\mathbf{s}}, \mathbf{a}', \mathbf{t}', \tau_{\mathbf{t}'})\}}$$

- Features can look at any part of any structure



Carnegie Mellon

Features

- Log-linear models allow “arbitrary” features, but in practice inference algorithms must be developed to support feature sets
- Many types of features appear in MT:
 - lexical word and phrase mappings
 - *N*-gram and syntactic language models
 - distortion/reordering
 - hierarchical phrase mappings
 - syntactic transfer rules
- We want to use all of these!



Carnegie Mellon

Outline

- Introduction
- Model
- **Quasi-Synchronous Grammar**
- Training and Decoding
- Experiments
- Conclusions and Future Work



Carnegie Mellon

Quasi-Synchronous Grammar

(Smith & Eisner 06)

- A **quasi-synchronous grammar** (QG) is a model of

$$p(\mathbf{t}, \tau_{\mathbf{t}}, \mathbf{a} \mid \mathbf{s}, \tau_{\mathbf{s}})$$



Carnegie Mellon

Quasi-Synchronous Grammar

(Smith & Eisner 06)

- A **quasi-synchronous grammar** (QG) is a model of

$$p(\mathbf{t}, \tau_{\mathbf{t}}, \mathbf{a} \mid \mathbf{s}, \tau_{\mathbf{s}})$$

- $\tau_{\mathbf{t}}$
 - To model target trees, any monolingual formalism can be used
 - We use a **quasi-synchronous dependency grammar** (QDG)



Carnegie Mellon

Quasi-Synchronous Grammar

(Smith & Eisner 06)

- A **quasi-synchronous grammar** (QG) is a model of

$$p(\mathbf{t}, \tau_{\mathbf{t}}, \mathbf{a} \mid \mathbf{s}, \tau_{\mathbf{s}})$$

- $\tau_{\mathbf{t}}$

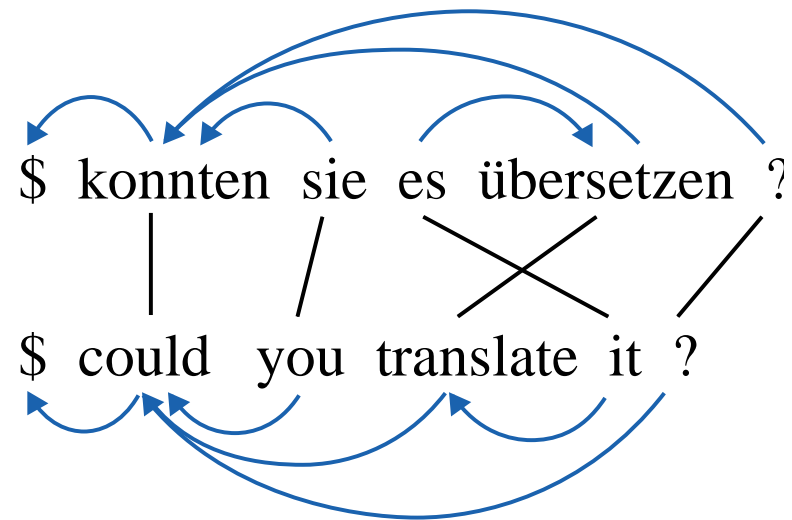
- To model target trees, any monolingual formalism can be used
- We use a **quasi-synchronous dependency grammar** (QDG)

- \mathbf{a}

- Each node in the target tree is aligned to zero or more nodes in the source tree (for a QDG, nodes = words)
- Constraints on the alignments \rightarrow synchronous grammar
- In QG, departures from synchrony are penalized *softly* using features

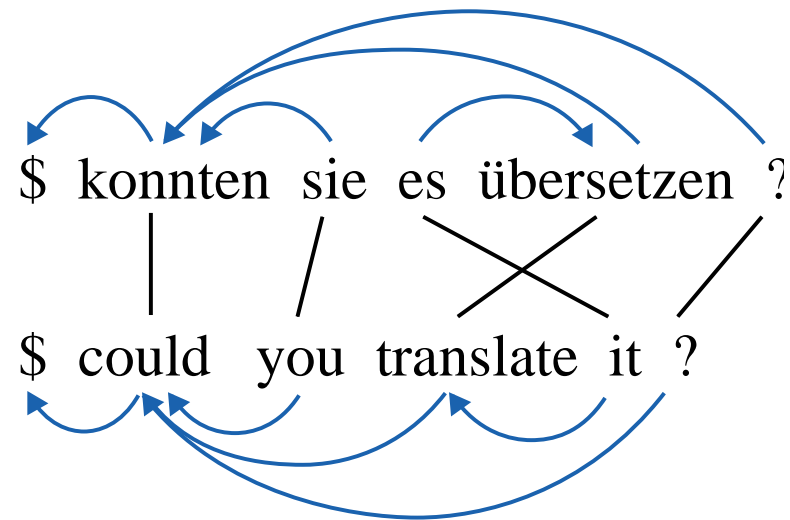


Carnegie Mellon



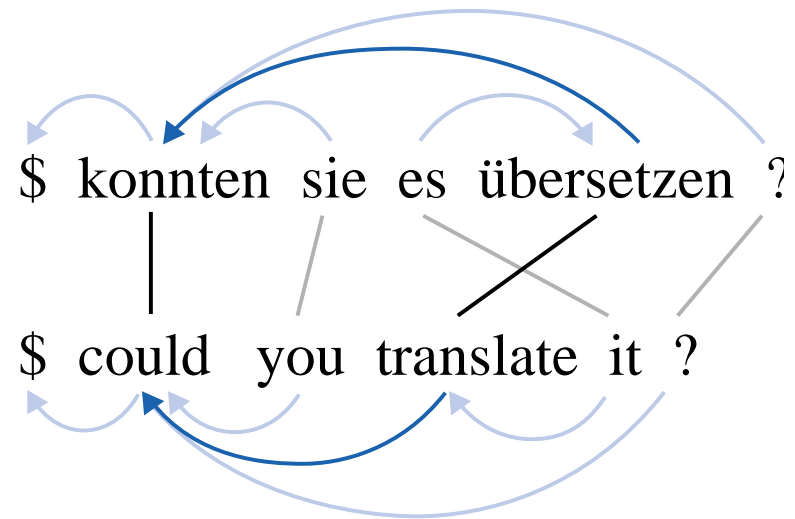
Carnegie Mellon

For every parent-child pair in the **target** sentence, what is the relationship of the **source** words they are linked to?



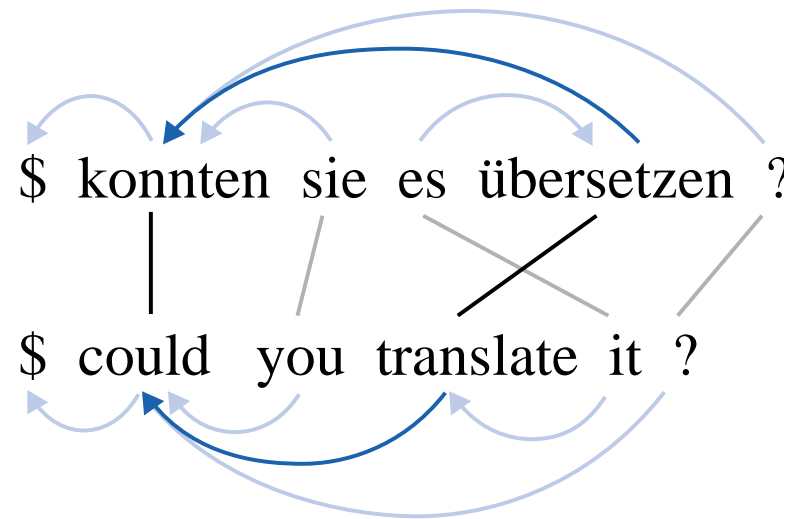
Carnegie Mellon

For every parent-child pair in the **target** sentence, what is the relationship of the **source** words they are linked to?



Carnegie Mellon

For every parent-child pair in the **target** sentence, what is the relationship of the **source** words they are linked to?



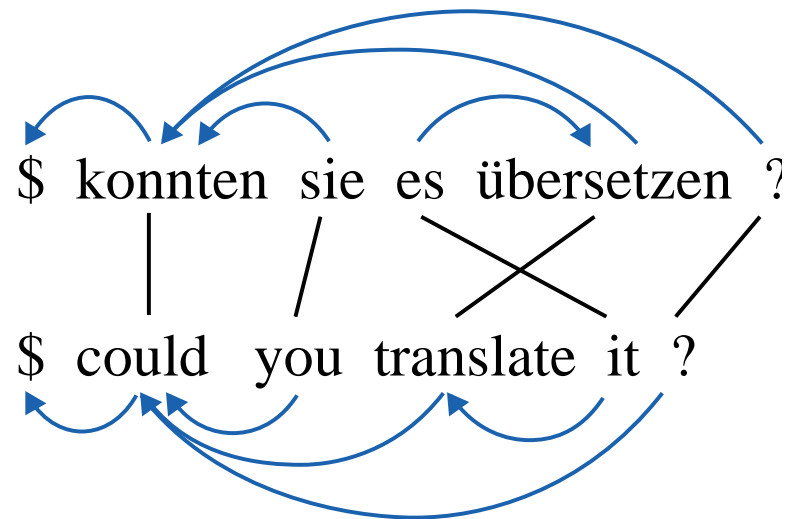
Parent-child



Carnegie Mellon

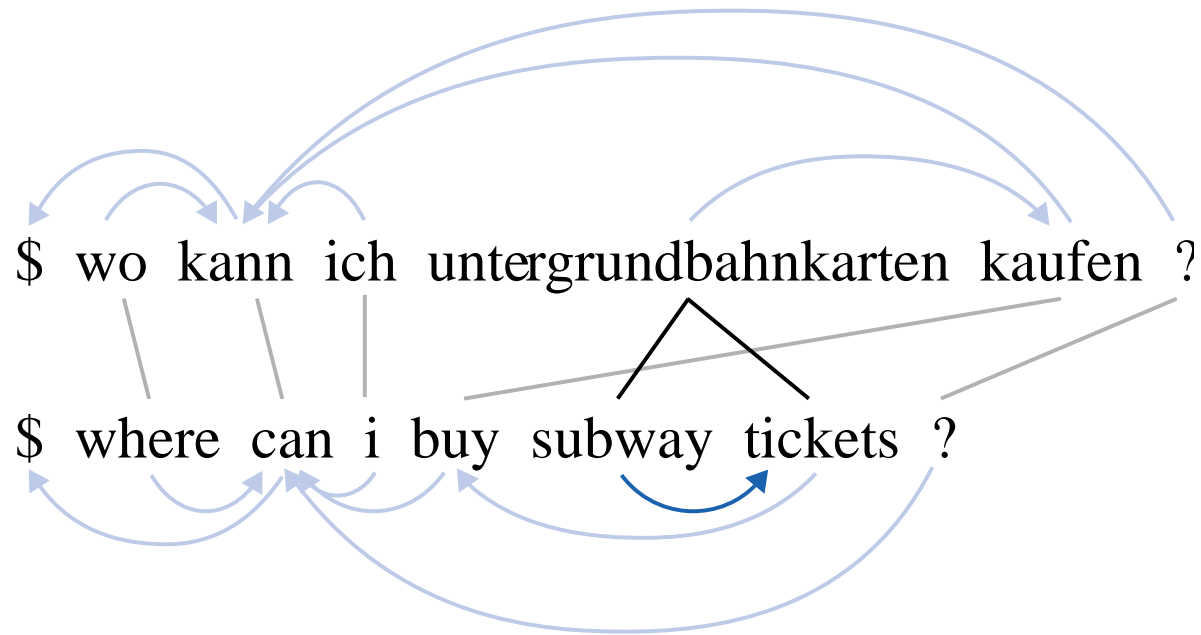
For every parent-child pair in the **target** sentence, what is the relationship of the **source** words they are linked to?

All “parent-child” configurations → synchronous dependency grammar



Carnegie Mellon

Many other configurations are possible:



Same node



Carnegie Mellon

Many other configurations are possible:

Parent-child

Child-parent

Same node

Sibling

Grandparent/child

Grandchild/parent

C-Command

Parent null

Child null

Both null

Other



Carnegie Mellon

Coverage Features

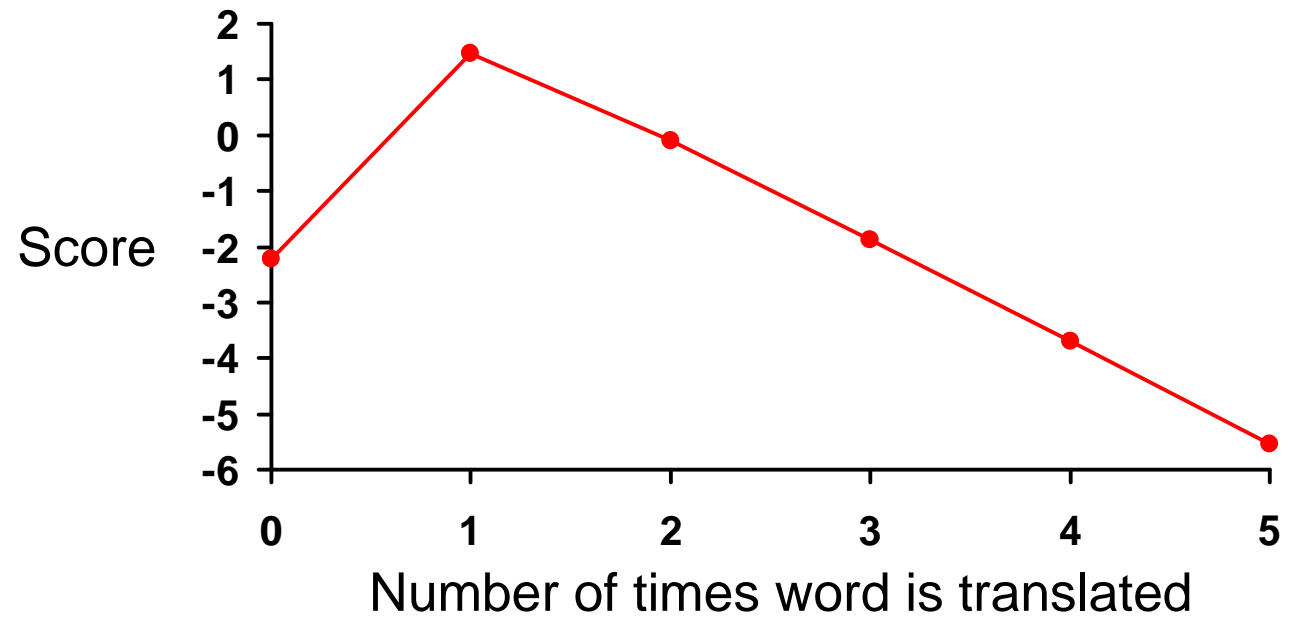
- There are no hard constraints to ensure that all source words get translated
- While QG has been used for several tasks, it has not previously been used for generation
- We add **coverage features** and learn their weights



Carnegie Mellon

Coverage Feature	Weight
Word never translated	-2.21
Word translated that was translated at least N times already:	
$N = 0$	1.48
$N = 1$	-3.04
$N = 2$	-0.22
$N = 3$	-0.05

Coverage Feature	Weight
Word never translated	-2.21
Word translated that was translated at least N times already:	
$N = 0$	1.48
$N = 1$	-3.04
$N = 2$	-0.22
$N = 3$	-0.05



Outline

- Introduction
- Model
- Quasi-Synchronous Grammar
- Training and Decoding
- Experiments
- Conclusions and Future Work



Carnegie Mellon

Decoding

- A QDG induces a monolingual grammar for a source sentence whose language consists of all possible translations
- Decoding:
 - Build a weighted lattice encoding the language of this grammar
 - Perform lattice parsing with a dependency grammar
 - Extension of dependency parsing algs for strings (Eisner 97)
 - Integrate non-local features via cube pruning/decoding (Chiang 07, Gimpel & Smith 09)



Carnegie Mellon

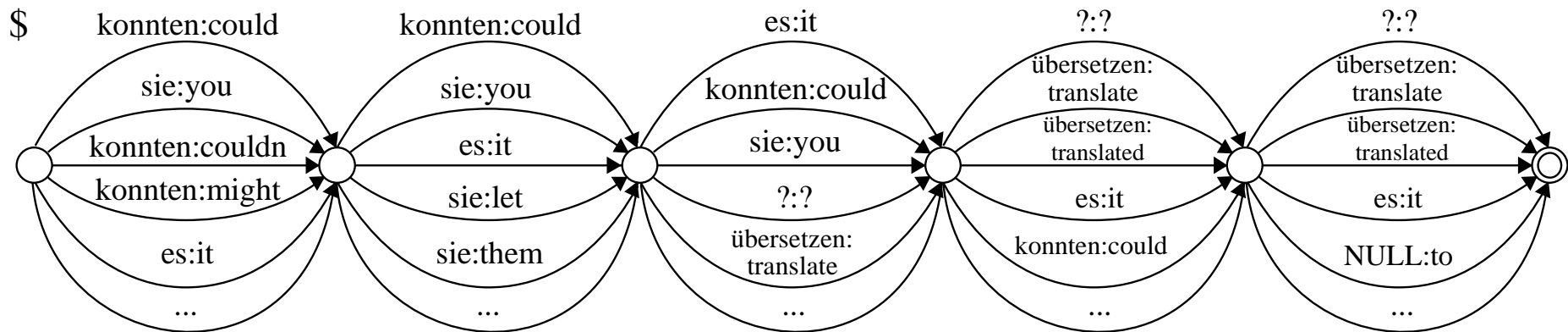
\$ konnten sie es übersetzen ?

could you translate it ?



Carnegie Mellon

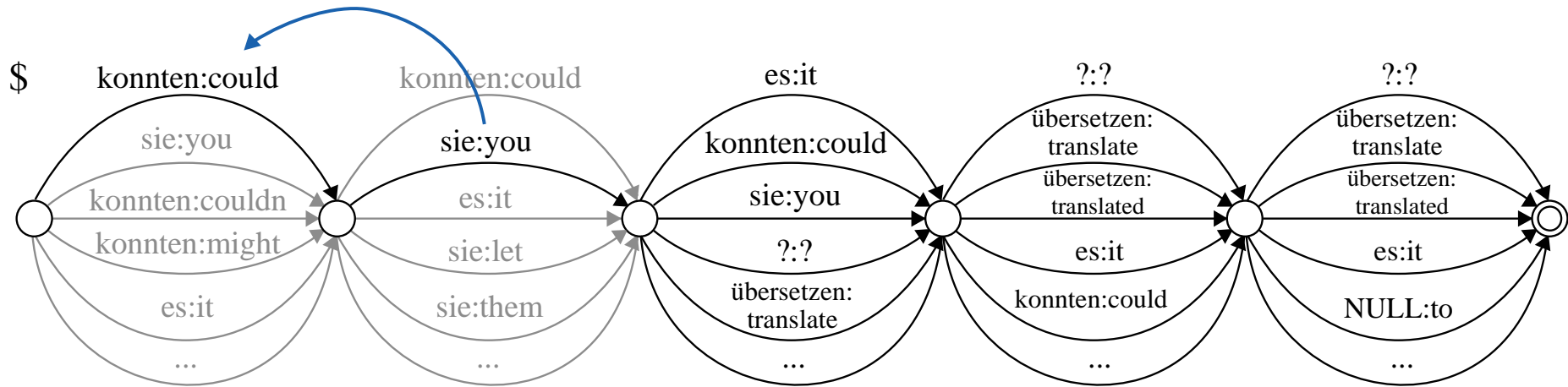
\$ konnten sie es übersetzen ?
 could you translate it ?



- ❑ Lattice arcs are weighted using lexical translation and distortion features
- ❑ Top 5 arcs shown in each bundle
- ❑ Hard limit on sentence length, multiple final states

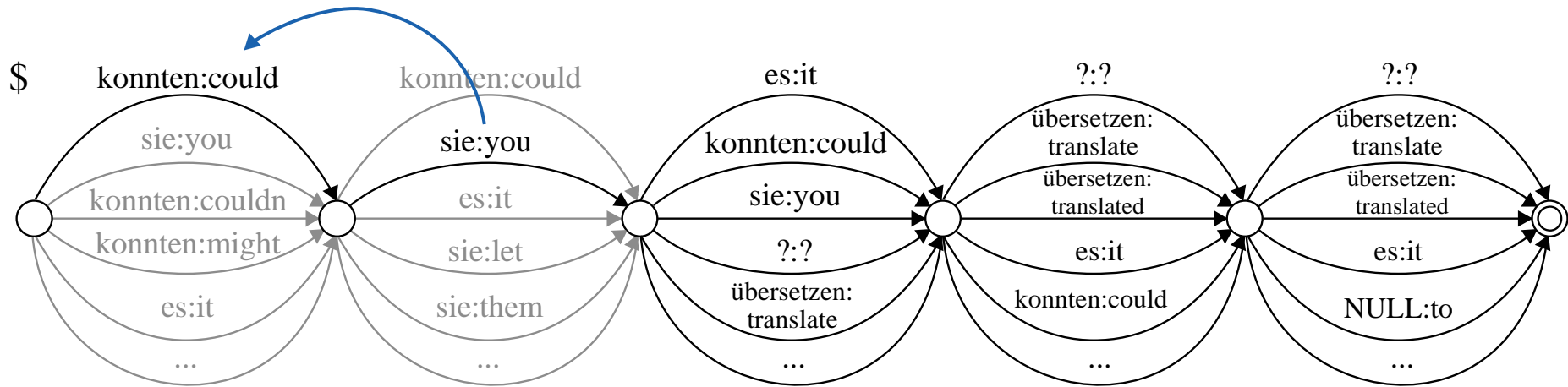
\$ konnten sie es übersetzen ?

could you translate it ?



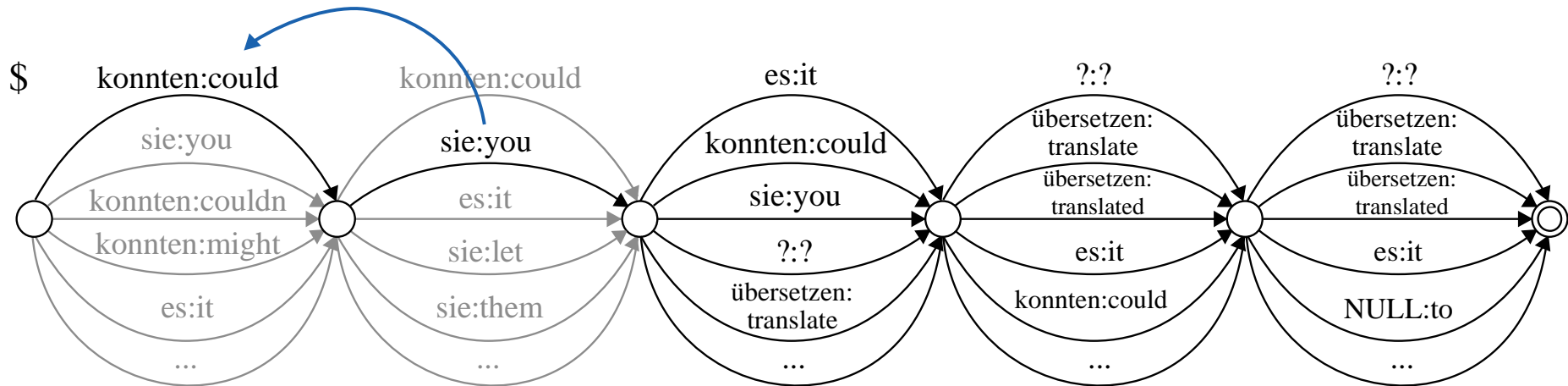
Carnegie Mellon

\$ konnten sie es übersetzen ?
 could you translate it ?



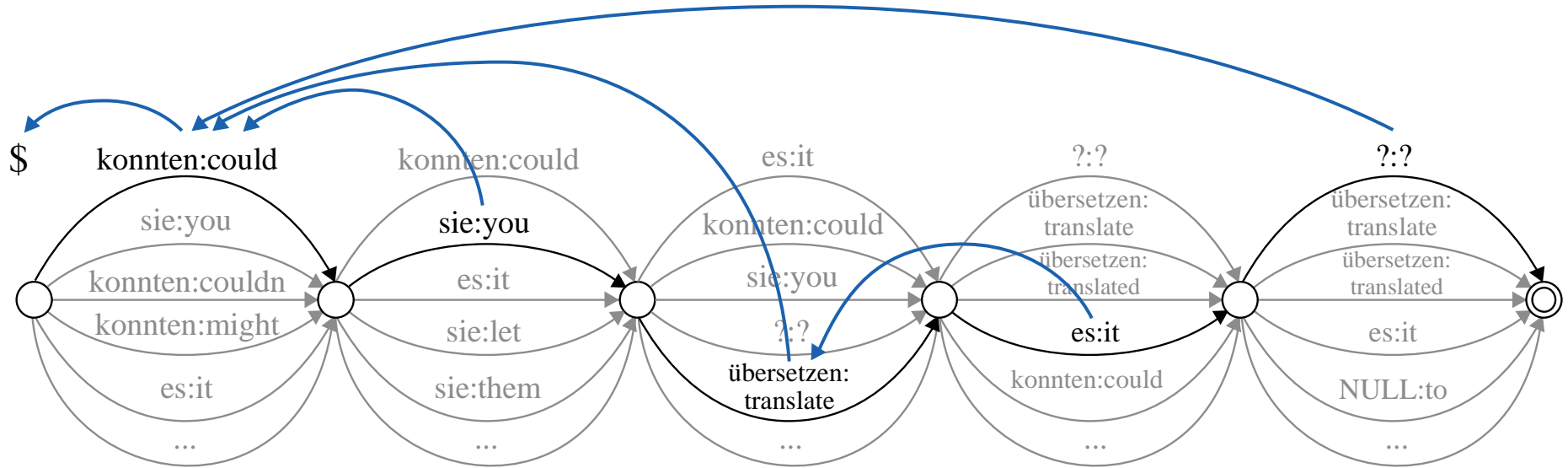
Bigram feature: "could you"

\$ konnten sie es übersetzen ?
 could you translate it ?



Bigram feature: “could you”
 Phrase features: “konnten sie” → “could you”

\$ konnten sie es übersetzen ?
 could you translate it ?



Carnegie Mellon

Training

- Recall that we use a single globally-normalized log-linear model:

$$p(\mathbf{t}, \tau_{\mathbf{t}}, \mathbf{a} \mid \mathbf{s}, \tau_{\mathbf{s}}) = \frac{\exp\{\boldsymbol{\theta}^\top \mathbf{g}(\mathbf{s}, \tau_{\mathbf{s}}, \mathbf{a}, \mathbf{t}, \tau_{\mathbf{t}})\}}{\sum_{\mathbf{a}', \mathbf{t}', \tau_{\mathbf{t}'}} \exp\{\boldsymbol{\theta}^\top \mathbf{g}(\mathbf{s}, \tau_{\mathbf{s}}, \mathbf{a}', \mathbf{t}', \tau_{\mathbf{t}'})\}}$$

- If all structures are given, this becomes a convex, supervised learning problem
- If a structure is not given, it can be marginalized out during training (or simply ignored during both training and testing)
- Here, we assume alignments are not given and marginalize them out during training



Carnegie Mellon

Training

- Standard approach is to optimize conditional likelihood

$$\begin{aligned} \text{LL}(\boldsymbol{\theta}) &= \sum_{i=1}^N \log p(\mathbf{t}^{(i)}, \tau_{\mathbf{t}}^{(i)} \mid \mathbf{s}^{(i)}, \tau_{\mathbf{s}}^{(i)}) \\ &= \sum_{i=1}^N \log \frac{\sum_{\mathbf{a}} \exp\{\boldsymbol{\theta}^\top \mathbf{g}(\mathbf{s}^{(i)}, \tau_{\mathbf{s}}^{(i)}, \mathbf{a}, \mathbf{t}^{(i)}, \tau_{\mathbf{t}}^{(i)})\}}{\sum_{\mathbf{t}, \tau_{\mathbf{t}}, \mathbf{a}} \exp\{\boldsymbol{\theta}^\top \mathbf{g}(\mathbf{s}^{(i)}, \tau_{\mathbf{s}}^{(i)}, \mathbf{a}, \mathbf{t}, \tau_{\mathbf{t}})\}} \end{aligned}$$

problem: must sum over words + trees + alignments!

Pseudo-likelihood

- Solution: optimize *pseudo-likelihood* (Besag, 1975) by making the following approximation:

$$p(\mathbf{t}, \tau_{\mathbf{t}} \mid \mathbf{s}, \tau_{\mathbf{s}}) \approx p(\mathbf{t} \mid \tau_{\mathbf{t}}, \mathbf{s}, \tau_{\mathbf{s}}) \times p(\tau_{\mathbf{t}} \mid \mathbf{t}, \mathbf{s}, \tau_{\mathbf{s}})$$

- The objective function becomes:

$$\begin{aligned} \text{PL}(\boldsymbol{\theta}) = & \sum_{i=1}^N \log \left(\sum_{\mathbf{a}} p(\mathbf{t}^{(i)}, \mathbf{a} \mid \tau_{\mathbf{t}}^{(i)}, \mathbf{s}^{(i)}, \tau_{\mathbf{s}}^{(i)}) \right) && \text{sum over} \\ & && \text{words + alignments} \\ & + \sum_{i=1}^N \log \left(\sum_{\mathbf{a}} p(\tau_{\mathbf{t}}^{(i)}, \mathbf{a} \mid \mathbf{t}^{(i)}, \mathbf{s}^{(i)}, \tau_{\mathbf{s}}^{(i)}) \right) && \text{sum over} \\ & && \text{trees + alignments} \end{aligned}$$

- Integrate non-local features via “cube summing” [Gimpel & Smith 09]



Carnegie Mellon

Outline

- Introduction
- Model
- Quasi-Synchronous Grammar
- Training and Decoding
- **Experiments**
- Conclusions and Future Work



Carnegie Mellon

Experiments

- One of our goals was an experimental platform to address questions like the following:
 - How do phrase features interact with syntactic features?
 - How do synchronous (isomorphism) constraints affect translation quality?
 - How do string-to-tree, tree-to-string, and tree-to-tree approaches compare in terms of runtime and translation quality?
 - Does a small number of feature templates work better than a large number of binary features?
 - How do MERT/MIRA compare with optimization of conditional likelihood?



Carnegie Mellon

Experiments

- One of our goals was an experimental platform to address questions like the following:
 - How do phrase features interact with syntactic features?
 - How do synchronous (isomorphism) constraints affect translation quality?
 - How do string-to-tree, tree-to-string, and tree-to-tree approaches compare in terms of runtime and translation quality?
 - Does a small number of feature templates work better than a large number of binary features?
 - How do MERT/MIRA compare with optimization of conditional likelihood?



Carnegie Mellon

Experimental Setup

■ Data

- German-English Basic Travel Expression Corpus (BTEC)
- Only sentences of length ≤ 15
- 80k sentences for training, 1k for tuning, 500 for testing

■ Features

- Parsed source and target text using Stanford parser
- Phrase extraction using Moses (max phrase length = 3)
- Trigram language model



Carnegie Mellon

Experiments

- This is not a state-of-the-art MT system
 - Moses obtains 68.4 BLEU and 85.2 METEOR on this dataset
 - Our best scores are 52 BLEU and 75 METEOR



Carnegie Mellon

Features

Lexical Translation
$p(s t)$ $p(t s)$

Language Model
BigramProbability TrigramProbability

Reordering
AbsoluteDistortion

Coverage
WordLeftUntranslated UsedWordAlreadyUsedNTimes (N in {0,1,2,3})

Phrase Translation
$p(s t)$ $p(t s)$ $lex(s t)$ $lex(t s)$

Target Dependency
words & word classes $\left\{ \begin{array}{l} p(root) \\ p(child parent, left) \\ p(child parent, right) \\ (+ 4 \text{ valence distributions}) \end{array} \right.$

QG Configuration
(14 binary features, one for each configuration)



Carnegie Mellon

Features

Lexical Translation

$$p(s | t)$$
$$p(t | s)$$

Language Model

BigramProbability
TrigramProbability

Reordering

AbsoluteDistortion

Coverage

WordLeftUntranslated
UsedWordAlreadyUsedNTimes
(N in {0,1,2,3})

Phrase Translation

$$p(s | t)$$

$$p(t | s)$$

$$lex(s | t)$$

$$lex(t | s)$$

Target Dependency

words & word classes $\left\{ \begin{array}{l} p(\text{root}) \\ p(\text{child} | \text{parent}, \text{left}) \\ p(\text{child} | \text{parent}, \text{right}) \\ (+ 4 \text{ valence distributions}) \end{array} \right.$

QG Configuration

(14 binary features, one for each configuration)



Carnegie Mellon

Feature Set Comparison: BLEU Scores

	No Syntax Features	Target Syntax Features Only	Source & Target Syntax Features
No Phrase Features	37.3	44.6	44.2
Phrase Features	46.8	49.7	51.4



Carnegie Mellon

QDG Configuration Comparison

	BLEU	METEOR
synchronous	40.1	69.5
+ nulls, root-any	41.1	69.3
+ child-parent, same-node	43.4	68.2
+ sibling	48.8	72.2
+ grandparent/child	50.2	73.7
+ c-command	51.6	74.4
+ other	51.4	74.7



Carnegie Mellon

Conclusions and Ongoing Work

- We have described an MT system based on quasi-synchronous grammar that can use features from many types of MT systems
- We reported on preliminary experiments comparing feature sets and synchronous dependency constraints
- Ongoing work in improving decoder efficiency, adding features, and conducting additional experiments



Carnegie Mellon

Thanks!



Carnegie Mellon