Word Reordering in Statistical Machine Translation with a POS-Based Distortion Model

Kay Rottmann (UKA), Stephan Vogel (CMU)

September 7, 2007

Kay Rottmann (UKA), Stephan Vogel (CMU) Word Reordering in Statistical Machine Translation with a POS

・ロン ・回 と ・ ヨ と ・ ヨ と

Outline

Motivation The Model Experiments Conclusion Translation Examples



Motivation

- Word Order Problem
- Current Approaches
- Goals
- The Model 2
 - Using POS Information
 - Learning the Rules
 - Application of the Rules
 - Reordering of Training Corpus

3 Experiments

- Setup
- Results

Conclusion 4

5 Translation Examples

・ 同 ト ・ ヨ ト ・ ヨ ト

Word Order Problem Current Approaches Goals

Problem of Word Order

• Different languages differ in word order

・ロン ・回と ・ヨン・

э

Word Order Problem Current Approaches Goals

Problem of Word Order

- Different languages differ in word order
- Differences within small context

Example: ADJ NN \rightarrow NN ADJ

An important agreement

Un acuerto importante

イロト イヨト イヨト イヨト

Word Order Problem Current Approaches Goals

Problem of Word Order

- Different languages differ in word order
- Differences within small context

Example: ADJ NN \rightarrow NN ADJ

An important agreement

Un acuerto importante

• Long range reorderings

Example: auxiliary verb and infinite verb Ich *werde* morgen nachmittag ... *ankommen* I *will arrive* tomorrow afternoon ...

イロト イヨト イヨト イヨト

Word Order Problem Current Approaches Goals

Current Approaches

• IBM constraints [BePP96], ITG [Wu96], lexicalised block oriented model [KAMCB⁺05] . . .

・ロン ・回と ・ヨン・

Word Order Problem Current Approaches Goals

Current Approaches

- IBM constraints [BePP96], ITG [Wu96], lexicalised block oriented model [KAMCB⁺05] ...
- Reordering of source sentence [ChCF06], [PoNe06], [CrMa06]

Word Order Problem Current Approaches Goals

Current Approaches

- IBM constraints [BePP96], ITG [Wu96], lexicalised block oriented model [KAMCB⁺05] ...
- Reordering of source sentence [ChCF06], [PoNe06], [CrMa06]
 - Reordering before translation process

Word Order Problem Current Approaches Goals

Current Approaches

- IBM constraints [BePP96], ITG [Wu96], lexicalised block oriented model [KAMCB⁺05] ...
- Reordering of source sentence [ChCF06], [PoNe06], [CrMa06]
 - Reordering before translation process
 - monotone decoding

Word Order Problem Current Approaches Goals

Current Approaches

- IBM constraints [BePP96], ITG [Wu96], lexicalised block oriented model [KAMCB⁺05] ...
- Reordering of source sentence [ChCF06], [PoNe06], [CrMa06]
 - Reordering before translation process
 - monotone decoding
 - more than one word order coded in lattice structure

・ロト ・回ト ・ヨト ・ヨト

Word Order Problem Current Approaches Goals

Current Approaches

- IBM constraints [BePP96], ITG [Wu96], lexicalised block oriented model [KAMCB⁺05] ...
- Reordering of source sentence [ChCF06], [PoNe06], [CrMa06]
 - Reordering before translation process
 - monotone decoding
 - more than one word order coded in lattice structure
- ullet \Rightarrow our work based on this approach

Word Order Problem Current Approaches Goals

Goals

• Restriction of search to make it fast

Kay Rottmann (UKA), Stephan Vogel (CMU) Word Reordering in Statistical Machine Translation with a POS

・ロン ・四 と ・ ヨ と ・ モ と

æ

Word Order Problem Current Approaches Goals



- Restriction of search to make it fast
- Correct reorderings in different contexts

・ロン ・回 と ・ ヨ と ・ ヨ と

3

Word Order Problem Current Approaches Goals

Goals

- Restriction of search to make it fast
- Correct reorderings in different contexts
- Better translations of long range reorderings

イロン 不同と 不同と 不同と

Using POS Information Learning the Rules Application of the Rules Reordering of Training Corpus

How the System works

Reorderings based on rules extracted prior to translation from corpus

・ロン ・回と ・ヨン・

э

Using POS Information Learning the Rules Application of the Rules Reordering of Training Corpus

How the System works

- Reorderings based on rules extracted prior to translation from corpus
- Use of POS-Tags for generalization
 - POS-Tagger are available for many languages

・ロン ・回 と ・ ヨ と ・ ヨ と

Using POS Information Learning the Rules Application of the Rules Reordering of Training Corpus

How the System works

- Reorderings based on rules extracted prior to translation from corpus
- Use of POS-Tags for generalization
 - POS-Tagger are available for many languages
- Assign probabilies to rules
 - as a guide for the decoding process

・ロン ・回と ・ヨン・

Using POS Information Learning the Rules Application of the Rules Reordering of Training Corpus

How the System works

- Reorderings based on rules extracted prior to translation from corpus
- Use of POS-Tags for generalization
 - POS-Tagger are available for many languages
- Assign probabilies to rules
 - as a guide for the decoding process
- Create a lattice with possible reorderings

イロン イヨン イヨン イヨン

Using POS Information Learning the Rules Application of the Rules Reordering of Training Corpus

How the System works

- Reorderings based on rules extracted prior to translation from corpus
- Use of POS-Tags for generalization
 - POS-Tagger are available for many languages
- Assign probabilies to rules
 - as a guide for the decoding process
- Create a lattice with possible reorderings
- Decoder finds best monotone translation path through the lattice

イロン 不同と 不同と 不同と

Using POS Information Learning the Rules Application of the Rules Reordering of Training Corpus

What is a Rule

- A rule consists of three parts:
 - Left hand side: Sequence of POS on the source side

イロン スポン イヨン イヨン

Using POS Information Learning the Rules Application of the Rules Reordering of Training Corpus

What is a Rule

- A rule consists of three parts:
 - Left hand side: Sequence of POS on the source side
 - Right hand side: Permutation on that word order

・ロン ・雪と ・ヨン・モン

Using POS Information Learning the Rules Application of the Rules Reordering of Training Corpus

What is a Rule

- A rule consists of three parts:
 - Left hand side: Sequence of POS on the source side
 - Right hand side: Permutation on that word order
 - Score for the rule: Relative frequency

イロン 不同と 不同と 不同と

Using POS Information Learning the Rules Application of the Rules Reordering of Training Corpus

What is a Rule

- A rule consists of three parts:
 - Left hand side: Sequence of POS on the source side
 - Right hand side: Permutation on that word order
 - Score for the rule: Relative frequency
- Example: ADJ NN \rightarrow 1 0 : 0.72

・ロン ・回 と ・ ヨ と ・ ヨ と

Using POS Information Learning the Rules Application of the Rules Reordering of Training Corpus

Context Dependency of Rules

• Left hand side is the POS-Sequence that needs to be reordered

イロン イヨン イヨン イヨン

Using POS Information Learning the Rules Application of the Rules Reordering of Training Corpus

Context Dependency of Rules

- Left hand side is the POS-Sequence that needs to be reordered
- Problem: different reorderings for the same POS sequence

He will come.

Er wird kommen.

He says that he will come.

Er sagt, dass *er kommen wird*.

イロト イヨト イヨト イヨト

Using POS Information Learning the Rules Application of the Rules Reordering of Training Corpus

Context Dependency of Rules

- Left hand side is the POS-Sequence that needs to be reordered
- Problem: different reorderings for the same POS sequence

He will come.

Er wird kommen.

He says that *he will come*.

Er sagt, dass er kommen wird.

- Idea: Use more complex left hand side that indicates the context \Rightarrow
 - $\bullet\,$ Usage of POS-Tags to the left and / or right of sequence
 - $\bullet\,$ Usage of words to the left and / or right of sequence
 - Usage of words as the sequence

イロン イヨン イヨン イヨン

Using POS Information Learning the Rules Application of the Rules Reordering of Training Corpus

Example Rules with Context Information

| source sequence | rule | freq. |
|--------------------------|------|-------|
| PDAT NN VVINF | 312 | 0.60 |
| VVFIN :: PDAT NN VVINF | 312 | 0.71 |
| moechte :: PDAT NN VVINF | 312 | 0.92 |

Table: Example rules for German to English translation with no context, with one tag of context to the left and one word of context to the left

イロト イヨト イヨト イヨト

Using POS Information Learning the Rules Application of the Rules Reordering of Training Corpus

Example Rules with Context Information

| source sequence | rule | freq. |
|--------------------------|------|-------|
| PDAT NN VVINF | 312 | 0.60 |
| VVFIN :: PDAT NN VVINF | 312 | 0.71 |
| moechte :: PDAT NN VVINF | 312 | 0.92 |

Table: Example rules for German to English translation with no context, with one tag of context to the left and one word of context to the left

• "Ich moechte diese Gelegenheit nutzen ,"

イロト イヨト イヨト イヨト

Using POS Information Learning the Rules Application of the Rules Reordering of Training Corpus

Example Rules with Context Information

| source sequence | rule | freq. |
|--------------------------|------|-------|
| PDAT NN VVINF | 312 | 0.60 |
| VVFIN :: PDAT NN VVINF | 312 | 0.71 |
| moechte :: PDAT NN VVINF | 312 | 0.92 |

Table: Example rules for German to English translation with no context, with one tag of context to the left and one word of context to the left

- "Ich moechte diese Gelegenheit nutzen ,"
- becomes "Ich moechte nutzen diese Gelegenheit ,"

イロン イヨン イヨン イヨン

Using POS Information Learning the Rules Application of the Rules Reordering of Training Corpus

Learning the Rules

- Use aligned corpus with a tagged source side
- whenever there is a crossing of alignments in a sentence
- store rules for different context types and count them

イロン イヨン イヨン イヨン

Using POS Information Learning the Rules Application of the Rules Reordering of Training Corpus

Learning the Rules

- Use aligned corpus with a tagged source side
- whenever there is a crossing of alignments in a sentence
- store rules for different context types and count them
- But only if the rule occurs without being part of a larger reordering that will be learned
 - This reduces the number of rules allows longer reorderings without getting problems in decoding time
 - Significant rules will still be extracted

Using POS Information Learning the Rules Application of the Rules Reordering of Training Corpus

Learning the Rules

- Use aligned corpus with a tagged source side
- whenever there is a crossing of alignments in a sentence
- store rules for different context types and count them
- But only if the rule occurs without being part of a larger reordering that will be learned
 - This reduces the number of rules allows longer reorderings without getting problems in decoding time
 - Significant rules will still be extracted
- Compute relative frequency for every rule

・ロト ・回ト ・ヨト ・ヨト

Using POS Information Learning the Rules Application of the Rules Reordering of Training Corpus

Learning the Rules

- Use aligned corpus with a tagged source side
- whenever there is a crossing of alignments in a sentence
- store rules for different context types and count them
- But only if the rule occurs without being part of a larger reordering that will be learned
 - This reduces the number of rules allows longer reorderings without getting problems in decoding time
 - Significant rules will still be extracted
- Compute relative frequency for every rule
- Throw away rules seen less than a given threshold

・ロト ・回ト ・ヨト ・ヨト

Using POS Information Learning the Rules Application of the Rules Reordering of Training Corpus

Building the Lattice (Basics)

• Start with monotone path of the sentence, weight of every $\mathsf{edge} = 1.0$

イロン イヨン イヨン イヨン

Using POS Information Learning the Rules Application of the Rules Reordering of Training Corpus

Building the Lattice (Basics)

- Start with monotone path of the sentence, weight of every $\mathsf{edge} = 1.0$
- Test for subsequences of the sentence, if a rule for that exists
 - Start with longest subsequences
 - adjust score of first edge according to monotone path
 - before testing rules that are shorter adjust score for monotone path

・ロン ・回 と ・ ヨ と ・ ヨ と

Using POS Information Learning the Rules Application of the Rules Reordering of Training Corpus

Building the Lattice (Basics)

- Start with monotone path of the sentence, weight of every $\mathsf{edge} = 1.0$
- Test for subsequences of the sentence, if a rule for that exists
 - Start with longest subsequences
 - adjust score of first edge according to monotone path
 - before testing rules that are shorter adjust score for monotone path
- BUT: This works only for one rule type!

イロン イヨン イヨン イヨン

Using POS Information Learning the Rules Application of the Rules Reordering of Training Corpus

Building the Lattice (Advanced)

• For more rule types: Combination is needed

Using POS Information Learning the Rules Application of the Rules Reordering of Training Corpus

Building the Lattice (Advanced)

- For more rule types: Combination is needed
- Use of all individual scores is bad
 - Same reorderings get different scores because of context
 - Scores will contradict each other
 - Optimization will lead to a preferred single type

イロト イヨト イヨト イヨト

Using POS Information Learning the Rules Application of the Rules Reordering of Training Corpus

Building the Lattice (Advanced)

- For more rule types: Combination is needed
- Use of all individual scores is bad
 - Same reorderings get different scores because of context
 - Scores will contradict each other
 - Optimization will lead to a preferred single type
- $\bullet \Rightarrow$ For same reorderings use max score of all rule types
- For monotone Path:
 - use minimum score over all individual scores for the monotone path

Using POS Information Learning the Rules Application of the Rules Reordering of Training Corpus

Reordering of the Training Corpus

- Phrases from reordered corpus were shown to perform better [PoNe06]
- Idea: phrases match the situation in the lattice better than before

・ロン ・回 と ・ ヨ と ・ ヨ と

Using POS Information Learning the Rules Application of the Rules Reordering of Training Corpus

Reordering of the Training Corpus

- Phrases from reordered corpus were shown to perform better [PoNe06]
- Idea: phrases match the situation in the lattice better than before
- Question: How should the training be corpus reordered?
- Usage of alignment information to monotonize alignment
 - new alignment should be nearly monotone

Using POS Information Learning the Rules Application of the Rules Reordering of Training Corpus

Reordering of the Training Corpus

- Phrases from reordered corpus were shown to perform better [PoNe06]
- Idea: phrases match the situation in the lattice better than before
- Question: How should the training be corpus reordered?
- Usage of alignment information to monotonize alignment
 - new alignment should be nearly monotone
- Usage of the rules to reorder corpus
 - better fits the decoding situation



Setup

• English \rightarrow Spanish (TC-Star 07)

- Training Corpus: Europarl Corpus 33M Words
- Developement Set: 1.2K Sentences / 79 OOV
- Test Set: 1.1K Sentences / 105 OOV
- 2 References



Setup

• English \rightarrow Spanish (TC-Star 07)

- Training Corpus: Europarl Corpus 33M Words
- Developement Set: 1.2K Sentences / 79 OOV
- Test Set: 1.1K Sentences / 105 OOV
- 2 References
- German \leftrightarrow English (WMT 06)
 - Training Corpus: Europarl Corpus 34M Words
 - Developement Set: 2K Sentences / (306 / 62) OOV
 - Test Set: 2K Sentences / (551 / 250) OOV
 - 1 Reference



Setup

• English \rightarrow Spanish (TC-Star 07)

- Training Corpus: Europarl Corpus 33M Words
- Developement Set: 1.2K Sentences / 79 OOV
- Test Set: 1.1K Sentences / 105 OOV
- 2 References
- German \leftrightarrow English (WMT 06)
 - Training Corpus: Europarl Corpus 34M Words
 - Developement Set: 2K Sentences / (306 / 62) OOV
 - Test Set: 2K Sentences / (551 / 250) OOV
 - 1 Reference
- Brill Tagger for English (36 Tags)
- Stuttgart Tree-Tagger for German (57 Tags)

Setup Results

Combination of all Ruletypes

• Addition of different context types to the rules

| System | $en \rightarrow es$ | $en \to de$ | $de\toen$ |
|---------------|---------------------|-------------|-----------|
| Baseline(RO3) | 48.51 | 17.69 | 23.70 |
| no Context | 49.52 | 17.78 | 24.79 |
| Combination | 49.58 | 18.27 | 24.85 |

イロン イヨン イヨン イヨン

3



Combination of all Ruletypes

• Addition of different context types to the rules

| System | $en\toes$ | $en \to de$ | $de \to en$ |
|---------------|-----------|-------------|-------------|
| Baseline(RO3) | 48.51 | 17.69 | 23.70 |
| no Context | 49.52 | 17.78 | 24.79 |
| Combination | 49.58 | 18.27 | 24.85 |

• Why is further improvement sometimes so low?



Combination of all Ruletypes

• Addition of different context types to the rules

| System | $en\toes$ | $en \to de$ | $de \to en$ |
|---------------|-----------|-------------|-------------|
| Baseline(RO3) | 48.51 | 17.69 | 23.70 |
| no Context | 49.52 | 17.78 | 24.79 |
| Combination | 49.58 | 18.27 | 24.85 |

- Why is further improvement sometimes so low?
 - Spanish and English Translations already very good

Setup Results

Combination of all Ruletypes

• Addition of different context types to the rules

| System | $en \rightarrow es$ | $en \to de$ | $de \to en$ |
|---------------|---------------------|-------------|-------------|
| Baseline(RO3) | 48.51 | 17.69 | 23.70 |
| no Context | 49.52 | 17.78 | 24.79 |
| Combination | 49.58 | 18.27 | 24.85 |

- Why is further improvement sometimes so low?
 - Spanish and English Translations already very good
 - AND: Phrases did not match lexical reorderings anymore

| System | $en\toes$ | $en \to de$ | $de\toen$ |
|------------------------|-----------|-------------|-----------|
| no Lexical Reorderings | 49.83 | 18.21 | 24.88 |

Setup Results

Reordering of Source Corpus

• Reordering via GIZA++ alignment information

・ロン ・回と ・ヨン・

э

Setup Results

Reordering of Source Corpus

\bullet Reordering via GIZA++ alignment information

| System | $en\toes$ | $en \to de$ | $de \to en$ |
|--------------------|-----------|-------------|-------------|
| Combination | 49.58 | 18.27 | 24.85 |
| no Lex Reorderings | 49.83 | 18.21 | 24.88 |
| all Rules GIZA++ | 49.78 | 18.23 | 24.09 |

イロト イヨト イヨト イヨト

Setup Results

Reordering of Source Corpus

• Reordering via GIZA++ alignment information

| System | $en\toes$ | $en\tode$ | $de \to en$ |
|--------------------|-----------|-----------|-------------|
| Combination | 49.58 | 18.27 | 24.85 |
| no Lex Reorderings | 49.83 | 18.21 | 24.88 |
| all Rules GIZA++ | 49.78 | 18.23 | 24.09 |

• Reordering via GIZA++ did not help for us!

• Phrases do not match decoding situation

イロト イポト イヨト イヨト

Setup Results

Reordering of Source Corpus

• Reordering via GIZA++ alignment information

| | System | $en\toes$ | $en \to de$ | $de \to en$ |
|---|--------------------|-----------|-------------|-------------|
| • | Combination | 49.58 | 18.27 | 24.85 |
| | no Lex Reorderings | 49.83 | 18.21 | 24.88 |
| | all Rules GIZA++ | 49.78 | 18.23 | 24.09 |

- Reordering via GIZA++ did not help for us!
 - Phrases do not match decoding situation
- Reordering: Most probable word order according to Reordering Rules

| System | $en\toes$ | $en \to de$ | $de\toen$ |
|-----------------|-----------|-------------|-----------|
| Rule Reordering | 49.75 | 18.42 | 25.06 |



• Addition of context leads to improved translation quality

・ロン ・聞と ・ほと ・ほと

3



- Addition of context leads to improved translation quality
- BUT: some context types help for some languages, some hurt performance for other languages

・ロン ・回と ・ヨン・



- Addition of context leads to improved translation quality
- BUT: some context types help for some languages, some hurt performance for other languages
- Reordering source side of training corpus before phrase extraction can help



- Addition of context leads to improved translation quality
- BUT: some context types help for some languages, some hurt performance for other languages
- Reordering source side of training corpus before phrase extraction can help
- BUT: reordered corpus has to be similar to decoding situation



- Addition of context leads to improved translation quality
- BUT: some context types help for some languages, some hurt performance for other languages
- Reordering source side of training corpus before phrase extraction can help
- BUT: reordered corpus has to be similar to decoding situation
- ullet pprox 1.3 improvement on English to Spanish



- Addition of context leads to improved translation quality
- BUT: some context types help for some languages, some hurt performance for other languages
- Reordering source side of training corpus before phrase extraction can help
- BUT: reordered corpus has to be similar to decoding situation
- ullet pprox 1.3 improvement on English to Spanish
- ullet pprox 0.7 improvement on English to German

・ロト ・回ト ・ヨト ・ヨト



- Addition of context leads to improved translation quality
- BUT: some context types help for some languages, some hurt performance for other languages
- Reordering source side of training corpus before phrase extraction can help
- BUT: reordered corpus has to be similar to decoding situation
- ullet pprox 1.3 improvement on English to Spanish
- ullet pprox 0.7 improvement on English to German
- ullet pprox 1.4 improvement on German to English

・ロト ・回ト ・ヨト ・ヨト

Translation Examples

• German Source: bessere Erkenntnisse und moderne Technik bieten die Chance , die Umwelt in Europas Staedten zu verbessern .

・ロン ・回 と ・ ヨ と ・ ヨ と

Translation Examples

- German Source: bessere Erkenntnisse und moderne Technik bieten die Chance , die Umwelt in Europas Staedten zu verbessern .
 - Baseline: better knowledge and modern technology offer the chance of the environment in Europe 's cities to improve .

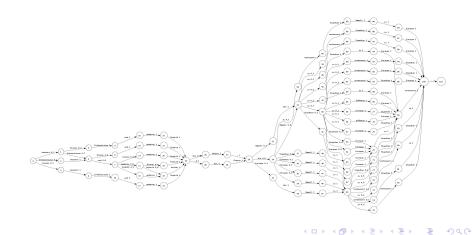
・ロン ・回と ・ヨン・

Translation Examples

- German Source: bessere Erkenntnisse und moderne Technik bieten die Chance , die Umwelt in Europas Staedten zu verbessern .
 - Baseline: better knowledge and modern technology offer the chance of the environment in Europe 's cities to improve .
 - Combination: better knowledge and modern technology offers the opportunity to improve the urban environment in Europe .



The Lattice



Kay Rottmann (UKA), Stephan Vogel (CMU) Word Reordering in

Word Reordering in Statistical Machine Translation with a POS

• Test on other language pairs (Arabic, Japanese, Farsi...)

• Test on other language pairs (Arabic, Japanese, Farsi...)

- Test on other language pairs (Arabic, Japanese, Farsi...)
- Additional internal reordering



- Test on other language pairs (Arabic, Japanese, Farsi...)
- Additional internal reordering
- Long range reorderings (more general)

・ロン ・回 と ・ ヨ と ・ ヨ と



- Test on other language pairs (Arabic, Japanese, Farsi...)
- Additional internal reordering
- Long range reorderings (more general)
- Dealing with languages without reliable POS-Tagger (using word clustering techniques)

Thank you for your attention

Kay Rottmann (UKA), Stephan Vogel (CMU) Word Reordering in Statistical Machine Translation with a POS

・ロ・・ (日・・ (日・・ (日・)

æ



- A. L. Berger, S. A. Della Pietra und V. J. Della Pietra.
 A maximum entropy approach to natural language processing. *Computational Linguistics*, 22(1), 1996, S. 39.
- B. Chen, M. Cettolo und M. Federico. Reordering rules for phrase-based statistical machine translation.

In Int. Workshop on Spoken Language Translation Evaluation Campaign on Spoken Language Translation, 2006, S. 1–15.

- Josep M. Crego und Jose B. Marino.
 Reordering Experiments for N-Gram-Based SMT.
 In Spoken Language Technology Workshop, Palm Beach, Aruba, 2006. S. 242–245.
- - P. Koehn, A. Axelrod, A. B. Mayne, C. Callison-Burch, M. Osborne und D. Talbot.



Edinburgh system description for the 2005 IWSLT speech translation evaluation.

In Proceedings of the International Workshop on Spoken Language Translation (IWSLT), Pittsburgh, PA, 2005.

M. Popovic und H. Ney.

POS-based word reorderings for statistical machine translation.

In *Proc. of the 5th Int. Conf. on Language Resources and Evaluation (LREC)*, Genoa, Italy, 2006. S. 1278.

D. Wu.

A polynomial-time algorithm for statistical machine translation.

Proc. 34th Annual Meeting of the Assoc. for Computational Linguistics, 1996, S. 152.

・ロン ・四マ ・ヨマ ・ヨマ