# Chunk-based Statistical Translation

Taro Watanabe†, Eiichiro Sumita and Hiroshi G. Okuno

`taro.watanabe@atr.co.jp`

ATR Spoken Language Translation Research Laboratories

# Contents

- Statistical Machine Translation

- Word Alignment Based Statistical Translation

- Chunk-based Statistical Translation

- Experiments

- Summary
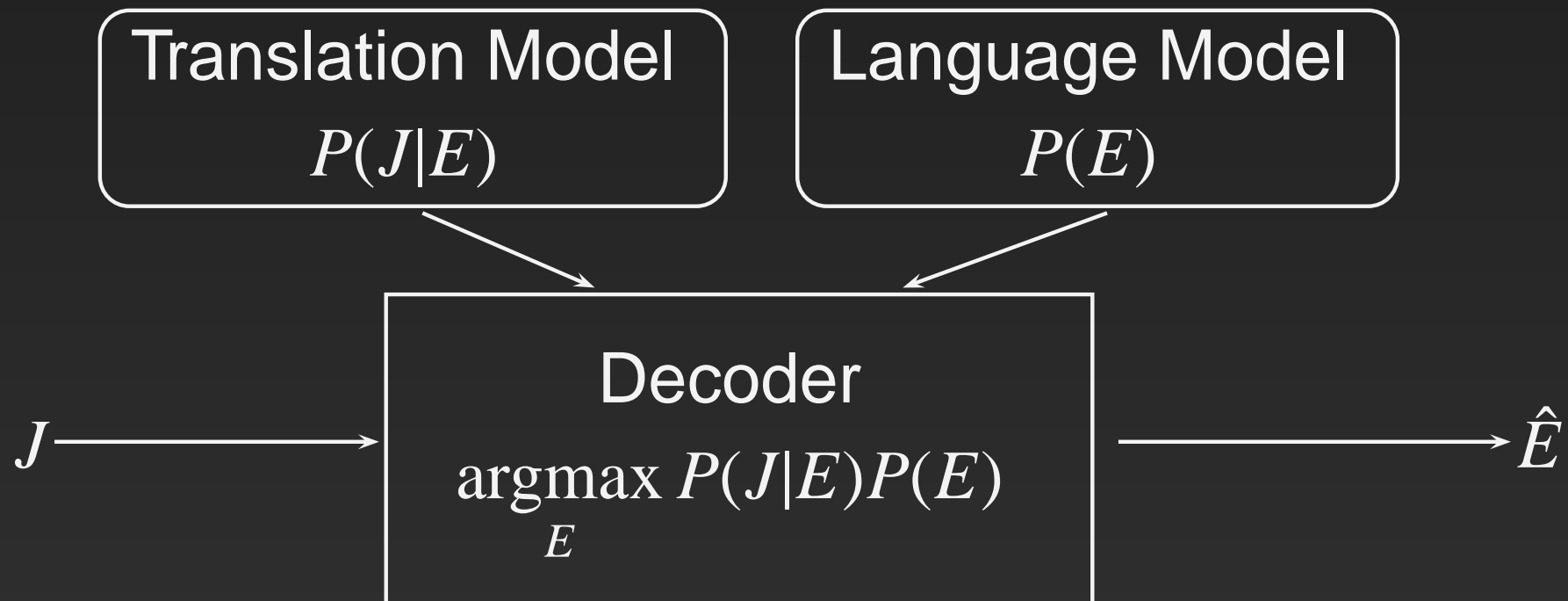
# Statistical Machine Translation

- Translation from $J$ into $E$

$$\hat{E} = \underset{E}{\text{argmax}}\, P(E|J)$$
$$= \underset{E}{\text{argmax}}\, P(E)P(J|E)$$

# Statistical Machine Translation

- Translation from $J$ into $E$

$$\hat{E} = \operatorname*{argmax}_{E} P(E|J)$$

$$= \operatorname*{argmax}_{E} P(E)P(J|E)$$

| Translation Model | Language Model |
|---|---|
| $P(J|E)$ | $P(E)$ |

Decoder

$$\operatorname*{argmax}_{E} P(J|E)P(E)$$

$J \longrightarrow$ $\longrightarrow \hat{E}$

# Word Alignment Based Statistical Translation

$$P(J|E) = \sum_A P(J, A|E)$$

$E =$ $\text{NULL}_0$ $\text{show}_1$ $\text{me}_2$ $\text{the}_3$ $\text{one}_4$ $\text{in}_5$ $\text{the}_6$ $\text{window}_7$

$J =$ $\text{uindo}_1$ $\text{no}_2$ $\text{shinamono}_3$ $\text{o}_4$ $\text{mise}_5$ $\text{tekudasai}_6$

$A = ($ 7 0 4 0 1 1 $)$

# Word Alignment Based Statistical Translation

$$P(J|E) = \sum_A P(J, A|E)$$

$E =$ NULL$_0$ show$_1$ me$_2$ the$_3$ one$_4$ in$_5$ the$_6$ window$_7$

$J =$ uindo$_1$ no$_2$ shinamono$_3$ o$_4$ mise$_5$ tekudasai$_6$

$A = ($ 7 0 4 0 1 1 $)$

- Generative Process of $P(J, A|E)$

# An Example — IBM Model 4

show$_1$ → show --→ show --→ mise

me$_2$ show NULL --→ no --→ no$_2$

the$_3$ one show → tekudasai shinamono$_3$

one$_4$ window NULL --→ o --→ o$_4$

in$_5$ one → shinamono mise$_5$

the$_6$ window → uindo tekudasai$_6$

window$_7$ uindo$_1$

**Fertility**

$n(2|E_1)$

$n(0|E_2)$

$n(0|E_3)$

...

**NULL**

$\binom{4}{2} p_0^{4-2} p_1^2$

**Lexicon**

$t(J_5|E_1)$

$t(J_6|E_1)$

$t(J_3|E_4)$

...

**Distortion**

$d_1(1 - \lceil \frac{3}{1} \rceil | E_4, J_1)$

$d_1(3 - \lceil \frac{5+6}{2} \rceil | E_1, J_3)$

$d_1(5 - \lceil \frac{2+4}{2} \rceil | \text{NULL}, J_5)$

$d_{>1}(6 - 5|J_6)$

# An Example — IBM Model 4



show$_1$ → show ⇢ show ⇢ mise
me$_2$ → show
the$_3$ → one → show → tekudasai
one$_4$ → window
in$_5$
the$_6$
window$_7$

NULL ⇢ no ⇢ no$_2$
NULL ⇢ o ⇢ o$_4$
one → shinamono
window → uindo

uindo$_1$
shinamono$_3$
mise$_5$
tekudasai$_6$

**Fertility**

$n(2|E_1)$

$n(0|E_2)$

$n(0|E_3)$

...

**NULL**

$\binom{4}{2} p_0^{4-2} p_1^2$

**Lexicon**

$t(J_5|E_1)$

$t(J_6|E_1)$

$t(J_3|E_4)$

...

**Distortion**

$d_1(1 - \lceil \frac{3}{1} \rceil | E_4, J_1)$

$d_1(3 - \lceil \frac{5+6}{2} \rceil | E_1, J_3)$

$d_1(5 - \lceil \frac{2+4}{2} \rceil | \text{NULL}, J_5)$

$d_{>1}(6 - 5 | J_6)$

# An Example — IBM Model 4



$show_1 \rightarrow show \rightarrow show \rightarrow mise$

$me_2 \quad show \quad NULL \rightarrow no$

$the_3 \quad one \quad show \rightarrow tekudasai$

$one_4 \quad window \quad NULL \rightarrow o$

$in_5 \qquad\qquad one \rightarrow shinamono$

$the_6 \qquad\qquad window \rightarrow uindo$

$window_7$

$uindo_1$
$no_2$
$shinamono_3$
$o_4$
$mise_5$
$tekudasai_6$

**Fertility**

$n(2|E_1)$

$n(0|E_2)$

$n(0|E_3)$

...

**NULL**

$\binom{4}{2} p_0^{4-2} p_1^2$

**Lexicon**

$t(J_5|E_1)$

$t(J_6|E_1)$

$t(J_3|E_4)$

...

**Distortion**

$d_1(1 - \lceil \frac{3}{1} \rceil | E_4, J_1)$

$d_1(3 - \lceil \frac{5+6}{2} \rceil | E_1, J_3)$

$d_1(5 - \lceil \frac{2+4}{2} \rceil | \text{NULL}, J_5)$

$d_{>1}(6 - 5 | J_6)$

# An Example — IBM Model 4



$$\text{show}_1 \rightarrow \text{show} \rightarrow \text{show} \rightarrow \text{mise} \qquad \text{uindo}_1$$

$$\text{me}_2 \qquad \text{show} \qquad \text{NULL} \rightarrow \text{no} \rightarrow \text{no}_2$$

$$\text{the}_3 \qquad \text{one} \qquad \text{show} \rightarrow \text{tekudasai} \qquad \text{shinamono}_3$$

$$\text{one}_4 \qquad \text{window} \qquad \text{NULL} \rightarrow \text{o} \rightarrow \text{o}_4$$

$$\text{in}_5 \qquad \qquad \text{one} \rightarrow \text{shinamono} \qquad \text{mise}_5$$

$$\text{the}_6 \qquad \qquad \text{window} \rightarrow \text{uindo} \qquad \text{tekudasai}_6$$

$$\text{window}_7$$

**Fertility**
$$n(2|E_1)$$
$$n(0|E_2)$$
$$n(0|E_3)$$
$$...$$

**NULL**
$$\binom{4}{2} p_0^{4-2} p_1^2$$

**Lexicon**
$$t(J_5|E_1)$$
$$t(J_6|E_1)$$
$$t(J_3|E_4)$$
$$...$$

**Distortion**
$$d_1(1 - \lceil \tfrac{3}{1} \rceil | E_4, J_1)$$
$$d_1(3 - \lceil \tfrac{5+6}{2} \rceil | E_1, J_3)$$
$$d_1(5 - \lceil \tfrac{2+4}{2} \rceil | \text{NULL}, J_5)$$
$$d_{>1}(6 - 5|J_6)$$

# An Example — IBM Model 4



show$_1$ → show → show → mise

me$_2$ → show

the$_3$ → one

one$_4$ → window

in$_5$

the$_6$

window$_7$

NULL → no

show → tekudasai

NULL → o

one → shinamono

window → uindo

uindo$_1$

no$_2$

shinamono$_3$

o$_4$

mise$_5$

tekudasai$_6$

**Fertility**

$n(2|E_1)$

$n(0|E_2)$

$n(0|E_3)$

...

**NULL**

$\binom{4}{2} p_0^{4-2} p_1^2$

**Lexicon**

$t(J_5|E_1)$

$t(J_6|E_1)$

$t(J_3|E_4)$

...

**Distortion**

$d_1(1 - \lceil \frac{3}{1} \rceil | E_4, J_1)$

$d_1(3 - \lceil \frac{5+6}{2} \rceil | E_1, J_3)$

$d_1(5 - \lceil \frac{2+4}{2} \rceil | \text{NULL}, J_5)$

$d_{>1}(6 - 5 | J_6)$

# Problems

- Strategy: Generate a set of words from each source word and reorder them.

# Problems

- Strategy: Generate a set of words from each source word and reorder them.

- Insertion/Deletion Modeling

  - Fertility Model to select deletion
  - A binomial distribution to determine insertion

# Problems

- Strategy: Generate a set of words from each source word and reorder them.

- Insertion/Deletion Modeling

  - Fertility Model to select deletion

  - A binomial distribution to determine insertion

- Local Alignment Modeling

  - Collection of Local Reordering $\longrightarrow$ Global Reodering

  - Long distance word alignment

# Chunk-based Statistical Translation

$$P(J|E) \;=\; \sum_{\mathcal{J}} \sum_{\mathcal{E}} P(J, \mathcal{J}, \mathcal{E}|E)$$

$\mathcal{J}, \mathcal{E}$: sequences of chunks ($|\mathcal{J}| = |\mathcal{E}|$)

# Chunk-based Statistical Translation

$$P(J|E) = \sum_{\mathcal{J}} \sum_{\mathcal{E}} P(J, \mathcal{J}, \mathcal{E}|E)$$

$$P(J, \mathcal{J}, \mathcal{E}|E) = \sum_{A} \sum_{\mathcal{A}} P(J, \mathcal{J}, A, \mathcal{A}, \mathcal{E}|E)$$
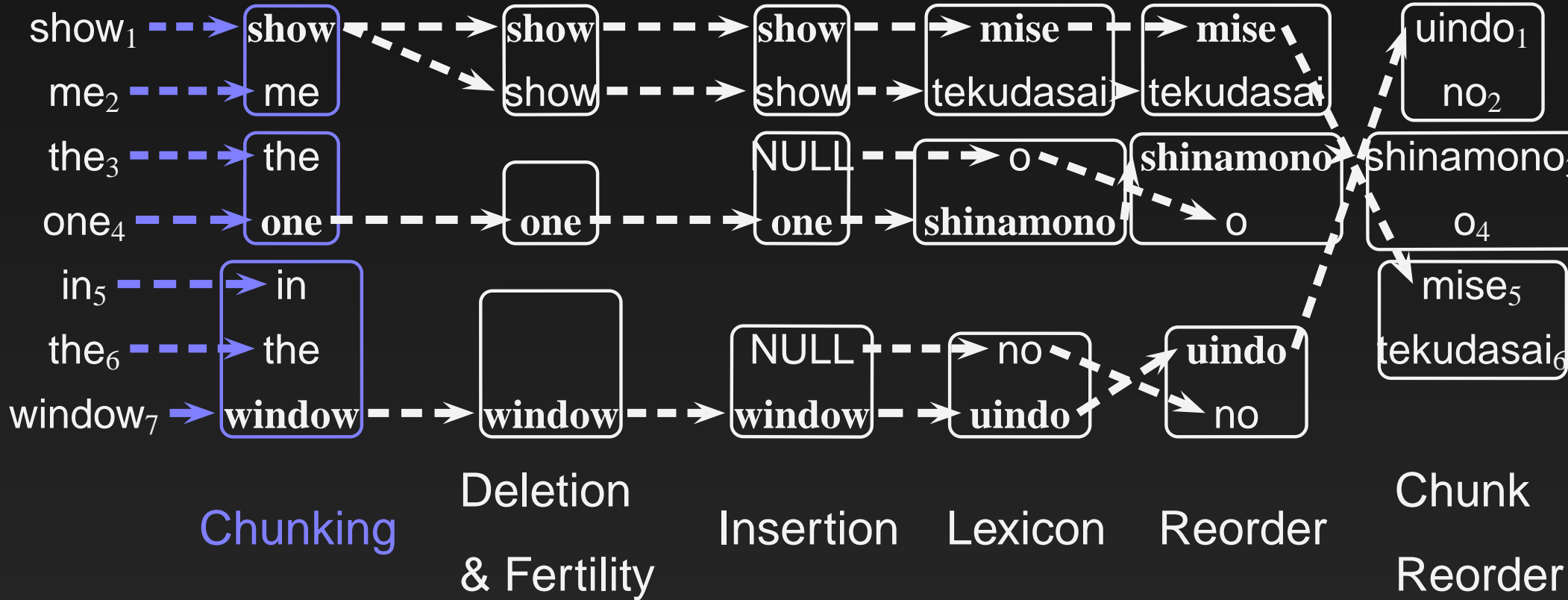
$\mathcal{J}, \mathcal{E}$: sequences of chunks $(|\mathcal{J}| = |\mathcal{E}|)$
$A$: chunk alignment
$\mathcal{A}$: word alignment

# Chunk-based Statistical Translation

$$P(J|E) = \sum_{\mathcal{J}} \sum_{\mathcal{E}} P(J, \mathcal{J}, \mathcal{E}|E)$$

$$P(J, \mathcal{J}, \mathcal{E}|E) = \sum_{A} \sum_{\mathcal{A}} P(J, \mathcal{J}, A, \mathcal{A}, \mathcal{E}|E)$$

$\mathcal{E} = $ $[\text{show}_1 \text{ me}_2]_1$ $[\text{the}_3 \text{ one}_4]_2$ $[\text{in}_5 \text{ the}_6 \text{ window}_7]_3$

$[\text{mise}_5 \text{ tekudasai}_6]$ $[\text{shinamono}_3 \text{ o}_4]$ $[\text{uindo}_1 \text{ no}_2]$

$\mathcal{J} = $ $[\text{uindo}_1 \text{ no}_2]_1$ $[\text{shinamono}_3 \text{ o}_4]_2$ $[\text{mise}_5 \text{ tekudasai}_6]_3$

$A = ($ 3 2 1 $)$

$\mathcal{A} = ($ $[7, 0]$ $[4, 0]$ $[1, 1]$ $)$

# Model Structure



| Chunking | Deletion & Fertility | Insertion | Lexicon | Reorder | Chunk Reorder |

# Model Structure



Chunking: Choose Chunk Size — $\prod_i \epsilon(\varphi_i | E_i)$

$\varphi_i$ = chunk size and if $\varphi_i > 0$ then, $E_i$ is a head word

# Model Structure



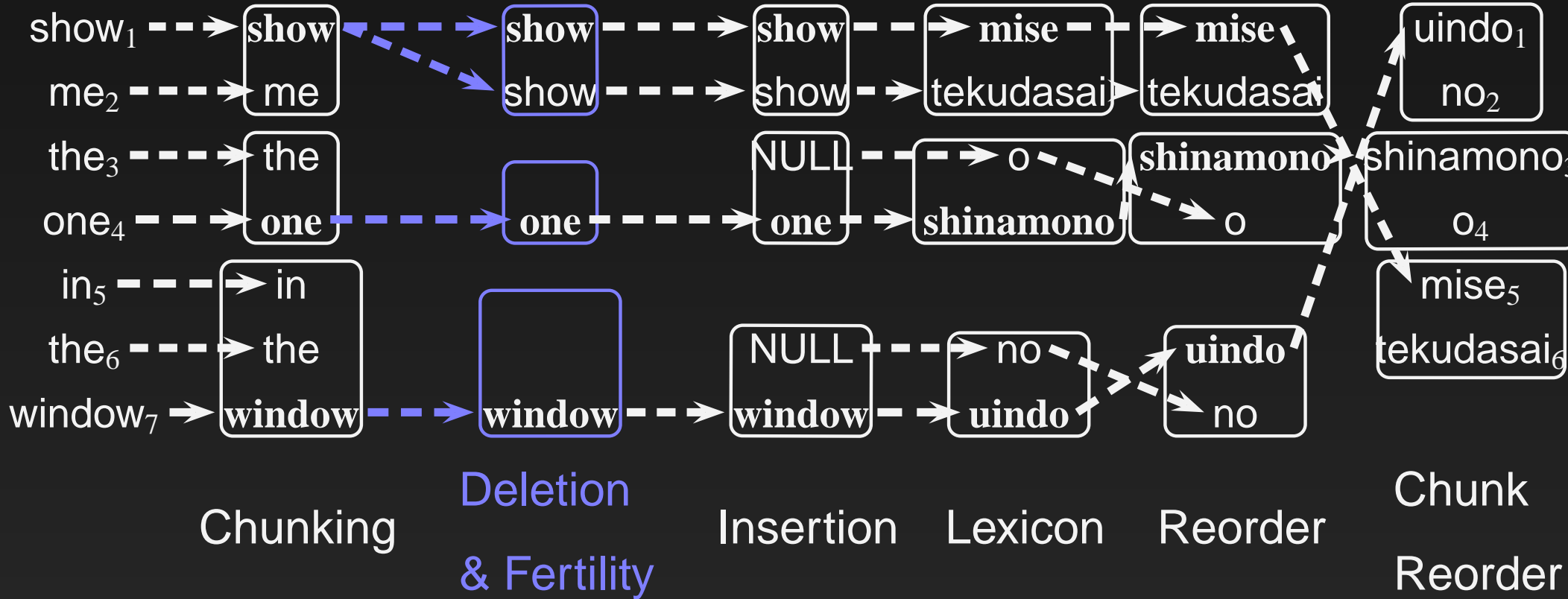Chunking: Associate Non-Head Words — $\prod_{i:\varphi_i=0} \eta(c(E_{h_i})|h_i - i, c(E_i))$
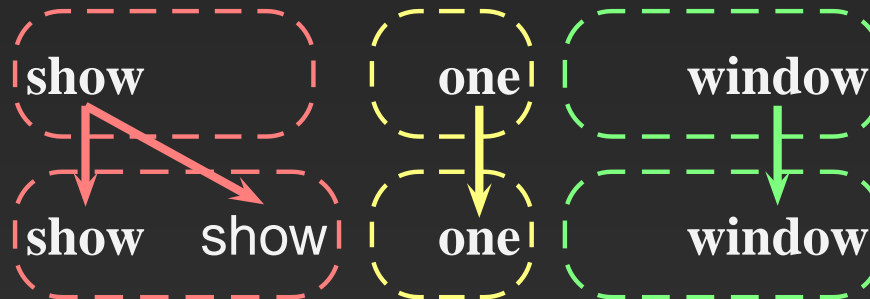
# Model Structure

| show$_1$ | → | **show** | → | **show** | → | **show** | → | **mise** | → | **mise** | uindo$_1$ |
| me$_2$ | → | me | | show | → | show | → | tekudasai | tekudasai | no$_2$ |
| the$_3$ | → | the | | | NULL | → | o | **shinamono** | shinamono$_3$ |
| one$_4$ | → | **one** | → | one | → | one | → | **shinamono** | o | o$_4$ |
| in$_5$ | → | in | | | | | | | mise$_5$ |
| the$_6$ | → | the | | NULL | → | no | **uindo** | tekudasai$_6$ |
| window$_7$ | → | **window** | → | **window** | → | **window** | → | **uindo** | no |

Chunking    Deletion & Fertility    Insertion    Lexicon    Reorder    Chunk Reorder

$$\text{Deletion} - \prod_{i:\varphi_i=0} \delta(d_i | c(E_i), c(E_{h_i}))$$

**show** me    the **one**    in the **window**

**show** me    the **one**    in the **window**

# Model Structure



$show_1$ $\dashrightarrow$ **show** $\dashrightarrow$ **show** $\dashrightarrow$ **show** $\dashrightarrow$ **mise** $\dashrightarrow$ **mise** $\;$ uindo$_1$

$me_2$ $\dashrightarrow$ me $\dashrightarrow$ show $\dashrightarrow$ show $\dashrightarrow$ tekudasai $\dashrightarrow$ tekudasai $\;$ no$_2$

$the_3$ $\dashrightarrow$ the $\qquad\qquad\qquad$ NULL $\dashrightarrow$ o $\dashrightarrow$ shinamono $\dashrightarrow$ shinamono$_3$

$one_4$ $\dashrightarrow$ **one** $\dashrightarrow$ one $\dashrightarrow$ one $\dashrightarrow$ **shinamono** $\dashrightarrow$ o $\;$ o$_4$

$in_5$ $\dashrightarrow$ in $\qquad\qquad\qquad\qquad\qquad\qquad\qquad$ mise$_5$

$the_6$ $\dashrightarrow$ the $\qquad$ NULL $\dashrightarrow$ no $\dashrightarrow$ **uindo** $\;$ tekudasai$_6$

$window_7$ $\rightarrow$ **window** $\dashrightarrow$ **window** $\dashrightarrow$ window $\dashrightarrow$ **uindo** $\dashrightarrow$ no

Chunking $\qquad$ **Deletion & Fertility** $\qquad$ Insertion $\quad$ Lexicon $\quad$ Reorder $\qquad$ Chunk Reorder

Fertility — $\prod_{i:\varphi_i>0} \nu(\phi_i|E_i)/\phi_i$

$\phi_i$ = # of words

| **show** | **one** | **window** |
| --- | --- | --- |
| **show** show | **one** | **window** |

# Model Structure



| | Chunking | Deletion & Fertility | Insertion | Lexicon | Reorder | Chunk Reorder |

$$\text{Insertion} \; — \; \prod_{i:\varphi_i>0} \iota(\phi'_i|c(E_i))$$

$\phi'_i = \#$ of NULL words

# Model Structure



Chunking
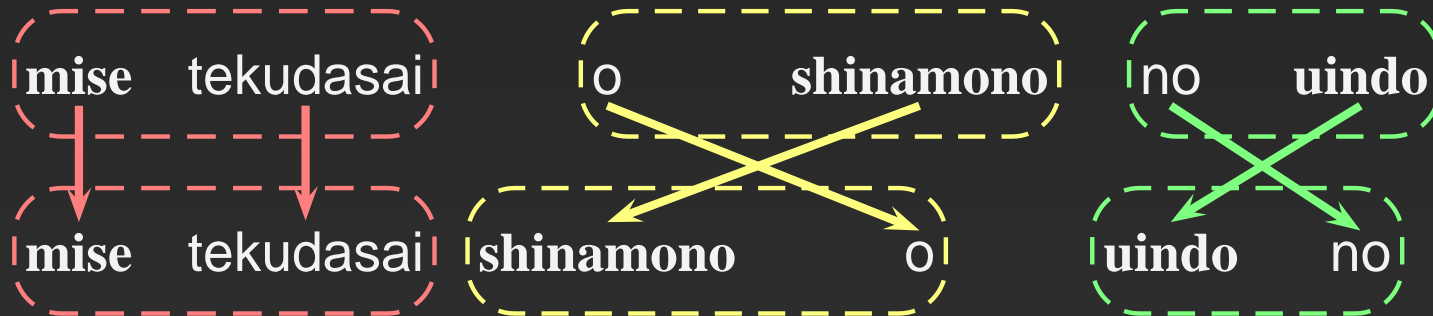
Deletion & Fertility

Insertion

Lexicon

Reorder

Chunk Reorder

$$\text{Lexical Transfer} - \prod_j \prod_k \tau(\mathcal{J}_{j,k} | E_{\mathcal{A}_{j,k}})$$
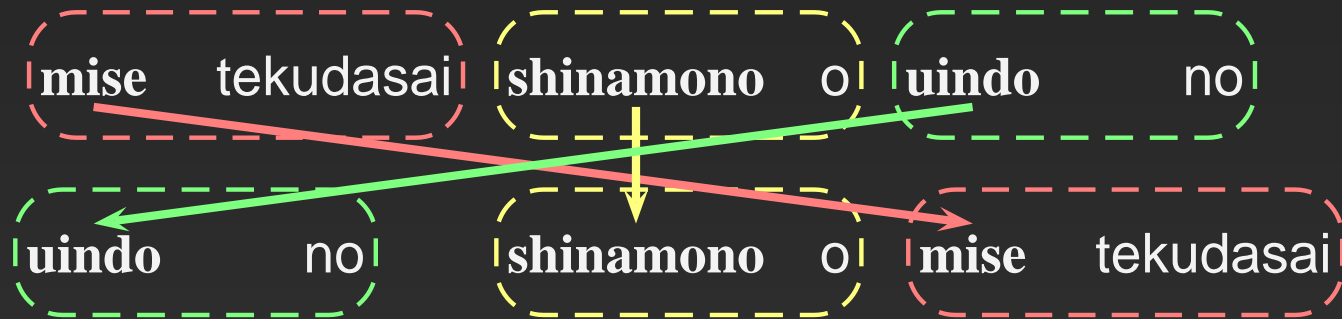
# Model Structure



Reorder — $\prod_j P(\mathcal{A}_j | \mathcal{E}_{A_j}, \mathcal{J}_j)$

# Model Structure



Chunking — Deletion & Fertility — Insertion — Lexicon — Reorder — Chunk Reorder

Chunk Reorder — $P(A|\mathcal{E}, \mathcal{J})$

# Characteristics

- String-to-String translation model with a hidden chunk layer

# Characteristics

- String-to-String translation model with a hidden chunk layer

- Bag-of-words to bag-of-words translation
  - ◆ Chunking – Translate – Reorder

# Characteristics

- String-to-String translation model
  with a hidden chunk layer

- Bag-of-words to bag-of-words translation
  - ♦ Chunking – Translate – Reorder

- Chunk-wise word insertion
  vs. Sentence-wise insertion

# Characteristics

- String-to-String translation model
  with a hidden chunk layer

- Bag-of-words to bag-of-words translation
  - Chunking – Translate – Reorder

- Chunk-wise word insertion
  vs. Sentence-wise insertion

- Chunking/Translate/Reorder
  by hypothesized "head" words

# Parameter Estimation

- EM-Algorithm

- E-step: for each pair $E$ and $J$

$$P(\mathcal{J}, A, \mathcal{A}, \mathcal{E}|J, E) = \frac{P(J, \mathcal{J}, A, \mathcal{A}, \mathcal{E}|E)}{\sum_{\mathcal{J}, A, \mathcal{A}, \mathcal{E}} P(J, \mathcal{J}, A, \mathcal{A}, \mathcal{E}|E)}$$

  Then, computes expectation

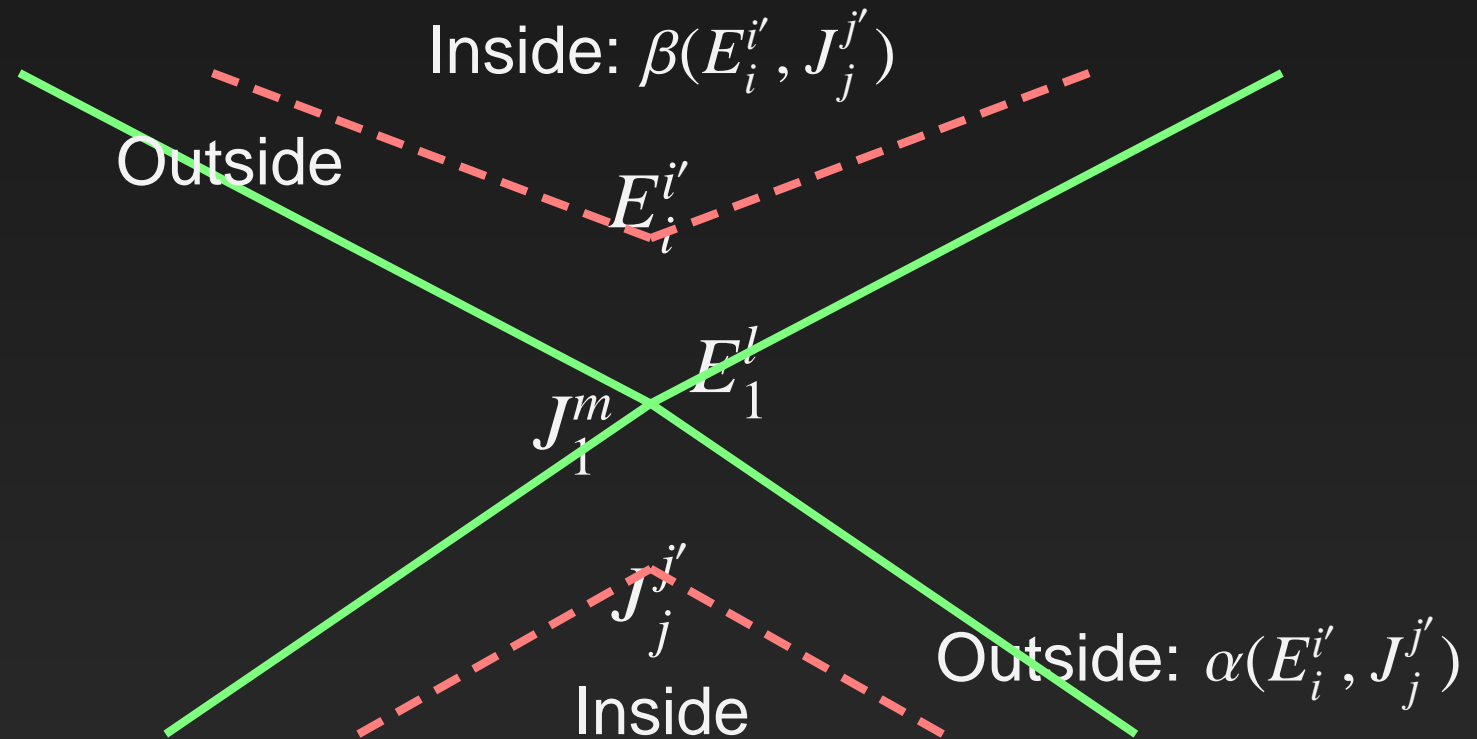- M-step: From expectation, induce parameters

# Some Tricks

- Computational problem

- Local maximum problem

# Some Tricks

- Computational problem
  - ♦ Inside-Outside Algorithm
  - ♦ Approximation
- Local maximum problem

# Some Tricks

- Computational problem
  - ♦ Inside-Outside Algorithm

Inside: $\beta(E_i^{i'}, J_j^{j'})$

Outside

$E_i^{i'}$

$E_1^l$

$J_1^m$

$J_j^{j'}$

Outside: $\alpha(E_i^{i'}, J_j^{j'})$

Inside

  - ♦ Approximation
- Local maximum problem

# Some Tricks

- Computational problem
    - ♦ Inside-Outside Algorithm
    - ♦ Approximation
        - All possible word alignment: $O(lmk^4(k+1)^k))$
        - All possible chunking/alignment: $O(2^l 2^m n!)$
- Local maximum problem

# Some Tricks

- Computational problem
  - ♦ Inside-Outside Algorithm
  - ♦ Approximation
    - All possible word alignment: $O(lmk^4(k+1)^k))$
    - All possible chunking/alignment: $O(2^l 2^m n!)$
    - Viterbi Chunking/Alignment + Neighbours
- Local maximum problem

# Some Tricks

- Computational problem
  - ♦ Inside-Outside Algorithm
  - ♦ Approximation
- Local maximum problem

# Some Tricks

- Computational problem
  - Inside-Outside Algorithm
  - Approximation

- Local maximum problem
  - Initial parameters from IBM Model 4
  - Smoothing

# Decoding

- Left-to-right generation breadth-first beam search
  - Generate possible output chunks for all possible input chunks
  - Generate hypothesized output by consuming input chunks in arbitrary order and combining possible output chunks in left-to-right order

# Decoding

- Left-to-right generation breadth-first beam search
  - Generate possible output chunks for all possible input chunks
  - Generate hypothesized output by consuming input chunks in arbitrary order and combining possible output chunks in left-to-right order

- Pruning
  - Beam size pruning
  - Example-based scoring

$$\log P_{tm}(J|E) + \log P_{lm}(E) + weight \times \sum_j freq(\mathcal{E}_{A_j}, \mathcal{J}_j)$$

  - Chunk-based translation model is a deficient model
  - Many model components

# Japanese-to-English Translation Experiments

## Basic Travel Expression Corpus

|  | Japanese | English |
|---|---|---|
| # of sentences | 171,894 | |
| # of words | 1,181,188 | 1,009,065 |
| vocabulary size | 20472 | 16232 |
| # of singletons | 82,06 | 5,854 |
| 3-gram perplexity | 23.7 | 35.8 |

- model4: IBM Model 4

- chunk3: Chunk-based Statistical Translation (chunk size $\leq$ 3)

- chunk3+: + Example-based scoring

# Sample Viterbi Chunking/Alignment

[ i * have ]    [ the * number ] [ of my * passport ]

[ *               ][ *          ] [    *         ]

[ i * have ]    [ a * stomach ache ][ please * give me ][ some * medicine ]

[     *     ]     [ *     ]       [ *     ]     [ *     ]

[ i $*$ have ][ a $*$ reservation ]   [ $*$ for ]   [ two $*$ nights ]   [ my $*$ name is ][ $*$ risa kobayashi ]

[  $*$  ]     [ $*$ ]    [   $*$   ][   $*$   ]  [  $*$  ]  [   $*$   ]

# Evaluation

**WER:** Word-error-rate, which penalizes the edit distance against reference translations.

**PER:** Position independent WER, which penalizes without considering positional disfluencies.

**BLEU:** BLEU score, which computes the ratio of n-gram for the translation results found in reference translations.

**SE:** Subjective evaluation ranks ranging from A to D (A:Perfect, B:Fair, C:Acceptable and D:Nonsense), judged by native speakers.

- Tested on 510 sentences
- 16 set of references for non-subjective evaluations

# Results

| Model | WER [%] | PER [%] | BLEU [%] | SE [%] | | |
|---|---|---|---|---|---|---|
| | | | | A | A+B | A+B+C |
| model4 | 43.3 | 37.2 | 46.5 | 59.2 | 74.1 | 80.2 |
| chunk3 | 40.9 | 36.1 | 48.4 | 59.8 | 73.5 | 78.8 |
| chunk3+ | 38.5 | 33.7 | 52.1 | 65.1 | 76.3 | 80.6 |

# Results

| Model | WER [%] | PER [%] | BLEU [%] | SE [%] | | |
|---|---|---|---|---|---|---|
| | | | | A | A+B | A+B+C |
| model4 | 43.3 | 37.2 | 46.5 | 59.2 | 74.1 | 80.2 |
| chunk3 | 40.9 | 36.1 | 48.4 | 59.8 | 73.5 | 78.8 |
| chunk3+ | 38.5 | 33.7 | 52.1 | 65.1 | 76.3 | 80.6 |

# Results

| Model | WER [%] | PER [%] | BLEU [%] | SE [%] A | A+B | A+B+C |
|---|---|---|---|---|---|---|
| model4 | 43.3 | 37.2 | 46.5 | 59.2 | 74.1 | 80.2 |
| chunk3 | 40.9 | 36.1 | 48.4 | 59.8 | 73.5 | 78.8 |
| chunk3+ | 38.5 | 33.7 | 52.1 | 65.1 | 76.3 | 80.6 |

# Results

| Model | WER [%] | PER [%] | BLEU [%] | SE [%] | | |
|---|---|---|---|---|---|---|
| | | | | A | A+B | A+B+C |
| model4 | 43.3 | 37.2 | 46.5 | 59.2 | 74.1 | 80.2 |
| chunk3 | 40.9 | 36.1 | 48.4 | 59.8 | 73.5 | 78.8 |
| chunk3+ | 38.5 | 33.7 | 52.1 | 65.1 | 76.3 | 80.6 |

# Sample Translations

| | |
|---|---|
| input: | |
| reference: | is this all the baggage from flight one five two |
| model4: | is this all you baggage for flight one five two |
| chunk3+: | is this all the baggage from flight one five two |
| input: | |
| reference: | may i have room service for breakfast please |
| model4: | please give me some room service please |
| chunk3+: | i 'd like room service for breakfast |
| input: | |
| reference: | hello i 'd like to change my reservation for march nineteenth |
| model4: | i 'd like to change my reservation for ninety days be march hello |
| chunk3+: | hello i 'd like to change my reservation on march nineteenth |
| input: | |
| reference: | wait a couple of minutes i 'm telephoning now |
| model4: | is this the line is busy now a few minutes |
| chunk3+: | i 'm on another phone now please wait a couple of minutes |

# Summary

- String-to-String translation model with hidden chunks

- More hidden variables
  $\longrightarrow$ More cost for training + decoding
  - Trainin Cost $\approx$ IBM Model 5 with pegging
  - Decoding Cost: moderate with Example-based scoring

- Quality Improvement: Slightly, but (probably) significant

# Summary

- String-to-String translation model with hidden chunks

- More hidden variables
  $\longrightarrow$ More cost for training + decoding
  - Trainin Cost $\approx$ IBM Model 5 with pegging
  - Decoding Cost: moderate with Example-based scoring

- Quality Improvement: Slightly, but (probably) significant

- Other approaches?

# Typology of Statistical Machine Translation

- Approach 1: Precomputation of Structure

- Approach 2: Structure-to-String

- Approach 3: Collection of Hierarchical FST

# Typology of Statistical Machine Translation

- Approach 1: Precomputation of Structure
- Approach 2: Structure-to-String
- Approach 3: Collection of Hierarchical FST

# Typology of Statistical Machine Translation

- **Approach 1: Precomputation of Structure**
  - Templates (Och et al. 1999)
  - Chunks from syntax-based phrase alignment (Watanabe et al. 2002)
  - Direct phrase induction (Marcu and Wong 2002)
    - Bias the training corpus by template, chunk or phrase
    - Works significantly better on observed word sequences, but not for unseen sequences
- Approach 2: Structure-to-String
- Approach 3: Collection of Hierarchical FST

# Typology of Statistical Machine Translation

- Approach 1: Precomputation of Structure

- Approach 2: Structure-to-String

- Approach 3: Collection of Hierarchical FST

# Typology of Statistical Machine Translation

- Approach 1: Precomputation of Structure

- Approach 2: Structure-to-String

  - ♦ Phrase-to-string Modeling (Wang 1998)

  - ♦ Syntax-to-string Modeling (Yamada and Knight 2001)

    - Bias the source part of a training corpus by "structure"
    - Computationally cheaper
    - Relies on the monolingual processing (parser or chunker)

- Approach 3: Collection of Hierarchical FST

# Typology of Statistical Machine Translation

- Approach 1: Precomputation of Structure

- Approach 2: Structure-to-String

- Approach 3: Collection of Hierarchical FST

# Typology of Statistical Machine Translation

- Approach 1: Precomputation of Structure

- Approach 2: Structure-to-String

- Approach 3: Collection of Hierarchical FST

  - ♦ (Alshawi et al. 2000)

    - Deterministic vs. Non-Deterministic
    - Faster decoding + less space
      vs. Slow decoding + pruning
    - Limited domain vs. Larger domain