

Breaking down the Language Barrier with Statistical Machine Translation: 2) Alignment/Phrase Extraction

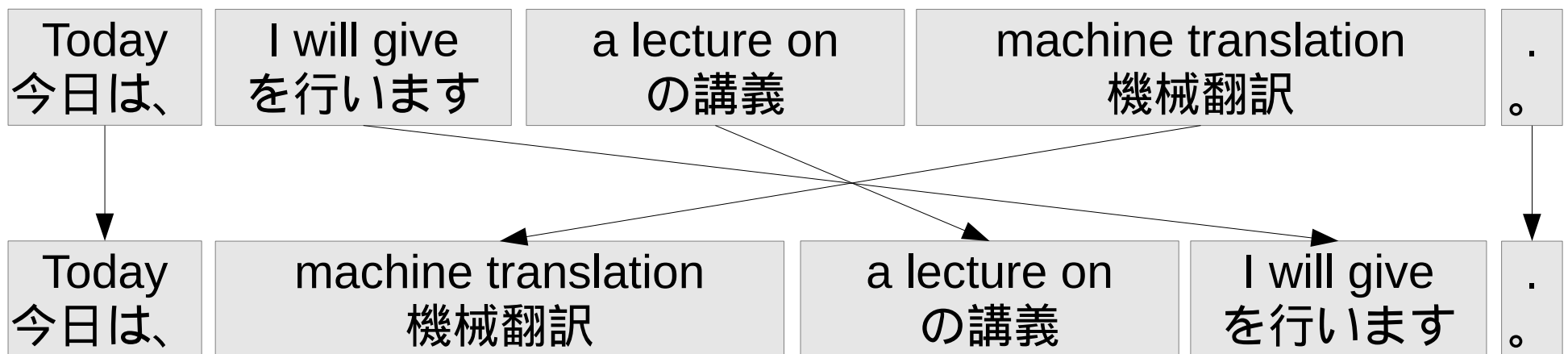
<http://www.phontron.com/class/sentan2014>

Advanced Research Seminar I/III
Graham Neubig
2014-1-30

How does machine translation work?

- Divide sentence into translatable patterns, reorder, combine

Today I will give a lecture on machine translation .



今日は、機械翻訳の講義を行います。

Problem

- There are millions of possible translations!

花子 が 太郎 に 会った

Hanako met Taro

Hanako met to Taro

Hanako ran in to Taro

Taro met Hanako

The Hanako met the Taro

- How do we tell which is better?

Statistical Machine Translation

- Translation model:

$P(\text{“今日”} \mid \text{“today”}) = \text{high}$

$P(\text{“今日は、”} \mid \text{“today”}) = \text{medium}$

$P(\text{“昨日”} \mid \text{“today”}) = \text{low}$

- Reordering Model:

$P(\begin{array}{cc} \text{鶏} & \text{食べる} \\ \text{が} & \\ \hline \text{chicken} & \text{eats} \end{array}) = \text{high}$

$P(\begin{array}{cc} \text{鶏} & \text{食べる} \\ \text{を} & \\ \hline \text{eats} & \text{chicken} \end{array}) = \text{high}$

$P(\begin{array}{cc} \text{鶏} & \text{食べる} \\ \text{が} & \\ \hline \text{eats} & \text{chicken} \end{array}) = \text{low}$

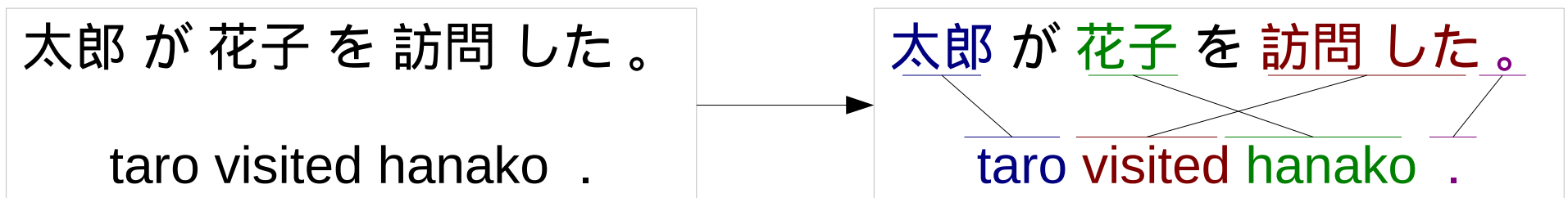
- Language Model:

$P(\text{“Taro met Hanako”}) = \text{high}$

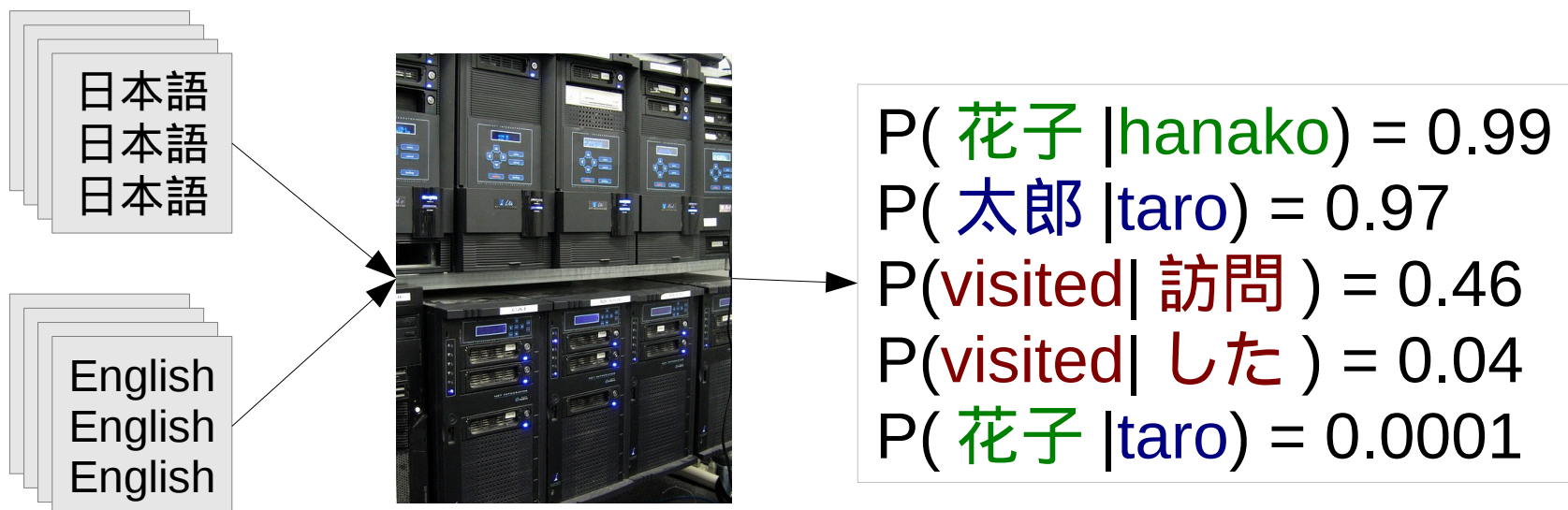
$P(\text{“the Taro met the Hanako”}) = \text{low}$

Alignment

- Find the correspondences between words



- Unsupervised probabilistic models are most common



Assignment

- I have given you code that:
 - 1: Calculates the **Dice coefficient** for each pair of words
 - 2: Decides the alignment in each direction using Dice coefficient statistics and the **Max Score** criterion
 - 3: Combines these two alignments using the “**intersection**” heuristic
 - 4: Measures alignment **F-measure**
- You should
 - 1: **Measure the alignment accuracy** on the provided data
 - 2: **Make one change** that improves the accuracy

Heuristic Alignment

How Do we Learn Alignments?

- For example, we go to an Italian restaurant w/ Japanese menu

チーズムース
Mousse di formaggi

タリアテッレ 4種のチーズソース
Tagliatelle al 4 formaggi

本日の鮮魚
Pesce del giorno

鮮魚のソテー お米とグリーンピース添え
Filetto di pesce su "Risi e Bisi"

ドルチェとチーズ
Dolce e Formaggi

- Try to find the patterns!

How Do we Learn Alignments?

- For example, we go to an Italian restaurant w/ Japanese menu

チーズムース
Mousse di formaggi

タリアテッレ 4種のチーズソース
Tagliatelle al 4 formaggi

本日の鮮魚
Pesce del giorno

鮮魚のソテー お米とグリーンピース添え
Filetto di pesce su “Risi e Bisi”

ドルチェとチーズ
Dolce e Formaggi

- Try to find the patterns!

Co-occurrence

- Simplest measure of correlation: **co-occurrence**

チーズ ムース
Mousse di formaggi

タリアテッレ 4 種のチーズソース
Tagliatelle al 4 formaggi

本日の鮮魚
Pesce del giorno

鮮魚のソテー お米とグリーンピース 添え
Filetto di pesce su "Risi e Bisi"

ドルチェとチーズ
Dolce e Formaggi

Occurrence

$c(\text{チーズ}) = 3$

$c(\text{の}) = 3$

$c(\text{と}) = 2$

...

$c(\text{formaggi}) = 3$

$c(\text{pesce}) = 2$

$c(\text{e}) = 2$

...

Co-occurrence

$c(\text{チーズ}, \text{formaggi}) = 3$

$c(\text{チーズ}, \text{mousse}) = 1$

$c(\text{チーズ}, \text{di}) = 1$

$c(\text{チーズ}, \text{tagliatelle}) = 1$

...

Problem with Co-occurrence

- Co-occurrence favors frequent words

the banker met a tall man
銀行員が背の高い男に会った

a man ran out of the room
男が部屋から飛び出た

the young boy is good at soccer
あの男の子はサッカーが上手だ

the statue of liberty
自由の女神

he enjoys the olympics
彼はオリンピックが大好きだ

Co-occurrence

$$c(\text{the}, \text{男}) = 3$$

$$c(\text{man}, \text{男}) = 2$$

Dice Coefficient

- The dice coefficient penalizes highly frequent words

$$\text{dice}(e, f) = \frac{2 * c(e, f)}{c(e) + c(f)}$$

the banker met a tall man
銀行員が背の高い男に会った

a man ran out of the room
男が部屋から飛び出た

the young boy is good at soccer
あの男の子はサッカーが上手だ

the statue of liberty
自由の女神

he enjoys the olympics
彼はオリンピックが大好きだ

Dice Coefficient

$$\text{dice}(\text{the}, \text{男}) = \\ (2 * 3) / (5 + 3) = 0.75$$

$$\text{dice}(\text{man}, \text{男}) = \\ (2 * 2) / (2 + 3) = 0.80$$

Scores → Alignments

- Now, we need a way to change dice coefficients to alignments

	historical	cold	outbreaks
歴代	0.596	0.018	0.250
の	0.002	0.003	0.000
風邪	0.020	0.909	0.037
大	0.007	0.002	0.085
流行	0.025	0.010	0.240

	historical	cold	outbreaks
歴代	●		
の	●		
風邪		●	
大			●
流行			●

Max-Score

- Choose the best target word for each source

	historical	cold	outbreaks	
歴代	0.596	0.018	0.250	→ historical
の	0.002	0.003	0.000	→ cold
風邪	0.020	0.909	0.037	→ cold
大	0.007	0.002	0.085	→ outbreaks
流行	0.025	0.010	0.240	→ outbreaks

Threshold

- Choose a threshold, align all words over threshold

	historical	cold	outbreaks
歴代	0.596	0.018	0.250
の	0.002	0.003	0.000
風邪	0.020	0.909	0.037
大	0.007	0.002	0.085
流行	0.025	0.010	0.240

$t > 0.1$

Competitive Linking

- Pick the highest scored alignment in order, but do not allow the same word be aligned to others

	historical	cold	outbreaks
歴代	0.596	0.018	0.250
の	0.002	0.003	0.000
風邪	0.020	0.909	0.037
大	0.007	0.002	0.085
流行	0.025	0.010	0.240

1. 風邪 → cold
2. 歴代 → historical
3. 流行 → outbreaks

Probabilistic Alignment: IBM Model 1

Probabilistic Alignment

- Create probabilistic model of two sentences

F= チーズ ムース

E= mousse di formaggi

$$P(F | E; M)$$

- The model M has some parameters (e.g. probabilities)

$$P(f= \text{チーズ} | e=\text{formaggi}) = 0.92$$

$$P(f= \text{チーズ} | e=\text{di}) = 0.001$$

$$P(f= \text{チーズ} | e=\text{mousse}) = 0.02$$

$$P(f= \text{ムース} | e=\text{formaggi}) = 0.07$$

$$P(f= \text{ムース} | e=\text{di}) = 0.002$$

$$P(f= \text{ムース} | e=\text{mousse}) = 0.89$$

- Probabilities are easier to understand (can make more refined models), easier to integrate

IBM Model One Idea

- Generate each word f_j in F according to
 - Pick one word index a_j at random ($P(a_j) = 1/(|E|+1)$), including special “NULL” word for unaligned words
 - Pick the word f_j according to $P(f | e_{a_j})$

Generate two words

チーズ

ムース

Choose: チーズ ($P(f|e) = 0.92$)

Choose: ムース ($P(f|e) = 0.89$)

Choose: $a_1 = 3$ ($P(a_1=3) = 0.25$)

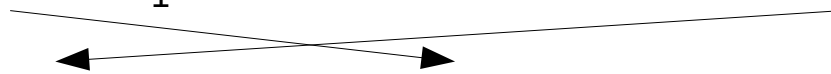
Choose: $a_2 = 1$ ($P(a_2=1) = 0.25$)

mousse

di

formaggi

NULL



IBM Model 1 Equation

- If we guess the alignment and output words, probability is

$$P(F, A|E) = \prod_{j=1}^J \frac{1}{I+1} P(f_j|e_{a_j})$$

Alignment Word translation

- We can also sum over all alignments

$$\begin{aligned} P(F|E) &= \sum_A \prod_{j=1}^J \frac{1}{I+1} P(f_j|e_{a_j}) \\ &= \prod_{j=1}^J \frac{1}{I+1} \sum_{i=1}^{I+1} P(f_j|e_i) \end{aligned}$$

Training Model 1

- We need to estimate our model parameters
- The most natural way to do so is with maximum likelihood

$$\hat{M} = \operatorname{argmax}_M P(F|E)$$

- How do we calculate the parameters to maximize likelihood?

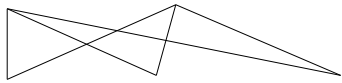
EM Algorithm

- IBM models trained using the **EM (Expectation-Maximization) algorithm**
- Idea:
 - **E Step:** Based on the model, estimate how many times each word e is translated into f
 - **M Step:** Based on the estimated counts, re-calculate the model parameters
- Every iteration, we increase the likelihood of the model

EM Example

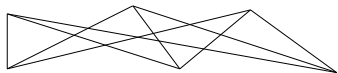
- Initialize: Simply count co-occurrence

チーズ ムース



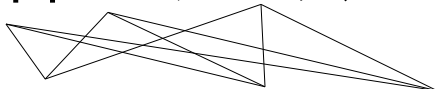
Mousse di formaggi

本日の 鮮魚



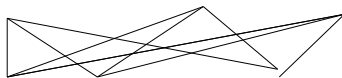
Pesce del giorno

本日の チーズ



Formaggi del giorno

ドルチェ と チーズ

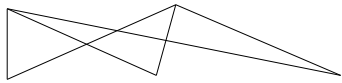


Dolce e Formaggi

EM Example

- M Step: Update the probabilities

チーズ ムース

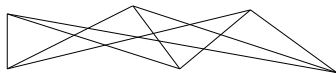


Mousse di formaggi

$$P(\text{チーズ} | \text{formaggi}) = 0.375$$

$$P(\text{ムース} | \text{formaggi}) = 0.125$$

本日の鮮魚



Pesce del giorno

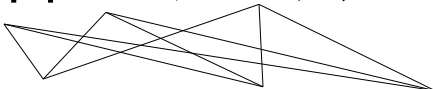
$$P(\text{本日} | \text{formaggi}) = 0.125$$

$$P(\text{の} | \text{formaggi}) = 0.125$$

$$P(\text{ドルチェ} | \text{formaggi}) = 0.125$$

$$P(\text{と} | \text{formaggi}) = 0.125$$

本日のチーズ



Formaggi del giorno

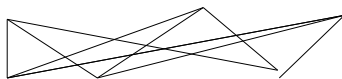
$$P(\text{本日} | \text{giorno}) = 0.33$$

$$P(\text{の} | \text{giorno}) = 0.33$$

$$P(\text{鮮魚} | \text{giorno}) = 0.16$$

$$P(\text{チーズ} | \text{giorno}) = 0.16$$

ドルチェ と チーズ



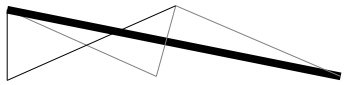
Dolce e Formaggi

...

EM Example

- E Step: Calculate connections between words

チーズ ムース



Mousse di formaggi

$$P(\text{チーズ} | \text{formaggi}) = 0.375$$

$$P(\text{ムース} | \text{formaggi}) = 0.125$$

本日の鮮魚



Pesce del giorno

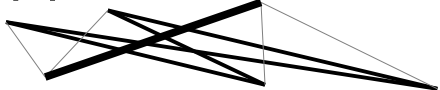
$$P(\text{本日} | \text{formaggi}) = 0.125$$

$$P(\text{の} | \text{formaggi}) = 0.125$$

$$P(\text{ドルチェ} | \text{formaggi}) = 0.125$$

$$P(\text{と} | \text{formaggi}) = 0.125$$

本日のチーズ



Formaggi del giorno

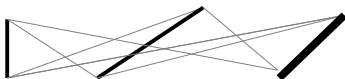
$$P(\text{本日} | \text{giorno}) = 0.33$$

$$P(\text{の} | \text{giorno}) = 0.33$$

$$P(\text{鮮魚} | \text{giorno}) = 0.16$$

$$P(\text{チーズ} | \text{giorno}) = 0.16$$

ドルチェ と チーズ



Dolce e Formaggi

...

EM Example

- M Step: Update the probabilities

~~チーズ ムース~~

Mousse di formaggi

$$P(\text{チーズ} | \text{formaggi}) = 0.9$$

$$P(\text{ムース} | \text{formaggi}) = 0.02$$

~~本日の鮮魚~~

Pesce del giorno

$$P(\text{本日} | \text{formaggi}) = 0.02$$

$$P(\text{の} | \text{formaggi}) = 0.02$$

$$P(\text{ドルチェ} | \text{formaggi}) = 0.02$$

$$P(\text{と} | \text{formaggi}) = 0.02$$

~~本日のチーズ~~

Formaggi del giorno

$$P(\text{本日} | \text{giorno}) = 0.48$$

$$P(\text{の} | \text{giorno}) = 0.48$$

$$P(\text{鮮魚} | \text{giorno}) = 0.02$$

$$P(\text{チーズ} | \text{giorno}) = 0.02$$

~~ドルチェ と チーズ~~

Dolce e Formaggi

...

EM Example

- E Step: Calculate connections between words

チーズ ムース

Mousse di formaggi

$$P(\text{チーズ} | \text{formaggi}) = 0.9$$

$$P(\text{ムース} | \text{formaggi}) = 0.02$$

本日の鮮魚

Pesce del giorno

$$P(\text{本日} | \text{formaggi}) = 0.02$$

$$P(\text{の} | \text{formaggi}) = 0.02$$

$$P(\text{ドルチェ} | \text{formaggi}) = 0.02$$

$$P(\text{と} | \text{formaggi}) = 0.02$$

本日のチーズ

Formaggi del giorno

$$P(\text{本日} | \text{giorno}) = 0.48$$

$$P(\text{の} | \text{giorno}) = 0.48$$

$$P(\text{鮮魚} | \text{giorno}) = 0.02$$

$$P(\text{チーズ} | \text{giorno}) = 0.02$$

ドルチェ と チーズ

Dolce e Formaggi

...

Initialization in Equations/Words

- Define our expected counts of two words x and y being aligned as

$$q(e = x, f = y)$$

- Initialize this to the co-occurrence counts

$$q(e = x, f = y) = c(e = x, f = y)$$

M Step in Equations/Words

- M Step: Calculate the model parameters
 - Simply calculate the maximum likelihood estimate of the conditional probability

$$P(f = y | e = x) = \frac{q(e = x, f = y)}{q(e = x)}$$

where

$$q(e = x) = \sum_y q(e = x, f = y)$$

E Step in Equations/Words

- E Step: Calculate expectations with model fixed
 - One sentence, probability of $a_j=i$ given model is:

$$P(a_j=i|F, E, M) = \frac{1}{I+1} P(f_j|e_i) \bigg/ \sum_{\tilde{i}=1}^{I+1} \frac{1}{I+1} P(f_j|e_{\tilde{i}})$$

↖ Current word ↖ All words

$$P(a_j=i|F, E, M) = P(f_j|e_i) \bigg/ \sum_{\tilde{i}=1}^{I+1} P(f_j|e_{\tilde{i}})$$

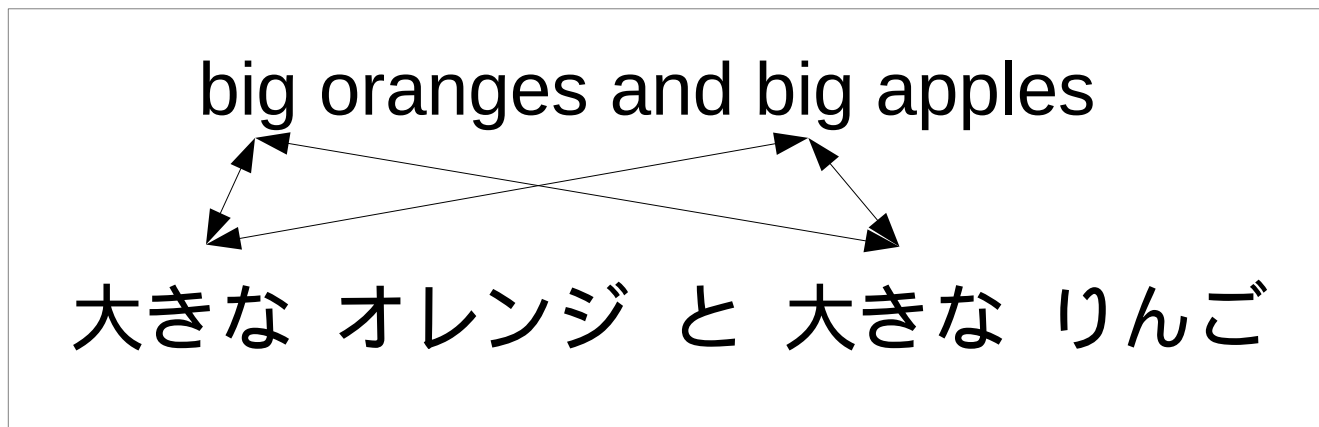
- Given this, calculate the expected count of translations (δ = Kronecker's delta, 1 when true, 0 otherwise)

$$q(e=x, f=y) = \sum_{E, F} \sum_{i=1}^{I+1} \sum_{j=1}^J P(a_j=i|F, E, M) \delta(e_i=x, f_j=y)$$

Probabilistic Alignment: IBM Models 2-5/HMM

Model 1 Problem

- Model 1 cannot handle word order!

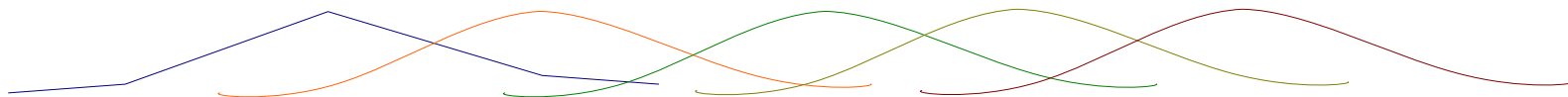


- Many new models proposed, all to handle this problem

Model 2 Idea

- “Words in both languages should be in about the same order”

big oranges and big apples

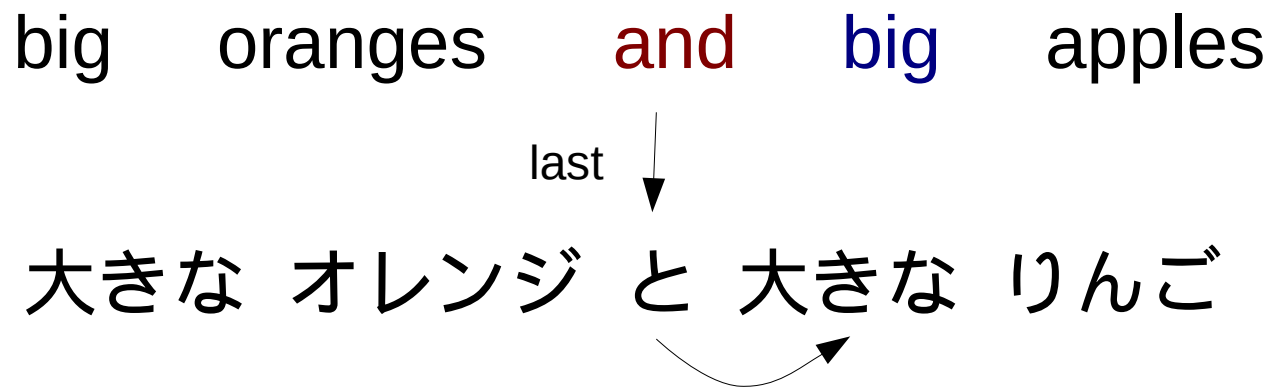


大きな オレンジ と 大きな りんご

- But, in most languages this is not true

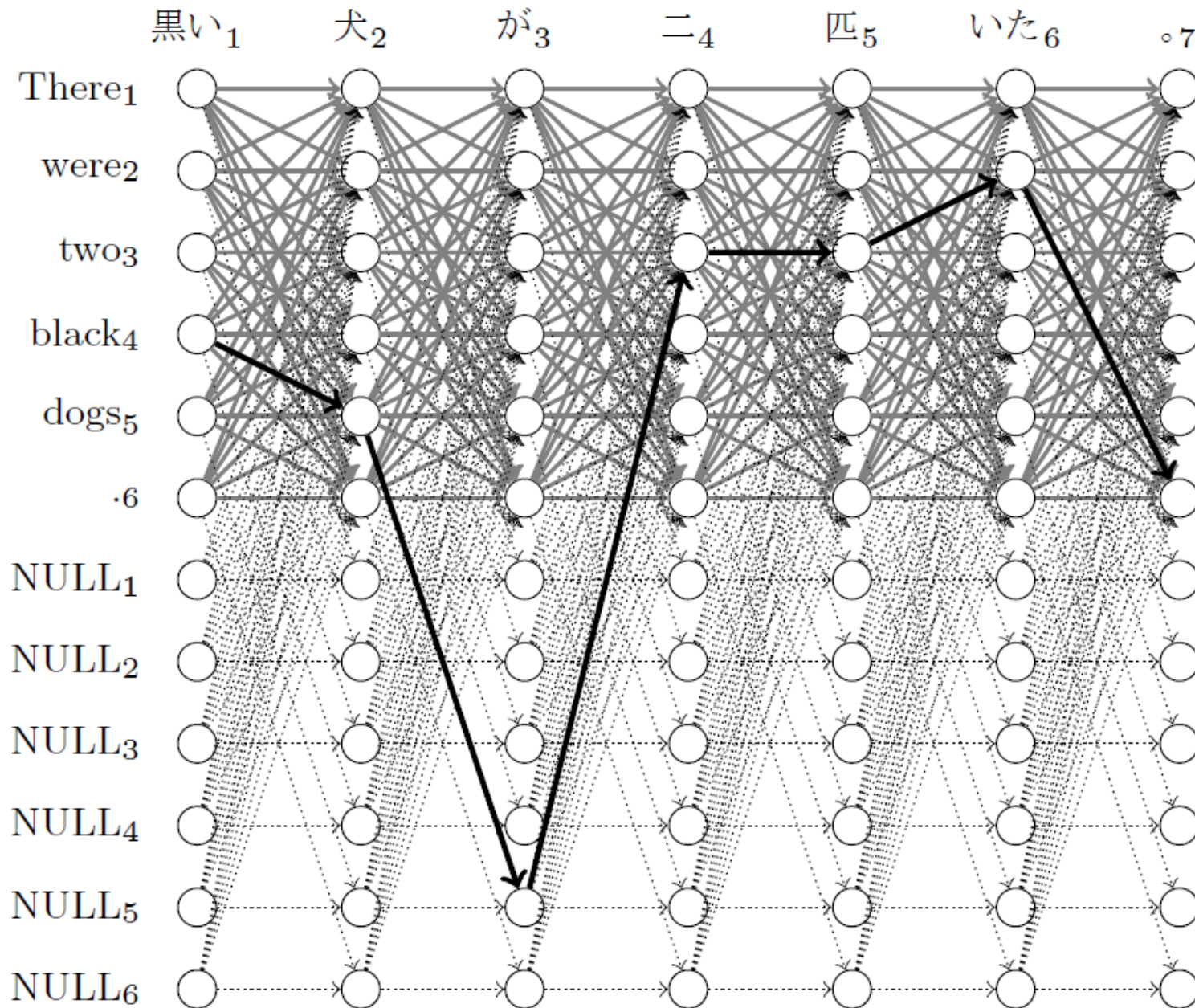
HMM Idea

- “Words that were previously aligned should give an idea about the next word”



- More effective than 2

HMM Graph



IBM Models 3-5

- Remaining IBM models are based on “fertility”
- “Given this word, how many words will it produce?”

Fertility 1

I

私
僕
俺

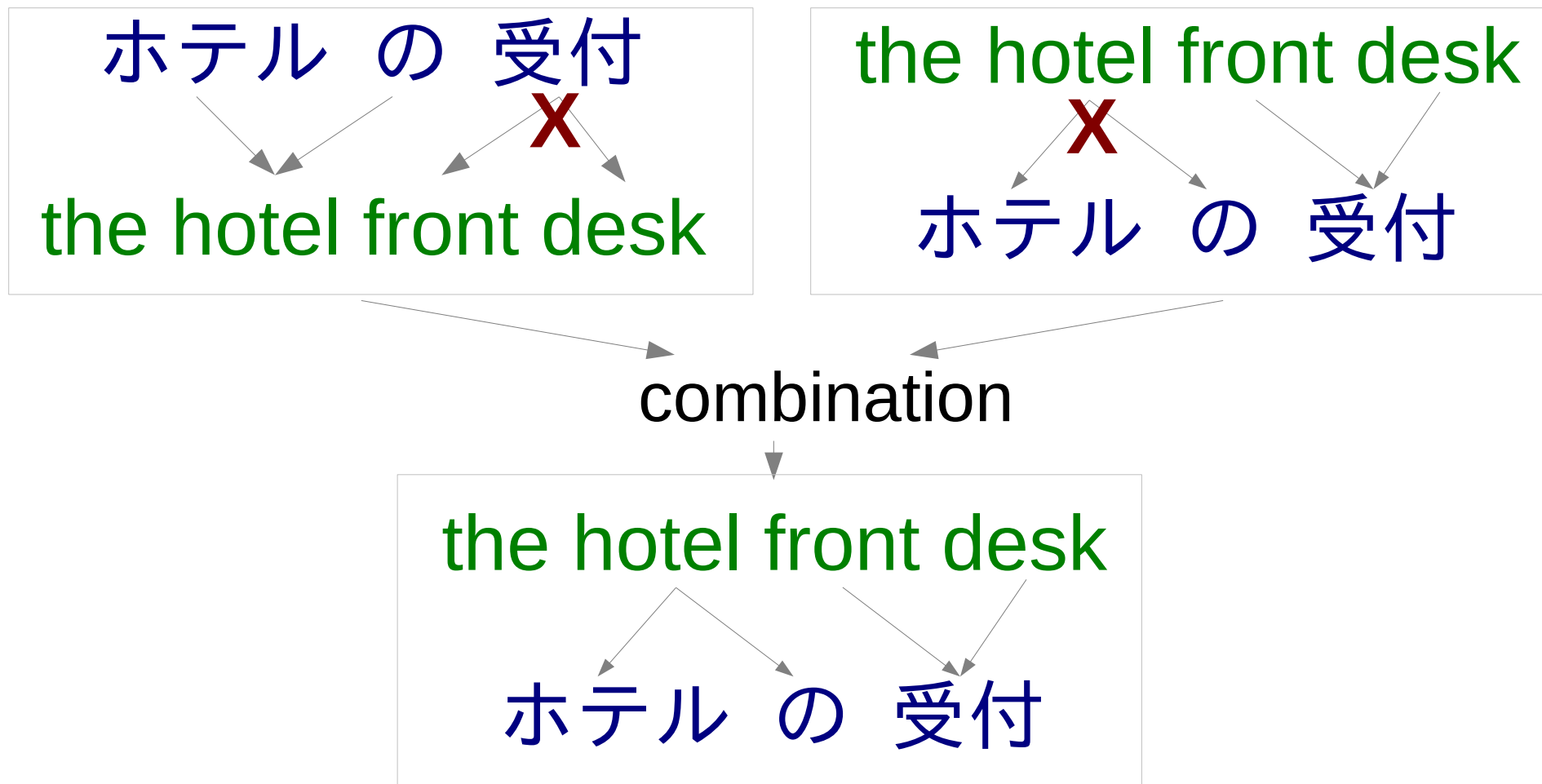
Fertility 3.5

adopted

採用 され た
養子 になっ た

Combining Alignments

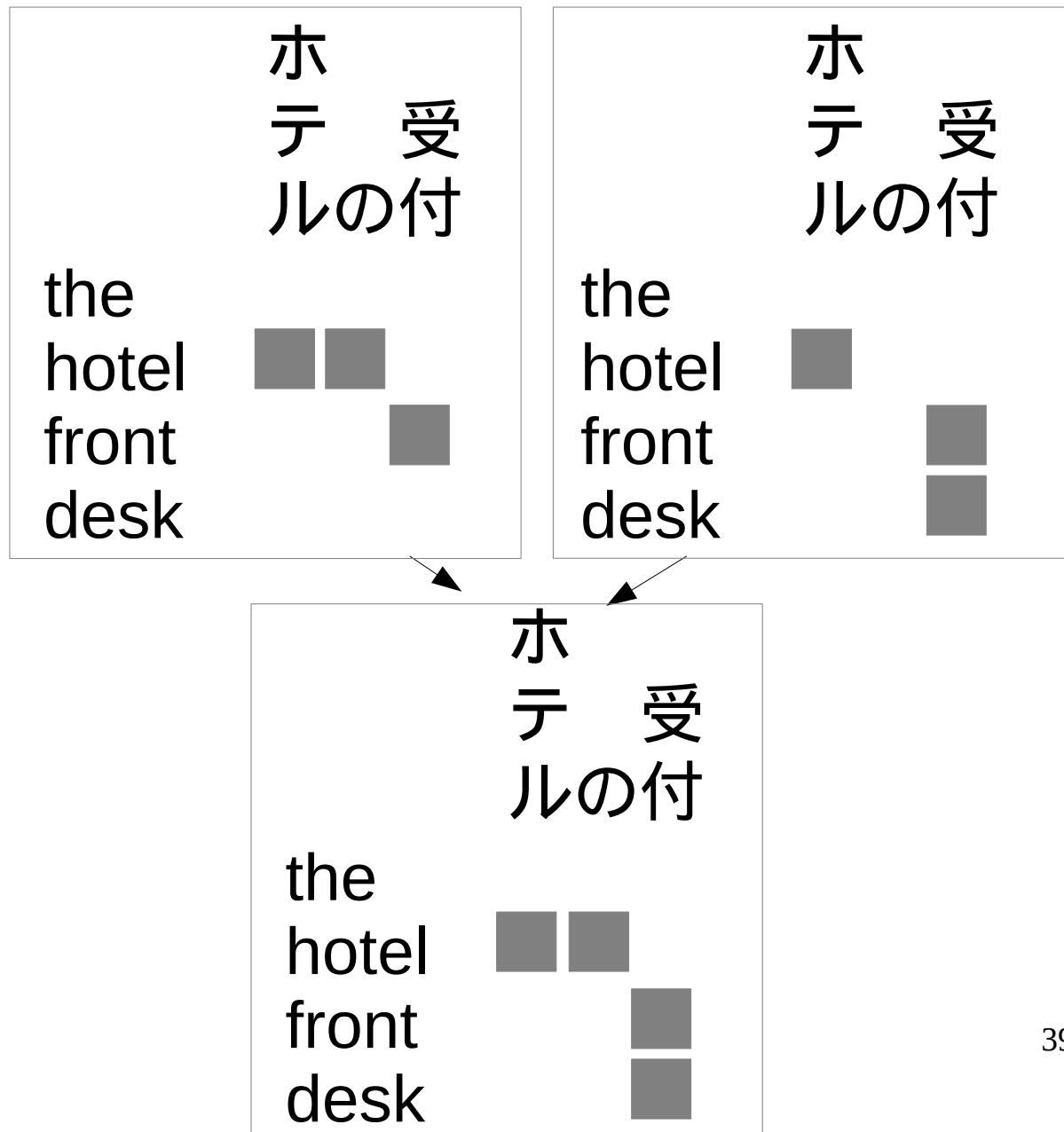
Combining 1-to-many Alignments



- Various heuristics for combination

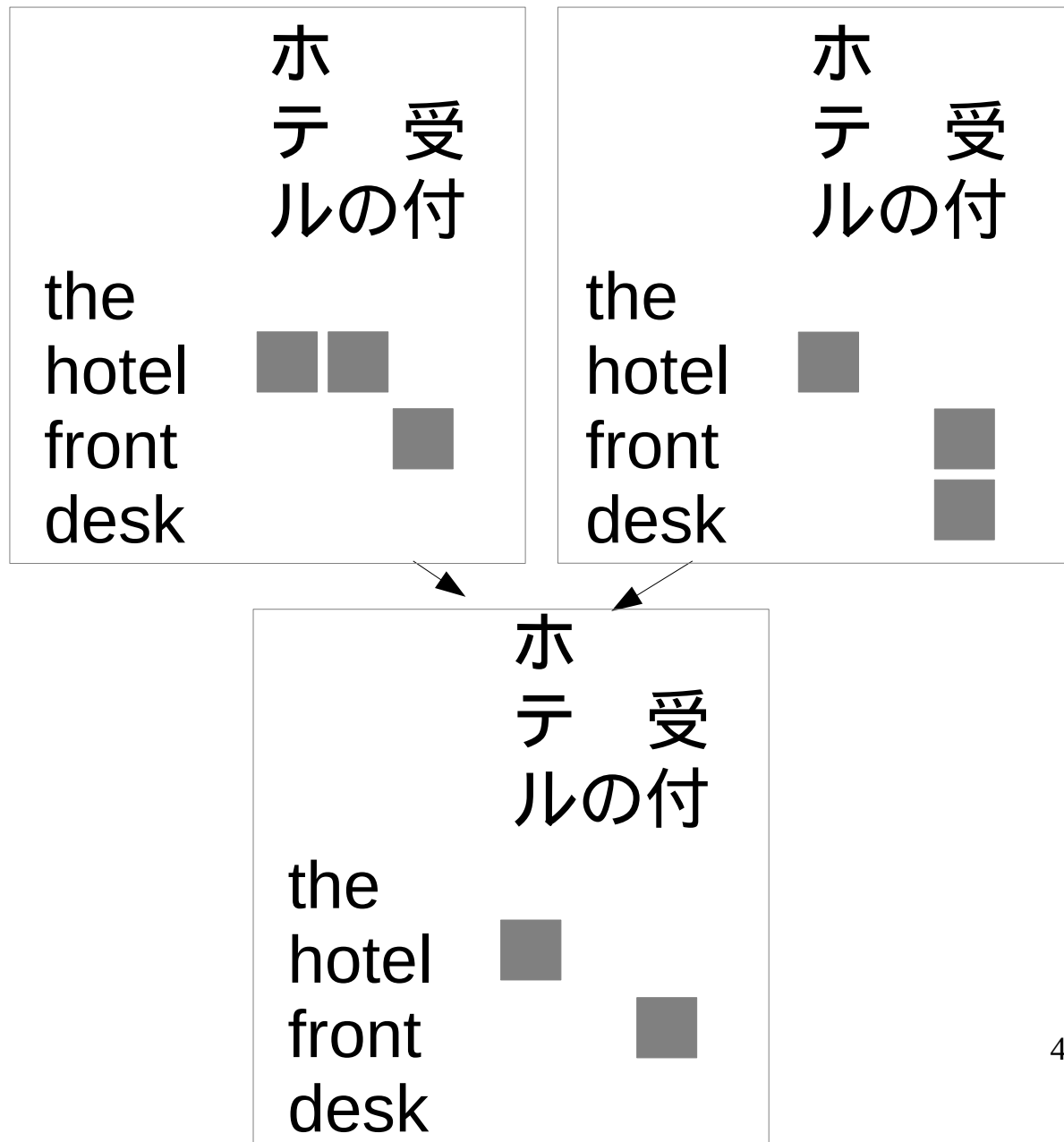
Union

- Union: Take all alignments in either direction



Intersection

- Intersection:
Only take
alignments in
one direction



Grow/Diag

- Grow: Use intersection, but add corner pieces included in the union



- Diag: Use intersection, but add corners

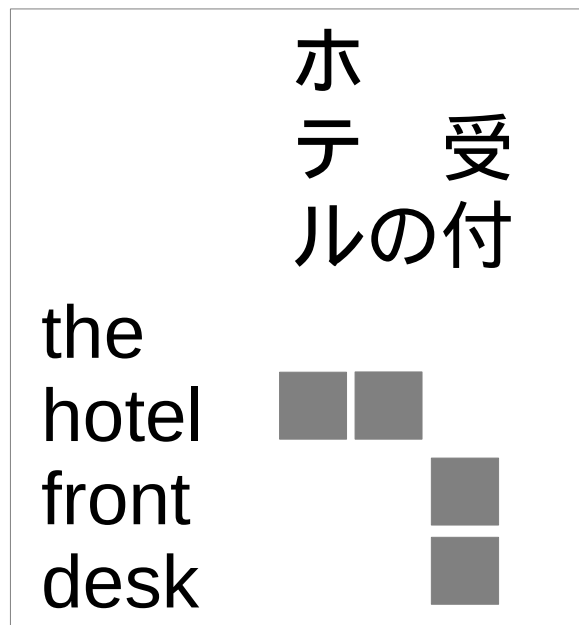


Evaluating Alignments

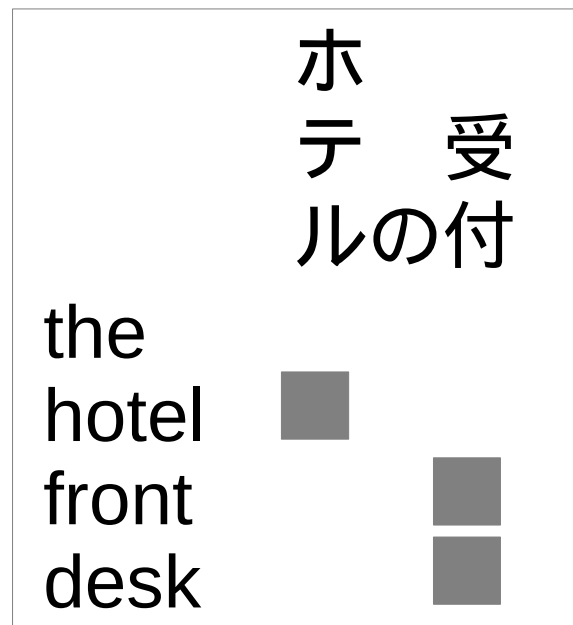
How Good are Our Alignments?

- Let's say we have two systems, which is better?

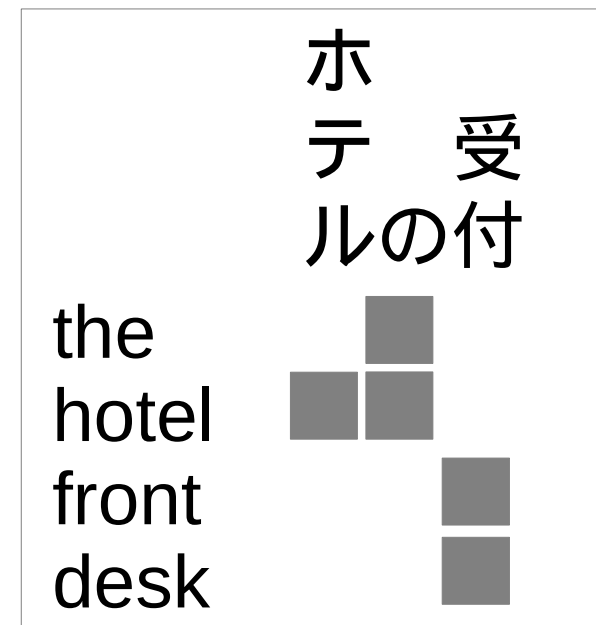
Reference



System A



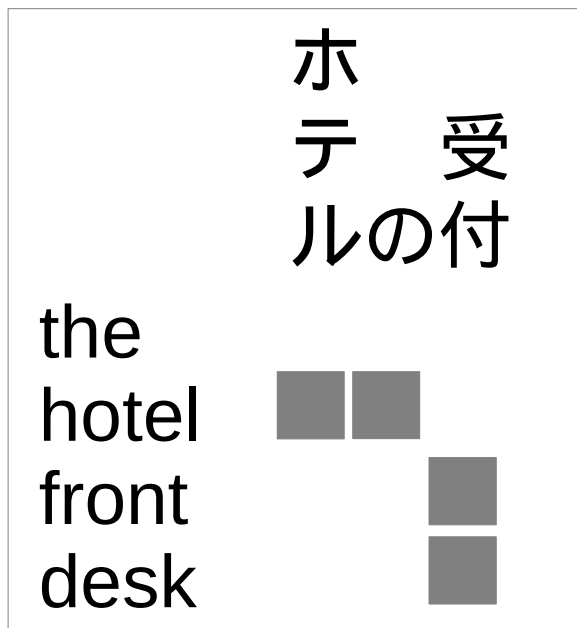
System B



Precision/Recall/F-Measure

- **Precision:** ratio of system alignments correct
- **Recall:** ratio of reference alignments correct
- **F-measure:** geometric mean of precision and recall

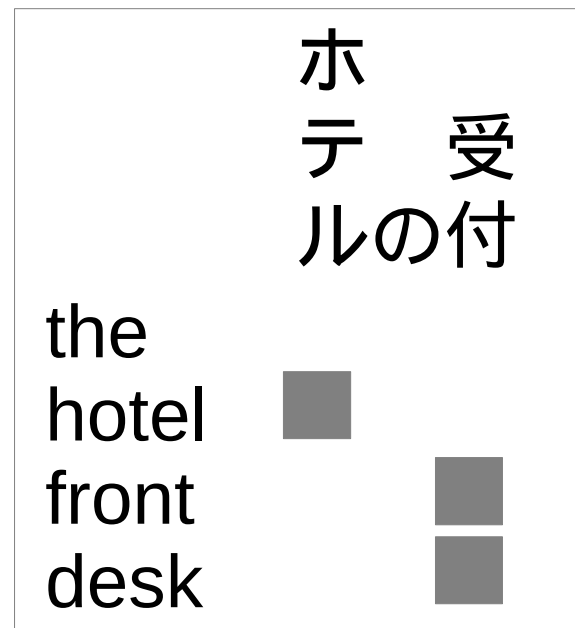
Reference (Handmade)



P=1.0

$F=2*1.0*0.75/(1.0+0.75)=0.85$

System A

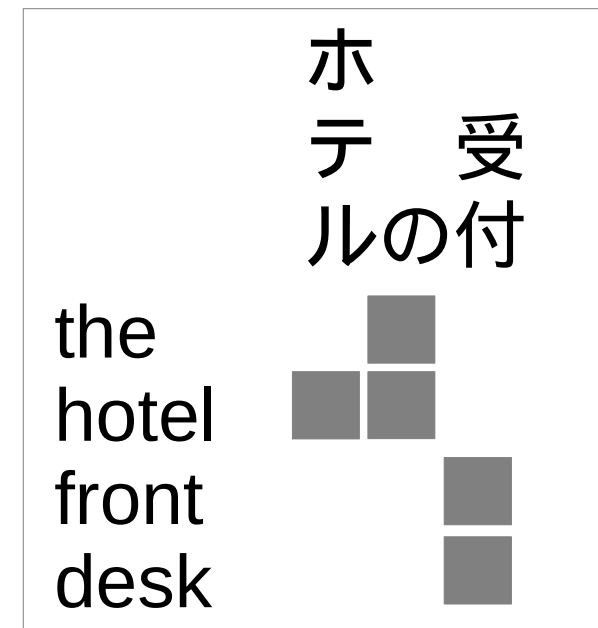


R=0.75

P=0.8

$F=2*0.8*1.0/(0.8+1.0)=0.88$

System B

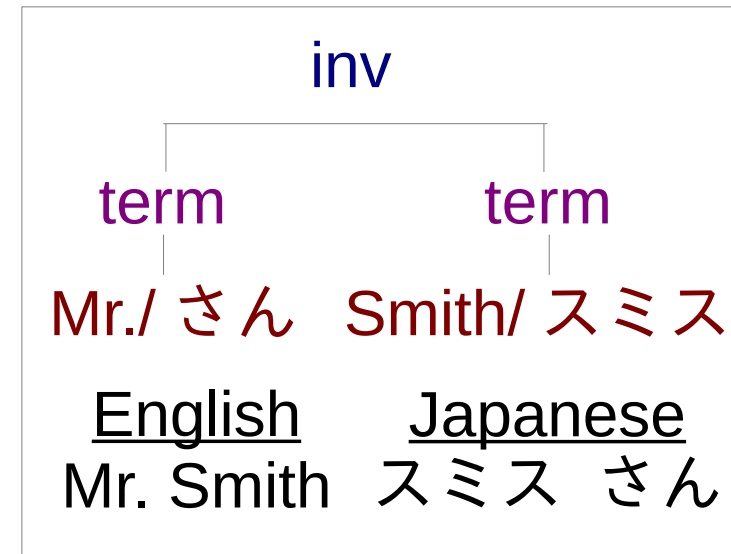
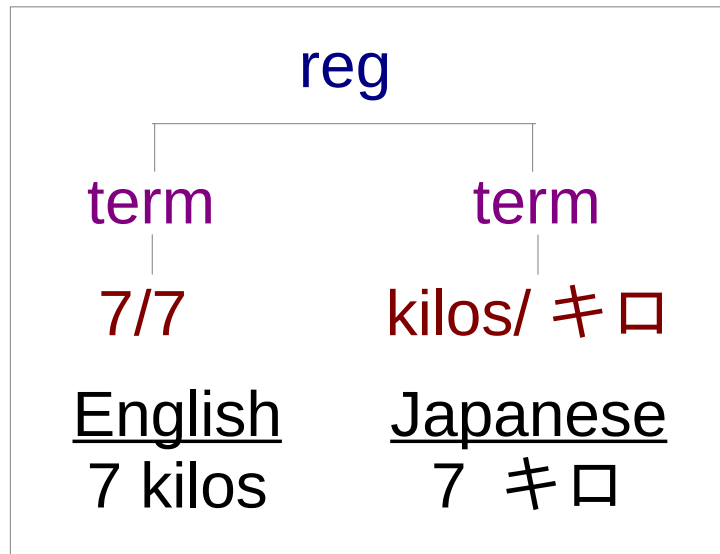


R=1.0

Advanced Topics

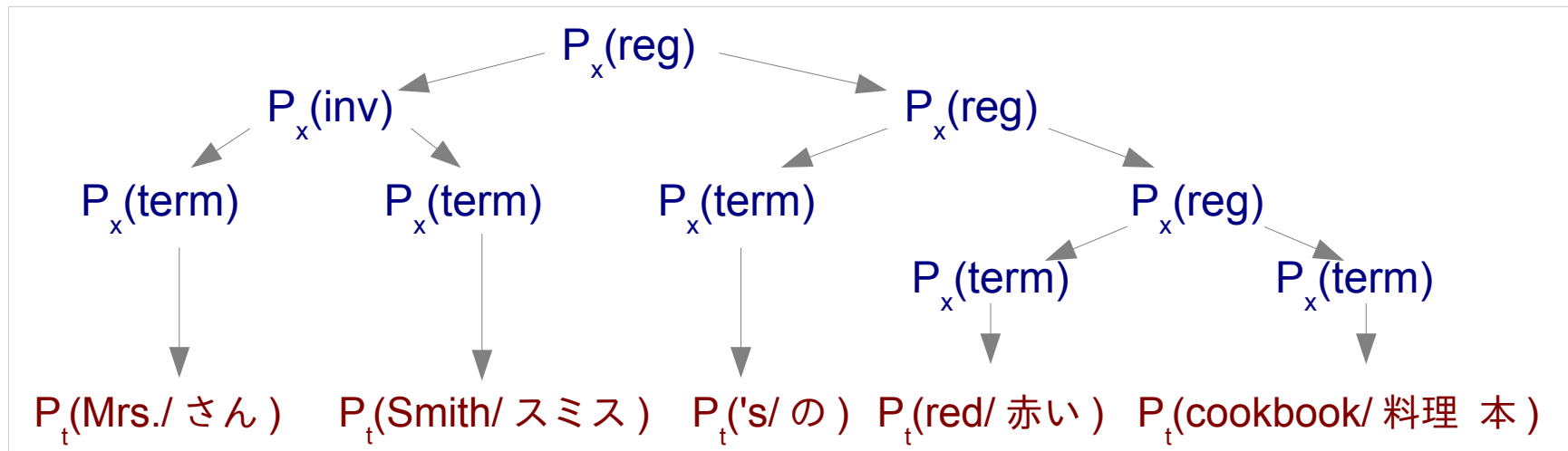
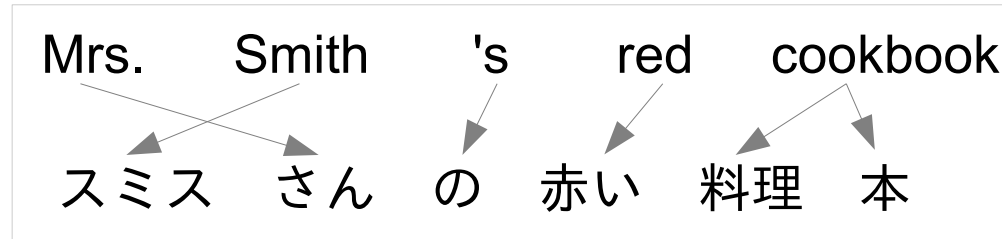
Inversion Transduction Grammar (ITG)

- A context-free grammar over two languages
 - **non-terminal** “regular (reg)” and “inverted (inv)”
 - one **pre-terminal** 「 term 」
 - **terminals** represent a phrase pair



Parsing with ITGs

- Define a probability over non-terminals/pre-terminals P_x , define another over terminals P_t



- Can use a variant of the CKY algorithm (standard algorithm for parsing context-free grammars)

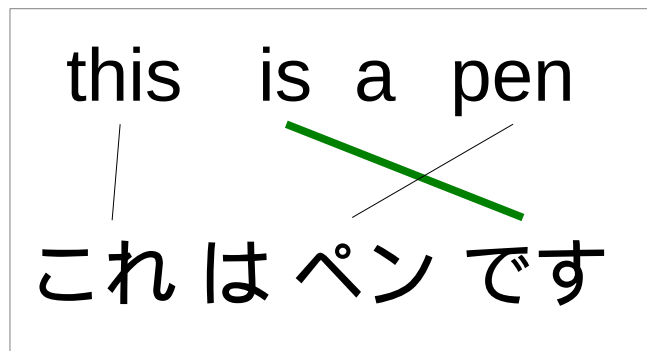
Advantages/Disadvantages of ITGs

- **Advantage:**
 - Simple model, no heuristics
 - Can handle many-to-many alignments in polynomial ($O(n^6)$) time
- **Disadvantage:**
 - Compared to other models, this is slow, especially for longer sentences

Supervised Alignment

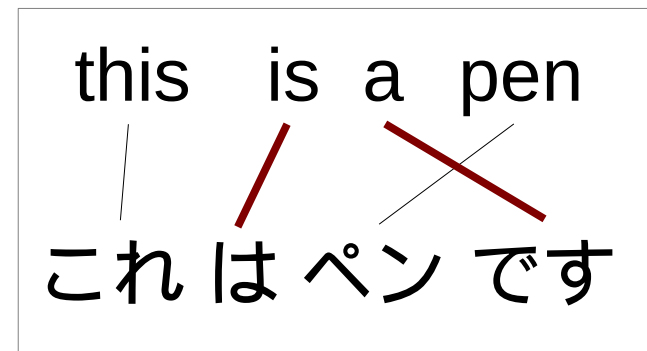
- Prepare a hand-made training corpus
- Use this corpus to train a model to “fix” unsupervised alignment

Reference



$c(\text{is, です})++$

Unsupervised

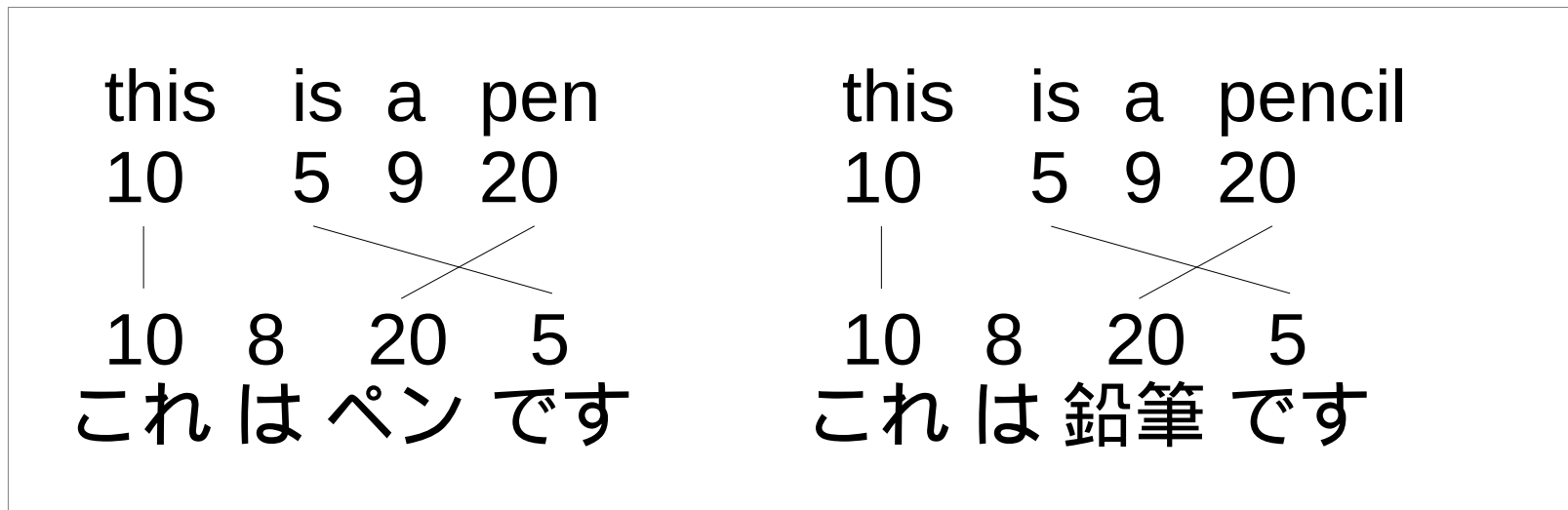


$c(\text{is, は})--$

$c(\text{a, です})--$

Class-Based Alignment

- Classes can be used to smooth the probability distribution

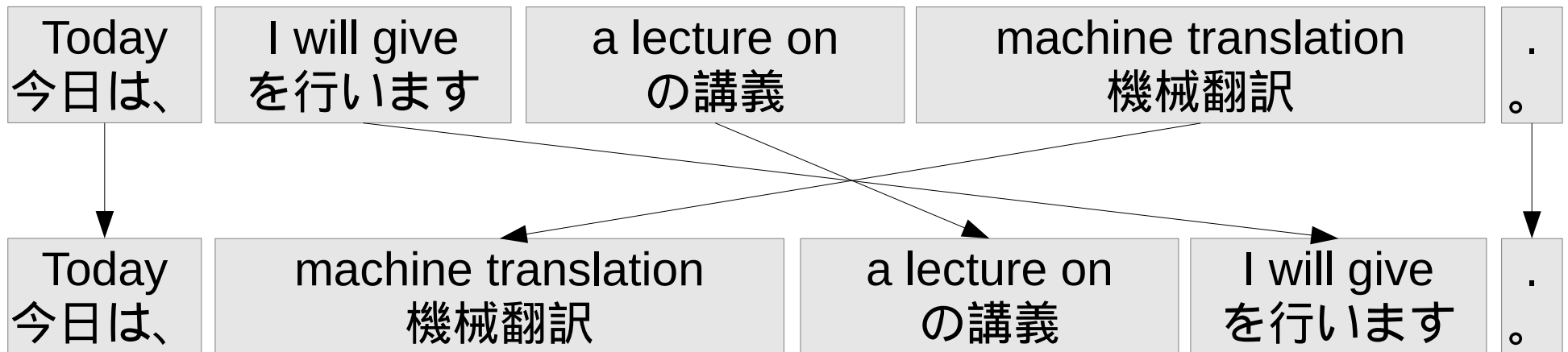


- Classes are learned automatically over both languages

Phrase Extraction/Scoring

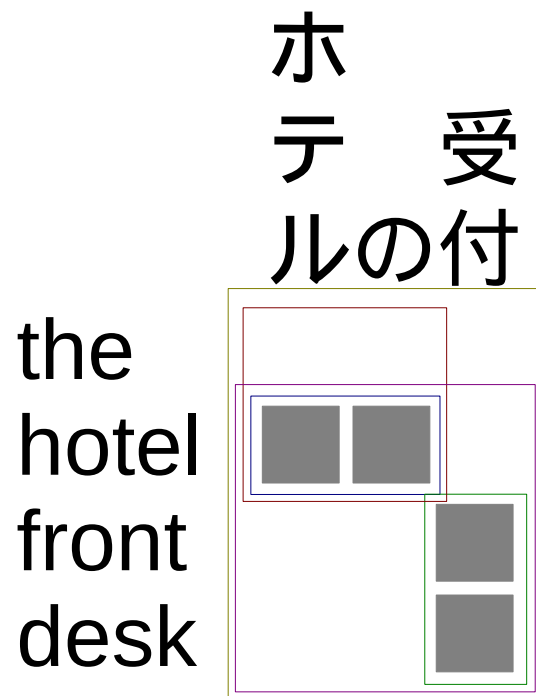
What is a “Phrase”

- In linguistics a “phrase” has a special meaning (“noun phrase, verb phrase”)
- In translation, it generally just means a string of words



Phrase Extraction

- We extract bilingual phrases from data



ホテルの → hotel

ホテルの → the hotel

受付 → front desk

ホテルの受付 → hotel front desk

ホテルの受付 → the hotel front desk

- Check all bilingual phrase pairs
- Extract only ones that have at least one good alignment, and no conflicts

Phrase Scoring

- We calculate probabilities for the phrases to measure how good a translation it is

- **Phrase translation probability (in both directions)**

$$P(\mathbf{f}|\mathbf{e}) = c(\mathbf{f},\mathbf{e})/c(\mathbf{e}) \quad P(\mathbf{e}|\mathbf{f}) = c(\mathbf{f},\mathbf{e})/c(\mathbf{f})$$

e.g.: $c(\text{ホテル の}, \text{the hotel}) / c(\text{the hotel})$

- **Lexical probabilities**

- Calculate phrase probabilities using IBM model 1 probabilities, and the words contained within
- Helps relieve sparsity problems

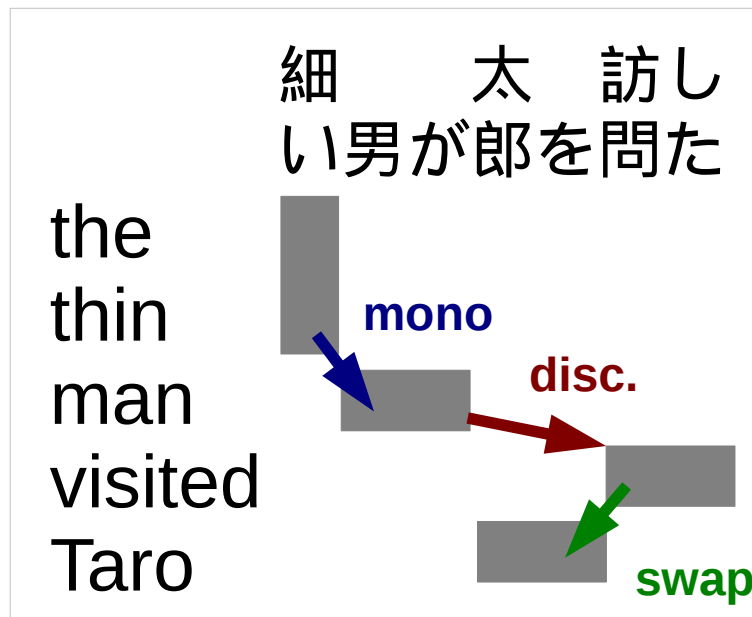
$$P(\mathbf{f}|\mathbf{e}) = \prod_f \frac{1}{|\mathbf{e}|} \sum_e P(\mathbf{f}|\mathbf{e})$$

e.g.:

$(P(\text{ホテル} | \text{the}) + P(\text{ホテル} | \text{hotel})) / 2 * (P(\text{の} | \text{the}) + P(\text{の} | \text{hotel})) / 2$

Reordering Varieties

- Probability of monotone, swap, discontinuous

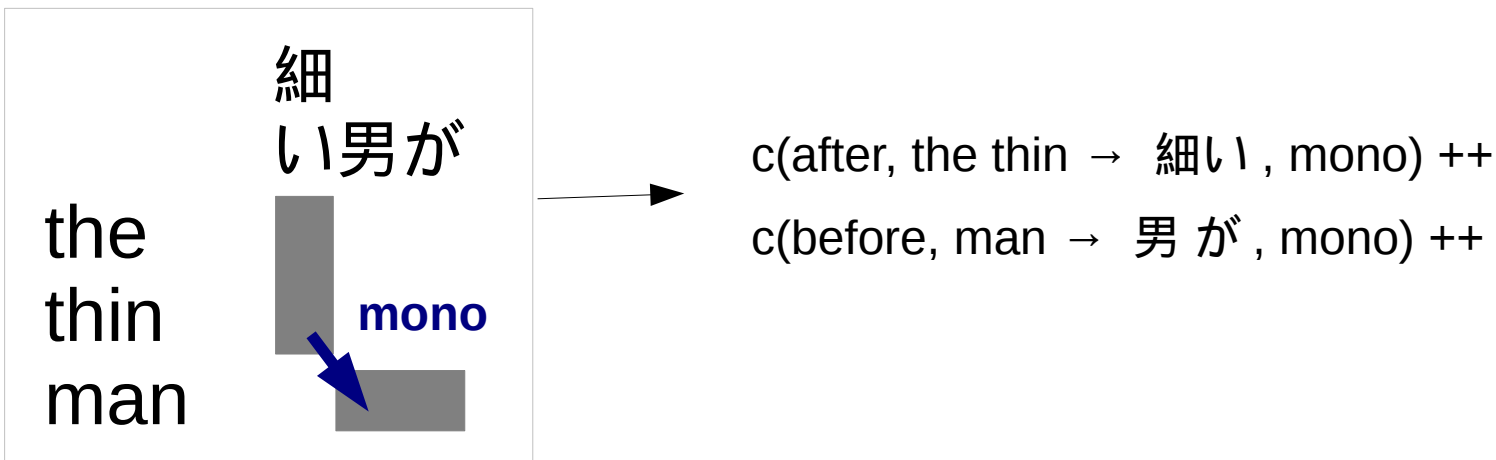


細い → the thin
high **monotone** probability

太郎 を → Taro
high **swap** probability

Reordering Probabilities

- Calculate reordering probabilities from data, like translation
- Condition on previous, next phrase



Assignment

Assignment

- I have given you code that:
 - 1: Calculates the **Dice coefficient** for each pair of words
 - 2: Decides the alignment in each direction using Dice coefficient statistics and the **Max Score** criterion
 - 3: Combines these two alignments using the “**intersection**” heuristic
 - 4: Measures alignment **F-measure**
- You should
 - 1: **Measure the alignment accuracy** on the provided data
 - 2: **Make one change** that improves the accuracy

Assignment Details

- Download the exercise from the web
- You can find a list of commands to run in `run-alignment.sh`
 - This is currently using smaller data (`kyoto-tunesmall`)
 - You can use larger data (`kyoto-tunetrain`) if you have a powerful computer (4GB memory, 5 minutes)
- Send your code, alignment f-measure before/after, a short description of the change, and a “username”
 - Due date: February 5th, 23:59
 - Address: neubig@is.naist.jp