# Machine Translation

Dan Klein
UC Berkeley

## Translation Task

- Text is both the input and the output.

- Input and output have roughly the same information content.

- Output is more predictable than a language modeling task.

- Lots of naturally occurring examples.

# Translation Examples

Republican leaders justified their policy by the need
to combat electoral fraud.

Die    Führungskräfte    der        Republikaner
 |          |              |              |
The    Executives        of the     republican

rechtfertigen    ihre    Politik     mit      der
      |            |        |          |        |
   justify       your    politics    With    of the

Notwendigkeit    ,    den    Wahlbetrug        zu
      |          |     |         |              |
    need         ,    the    election fraud    to

bekämpfen     .
    |         |
  fight       .

# Variety in Translations?

A small planet, whose is as big as could destroy a middle sized city, passed by the earth with a distance of 463 thousand kilometers. This was not found in advance. The astronomists got to know this incident 4 days later. This small planet is 50m in diameter. The astonomists are hard to find it for it comes from the direction of sun.

A volume enough to destroy a medium city small planet is big, flit earth within 463,000 kilometres of close however were not in advance discovered, astronomer just knew this matter after four days. This small planet diameter is about 50 metre, from the direction at sun, therefore astronomer very hard to discovers it.

An asteroid that was large enough to destroy a medium-sized city, swept across the earth at a short distance of 463,000 kilometers, but was not detected early. Astronomers learned about it four days later. The asteroid is about 50 meters in diameter and comes from the direction of the sun, making it difficult for astronomers to spot it.

From https://catalog.ldc.upenn.edu/LDC2003T17

# Evaluation

# BLEU Score

BLEU score: geometric mean of 1-, 2-, 3-, and 4-gram precision vs. a reference, multiplied by brevity penalty (harshly penalizes translations shorter than the reference).

$$\text{Matched}_i = \sum_{t_i} \min\left\{C_h(t_i), \max_j C_j(t_i)\right\}$$

If "of the" appears twice in hypothesis h but only at most once in a reference, then only the first is "correct"

$$P_i = \frac{\text{Matched}_i}{H_i}$$

"Clipped" precision of n-gram tokens

$$B = \exp\left\{\min\left(0, \frac{n-L}{n}\right)\right\}$$

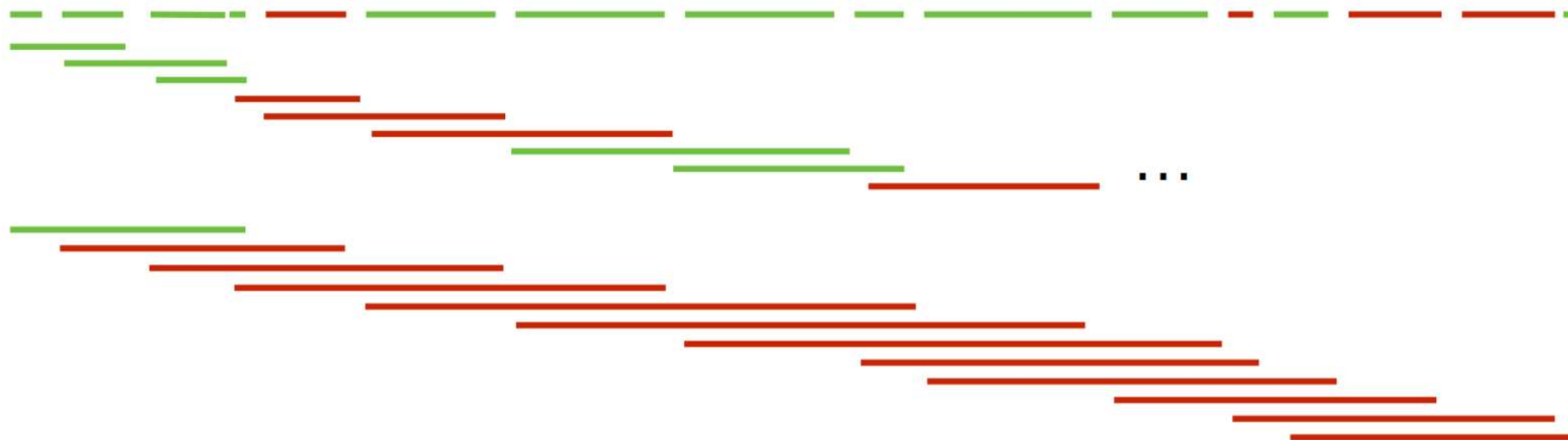Brevity penalty only matters if the hypothesis **corpus** is shorter than the sum of (shortest) references.

$$\text{BLUE} = B\left(\prod_{i=1}^{4} P_i\right)^{\frac{1}{4}}$$

BLEU is a mean of clipped precisions, scaled down by the brevity penalty.

# Evaluation with BLEU

In this sense, the measures will partially undermine the American democratic system.

In this sense, these measures partially undermine the democratic system of the United States.



BLEU = 26.52, 75.0/40.0/21.4/7.7 (BP=1.000, ratio=1.143, hyp_len=16, ref_len=14)

(Papineni et al., 2002) BLEU: a method for automatic evaluation of machine translation.

# Corpus BLEU Correlations with Average Human Judgments

These are ecological correlations over multiple segments; segment-level BLEU scores are noisy.

Commercial machine translation providers seem to all perform human evaluations of some sort.

(Ma et al., 2019) Results of the WMT19 Metrics Shared Task: Segment-Level and Strong MT Systems Pose Big Challenges
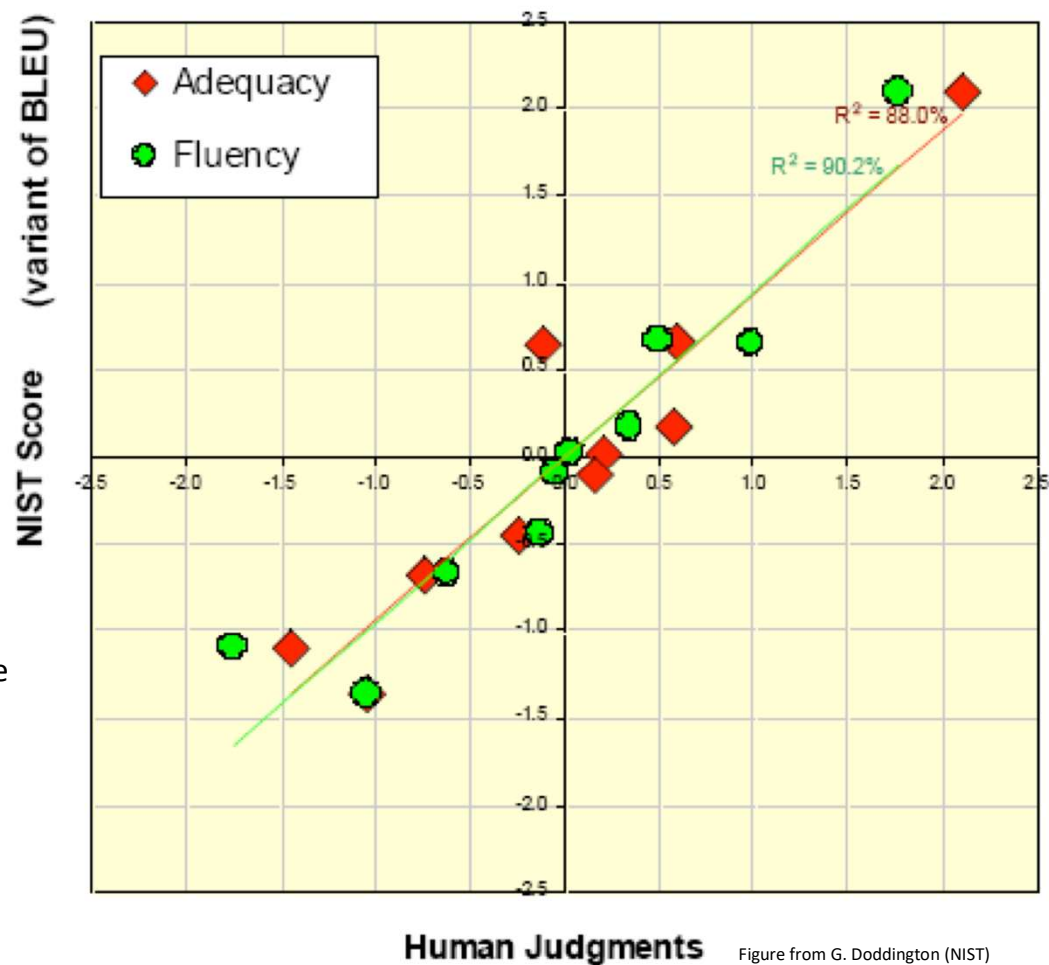


Figure from G. Doddington (NIST)

# Human Evaluations

**Direct assessment**: adequacy & fluency

- Monolingual: Ask humans to compare machine translation to a human-generated reference. (Easier to source annotators)

- Bilingual: Ask humans to compare machine translation to the source sentence that was translated. (Compares to human quality)

- Annotators can assess segments (sentences) or whole documents.

- Segments can be assessed with or without document context.

**Ranking assessment**:

- Raters are presented with 2 or more translations.

- A human-generated reference may be provided, along with the source.

- "In a pairwise ranking experiment, human raters assessing adequacy and fluency show a stronger preference for human over machine translation when evaluating documents as compared to isolated sentences." (Laubli et al., 2018)

**Editing assessment**: How many edits required to reach human quality

(Laubli et al., 2018) Has Machine Translation Achieved Human Parity? A Case for Document-level Evaluation

(Akhbardeh et al., 2021) Findings of the 2021 Conference on Machine Translation

# Translationese and Evaluation

Translated text can: (Baker et al., 1993; Graham et al., 2019)

- be more explicit than the original source

- be less ambiguous

- be simplified (lexically, syntactically, and stylistically)

- display a preference for conventional grammaticality

- avoid repetition

- exaggerate target language features

- display features of the source language

"If we consider only original source text (i.e. not translated from another language, or translationese), then we find evidence showing that human parity has not been achieved."
(Toral et al., 2018)

(Baker et al., 1993) Corpus linguistics and transla- tion studies: Implications and applications.
(Graham et al., 2019) Translationese in Machine Translation Evaluation.
(Toral et al, 2018) Attaining the Unattainable? Reassessing Claims of Human Parity in Neural Machine Translation

# How are We Doing?  Example: WMT 2019 Evaluation

2019 segment-in-context direct assessment (Barrault et al, 2019):

✓ German to English: many systems are tied with human performance;

× English to Chinese: all systems are outperformed by the human translator;

× English to Czech: all systems are outperformed by the human translator;

× English to Finnish: all systems are outperformed by the human translator;

✓ English to German: Facebook-FAIR achieves super-human translation performance; several systems are tied with human performance;

× English to Gujarati: all systems are outperformed by the human translator;

× English to Kazakh: all systems are outperformed by the human translator;

× English to Lithuanian: all systems are outperformed by the human translator;

✓ English to Russian: Facebook-FAIR is tied with human performance.

(Barrault et al, 2019) Findings of the 2019 Conference on Machine Translation (WMT19)

# Statistical Machine Translation
## (1990 - 2015)

When I look at an article in Russian, I say: "This is really written in English, but it has been coded in some strange symbols. I will now proceed to decode."

Warren Weaver (1949)

# Levels of Transfer: Vauquois Triangle (1968)



interlingua

semantics — semantics

syntax — syntax

phrases — phrases

words — words

SOURCE          TARGET

Yo lo haré mañana
I will do it tomorrow
NP
VP

$$P\left( \begin{array}{c} \text{VP} \\ \text{MD} \quad \text{VP} \\ \text{VB} \quad \text{PRN} \quad \text{NP} \\ \text{will} \quad \text{do} \quad \text{it} \end{array} \middle| \begin{array}{c} \text{VP} \\ \text{lo haré} \quad \text{NP} \end{array} \right) = 0.8$$

Yo lo haré mañana
I will do it tomorrow

| English (E) | P( E | lo haré ) |
|---|---|
| will do it | 0.8 |
| will do so | 0.2 |

Yo lo haré mañana

I will do it tomorrow

| English (E) | P( E | mañana ) |
|---|---|
| tomorrow | 0.7 |
| morning | 0.3 |

# Data-Driven Machine Translation

*Target language corpus gives examples of well-formed sentences*

| I will get to it later | See you later | He will do it |
|---|---|---|

*Parallel corpus gives translation examples*

| I will do it gladly |
| Yo lo haré de muy buen grado |

| You will see later |
| Después lo veras |

*Machine translation system:*

**Source language**

| Yo lo haré después |

NOVEL SENTENCE

Model of translation

**Target language**

| I will do it later |

# Stitching Together Fragments

*Parallel corpus gives translation examples*

S
NP        VP
     MD        VP
PRP    |    VB PRP ADV

I will do it gladly

Yo lo haré de muy buen grado

S
NP        VP
     MD        VP
PRP    |    VB  ADV

You will see later

Después lo veras

*Machine translation system:*

| S | |
|---|---|
| | ADV |
| Yo lo haré | después |

Model of translation

| S | |
|---|---|
| | ADV |
| I will do it | later |

# Evolution of the Noisy Channel Model

$$P(e|f) \propto P(f|e) \cdot P(e)$$

$$P(e|f) \propto P(f|e)^{\phi_{\text{tm}}} \cdot P(e)^{\phi_{\text{lm}}}$$

$$P(e|f) \propto \exp\left\{\sum_i w_i \cdot f_i(e, f)\right\}$$

Chosen to minimize loss

E.g., log P(e)

# Word Alignment and Phrase Extraction

# Extracting Translation Rules

# Counting Aligned Phrases

d'assister à la reunion et ||| to attend the meeting and

assister à la reunion ||| attend the meeting

la reunion et ||| the meeting and

nous ||| we

...

- Relative frequencies are the most important features in a phrase-based or syntax-based model.

- Scoring a phrase under a lexical model is the second most important feature.

- Estimation does not involve choosing among segmentations of a sentence into phrases.

# Translation Options

| er | geht | ja | nicht | nach | hause |
|---|---|---|---|---|---|

| he | is | yes | not | after | house |
| it | are | is | do not | to | home |
| , it | goes | , of course | does not | according to | chamber |
| , he | go | , | is not | in | at home |

| it is | | not | | home | |
| he will be | | is not | | under house | |
| it goes | | does not | | return home | |
| he goes | | do not | | do not | |

| | is | | to | |
| | are | | following | |
| | is after all | | not after | |
| | does | | not to | |

| | not | |
| | is not | |
| | are not | |
| | is not a | |

- Many translation options to choose from
  - in Europarl phrase table: 2727 matching phrase pairs for this sentence
  - by pruning to the top 20 per phrase, 202 translation options remain

# Decoding: Find Best Path

**er     geht     ja     nicht     nach     hause**

# Phrase-Based Decoding

这 | 7人 | 中包括 | 来自 | 法国 | 和 | 俄罗斯 | 的 | 宇航 | 员 | |
.

| the | 7 people | including | by some | | and | the russian | the | the astronauts | | , |
|-----|----------|-----------|---------|--|-----|-------------|-----|----------------|--|---|
| it | 7 people included | | by france | | and the | the russian | | international astronautical | of rapporteur . | |
| this | 7 out | including the | from | the french | and the russian | | the fifth | | . | |
| these | 7 among | including from | | the french and | | of the russian | of | space | members | . |
| that | 7 persons | including from the | | of france | and to | russian | of the | aerospace | members . | |
| | 7 include | | from the | of france and | | russian | | astronauts | | . the |
| | 7 numbers include | from france | | and russian | | of astronauts who | | | . " |
| | 7 populations include | those from france | | and russian | | astronauts . | | |
| | 7 deportees included | come from | france | and russia | | in | astronautical | personnel | ; |
| | 7 philtrum | including those from | france and | russia | a space | | member | |
| | | including representatives from | france and the | russia | | astronaut | | |
| | | include | came from | france and russia | | by cosmonauts | | |
| | | include representatives from | french | and russia | | cosmonauts | | |
| | | include | came from france | and russia 's | | cosmonauts . | | |
| | | includes | coming from | french and | russia 's | | cosmonaut | |
| | | | | french and russian | | 's | astronavigation | member . |
| | | | | french | and russia | astronauts | | |
| | | | | | and russia 's | | | special rapporteur |
| | | | | | , and | russia | | | rapporteur |
| | | | | | , and russia | | | rapporteur . |
| | | | | | , and russia | | | |
| | | | | | or | russia 's | | |

# Machine Translation

Dan Klein
UC Berkeley

*Many slides from John DeNero and Philip Koehn*

# Word Alignments

# Word Alignment

Given a sentence pair, which words correspond to each other?

# Word Alignment?



Is the English word does aligned to
the German wohnt (verb) or nicht (negation) or neither?

# Word Alignment?



How do the idioms kicked the bucket and biss ins grass match up?
Outside this exceptional context, bucket is never a good translation for grass

# Lexical Translation / Word Alignment Models

# Unsupervised Word Alignment

- Input: a *bitext*: pairs of translated sentences

> **nous acceptons votre opinion .**
>
> **we accept your view .**

- Output: *alignments*: pairs of translated words

  - When words have unique sources, can represent as a (forward) alignment function a from French to English positions

## Word Alignment

- Even today models are often built on the IBM alignment models

- Create probabilistic word-level translation models

- The models incorporate latent (unobserved) word alignments

- Optimize the probability of the observed words

- Use the imputed alignments to reveal word-level correspondence

- Throw out the translation models themselves

# Alignment

- In a parallel text (or when we translate), we align words in one language with the words in the other

<br>

<div align="center">

| 1 | 2 | 3 | 4 |
|:---:|:---:|:---:|:---:|
| das | Haus | ist | klein |
| \| | \| | \| | \| |
| the | house | is | small |
| 1 | 2 | 3 | 4 |

</div>

<br>

- Word positions are numbered 1–4

# Alignment Function

- Formalizing alignment with an alignment function

- Mapping an English target word at position $i$ to a German source word at position $j$ with a function $a : i \rightarrow j$

- Example
$$a : \{1 \rightarrow 1, 2 \rightarrow 2, 3 \rightarrow 3, 4 \rightarrow 4\}$$

# Reordering

Words may be reordered during translation



$$a : \{1 \rightarrow 3, 2 \rightarrow 4, 3 \rightarrow 2, 4 \rightarrow 1\}$$

# One-to-Many Translation

A source word may translate into multiple target words



$$a : \{1 \rightarrow 1, 2 \rightarrow 2, 3 \rightarrow 3, 4 \rightarrow 4, 5 \rightarrow 4\}$$

# Dropping Words

Words may be dropped when translated
(German article das is dropped)



$$a : \{1 \rightarrow 2, 2 \rightarrow 3, 3 \rightarrow 4\}$$

# Inserting Words

- Words may be added during translation

    - The English just does not have an equivalent in German
    - We still need to map it to something: special NULL token

| 0 | 1 | 2 | 3 | 4 |
|---|---|---|---|---|
| NULL | das | Haus | ist | klein |

the   house   is   just   small

1   2   3   4   5

$$a : \{1 \to 1, 2 \to 2, 3 \to 3, 4 \to 0, 5 \to 4\}$$

# IBM Model 1: Allocation

# IBM Model 1 (Brown 93)

- Alignments: a hidden vector called an *alignment* specifies which English source is responsible for each French target word.

$$a = a_1 \ldots a_J$$



And$_1$    the$_2$    program$_3$    has$_4$    been$_5$    implemented$_6$

$a_1 = 2$    $a_2 = 3$    $a_3 = 4$    $a_4 = 5$    $a_5 = 6$   $a_6 = 6$   $a_7 = 6$

Le$_1$    programme$_2$    a$_3$    été$_4$    mis$_5$    en$_6$    application$_7$

$$P(f, a | e) = \prod_j P(a_j = i) P(f_j | e_i)$$

$$= \prod_j \frac{1}{I+1} P(f_j | e_i)$$

$$P(f | e) = \sum_a P(f, a | e)$$

# Example

| das | |
|---|---|
| $e$ | $t(e\|f)$ |
| the | 0.7 |
| that | 0.15 |
| which | 0.075 |
| who | 0.05 |
| this | 0.025 |

| Haus | |
|---|---|
| $e$ | $t(e\|f)$ |
| house | 0.8 |
| building | 0.16 |
| home | 0.02 |
| household | 0.015 |
| shell | 0.005 |

| ist | |
|---|---|
| $e$ | $t(e\|f)$ |
| is | 0.8 |
| 's | 0.16 |
| exists | 0.02 |
| has | 0.015 |
| are | 0.005 |

| klein | |
|---|---|
| $e$ | $t(e\|f)$ |
| small | 0.4 |
| little | 0.4 |
| short | 0.1 |
| minor | 0.06 |
| petty | 0.04 |

$$p(e,a|f) = \frac{\epsilon}{4^3} \times t(\text{the}|\text{das}) \times t(\text{house}|\text{Haus}) \times t(\text{is}|\text{ist}) \times t(\text{small}|\text{klein})$$

$$= \frac{\epsilon}{4^3} \times 0.7 \times 0.8 \times 0.8 \times 0.4$$

$$= 0.0028\epsilon$$

# Expectation Maximization

# EM Algorithm

- Incomplete data

  - if we had *complete data*, would could estimate *model*
  - if we had *model*, we could fill in the *gaps in the data*

- Expectation Maximization (EM) in a nutshell

  1. initialize model parameters (e.g. uniform)
  2. assign probabilities to the missing data
  3. estimate model parameters from completed data
  4. iterate steps 2–3 until convergence

# EM Algorithm

```
... la maison ... la maison blue ... la fleur ...



... the house ... the blue house ... the flower ...
```

- Initial step: all alignments equally likely

- Model learns that, e.g., la is often aligned with the

# EM Algorithm

... la maison ... la maison blue ... la fleur ...

... the house ... the blue house ... the flower ...

- After one iteration

- Alignments, e.g., between la and the are more likely

# EM Algorithm



```
... la maison ... la maison bleu ... la fleur ...



... the house ... the blue house ... the flower ...
```

- After another iteration

- It becomes apparent that alignments, e.g., between fleur and flower are more likely (pigeon hole principle)

# EM Algorithm

... la maison ... la maison bleu ... la fleur ...

/ / | | | X | |

... the house ... the blue house ... the flower ...

- Convergence

- Inherent hidden structure revealed by EM

# EM Algorithm

... la maison ... la maison bleu ... la fleur ...

... the house ... the blue house ... the flower ...

$$p(\text{la}\,|\,\text{the}) = 0.453$$
$$p(\text{le}\,|\,\text{the}) = 0.334$$
$$p(\text{maison}\,|\,\text{house}) = 0.876$$
$$p(\text{bleu}\,|\,\text{blue}) = 0.563$$
$$\ldots$$

- Parameter estimation from the aligned corpus

# IBM Model 1 and EM

- EM Algorithm consists of two steps

- Expectation-Step: Apply model to the data

    - parts of the model are hidden (here: alignments)
    - using the model, assign probabilities to possible values

- Maximization-Step: Estimate model from data

    - take assign values as fact
    - collect counts (weighted by probabilities)
    - estimate model from counts

- Iterate these steps until convergence

# IBM Model 1 and EM

- We need to be able to compute:

  - Expectation-Step: probability of alignments

  - Maximization-Step: count collection

# IBM Model 1 and EM

- **Probabilities**

$$p(\text{the}|\text{la}) = 0.7 \qquad p(\text{house}|\text{la}) = 0.05$$
$$p(\text{the}|\text{maison}) = 0.1 \qquad p(\text{house}|\text{maison}) = 0.8$$

- **Alignments**



$$p(\mathbf{e}, a|\mathbf{f}) = 0.56 \qquad p(\mathbf{e}, a|\mathbf{f}) = 0.035 \qquad p(\mathbf{e}, a|\mathbf{f}) = 0.08 \qquad p(\mathbf{e}, a|\mathbf{f}) = 0.005$$

$$p(a|\mathbf{e}, \mathbf{f}) = 0.824 \qquad p(a|\mathbf{e}, \mathbf{f}) = 0.052 \qquad p(a|\mathbf{e}, \mathbf{f}) = 0.118 \qquad p(a|\mathbf{e}, \mathbf{f}) = 0.007$$

- **Counts**

$$c(\text{the}|\text{la}) = 0.824 + 0.052 \qquad c(\text{house}|\text{la}) = 0.052 + 0.007$$
$$c(\text{the}|\text{maison}) = 0.118 + 0.007 \qquad c(\text{house}|\text{maison}) = 0.824 + 0.118$$

# Convergence

das   Haus        das   Buch       ein   Buch

the   house        the   book        a   book

| $e$ | $f$ | initial | 1st it. | 2nd it. | 3rd it. | ... | final |
|---|---|---|---|---|---|---|---|
| the | das | 0.25 | 0.5 | 0.6364 | 0.7479 | ... | 1 |
| book | das | 0.25 | 0.25 | 0.1818 | 0.1208 | ... | 0 |
| house | das | 0.25 | 0.25 | 0.1818 | 0.1313 | ... | 0 |
| the | buch | 0.25 | 0.25 | 0.1818 | 0.1208 | ... | 0 |
| book | buch | 0.25 | 0.5 | 0.6364 | 0.7479 | ... | 1 |
| a | buch | 0.25 | 0.25 | 0.1818 | 0.1313 | ... | 0 |
| book | ein | 0.25 | 0.5 | 0.4286 | 0.3466 | ... | 0 |
| a | ein | 0.25 | 0.5 | 0.5714 | 0.6534 | ... | 1 |
| the | haus | 0.25 | 0.5 | 0.4286 | 0.3466 | ... | 0 |
| house | haus | 0.25 | 0.5 | 0.5714 | 0.6534 | ... | 1 |

# Perplexity

- How well does the model fit the data?

- Perplexity: derived from probability of the training data according to the model

$$\log_2 PP = -\sum_s \log_2 p(\mathbf{e}_s|\mathbf{f}_s)$$

- Example ($\epsilon$=1)

|  | initial | 1st it. | 2nd it. | 3rd it. | ... | final |
|---|---|---|---|---|---|---|
| $p(\text{the haus}|\text{das haus})$ | 0.0625 | 0.1875 | 0.1905 | 0.1913 | ... | 0.1875 |
| $p(\text{the book}|\text{das buch})$ | 0.0625 | 0.1406 | 0.1790 | 0.2075 | ... | 0.25 |
| $p(\text{a book}|\text{ein buch})$ | 0.0625 | 0.1875 | 0.1907 | 0.1913 | ... | 0.1875 |
| perplexity | 4095 | 202.3 | 153.6 | 131.6 | ... | 113.8 |

# Problems with Model 1

- There's a reason they designed models 2-5!

- Problems: alignments jump around, align everything to rare words

- Experimental setup:
  - Training data: 1.1M sentences of French-English text, Canadian Hansards
  - Evaluation metric: alignment error Rate (AER)
  - Evaluation data: 447 hand-aligned sentences

# IBM Model 2: Global Monotonicity

# Monotonic Translation

Japan shaken by two new quakes

Le Japon secoué par deux nouveaux séismes

# Local Order Change

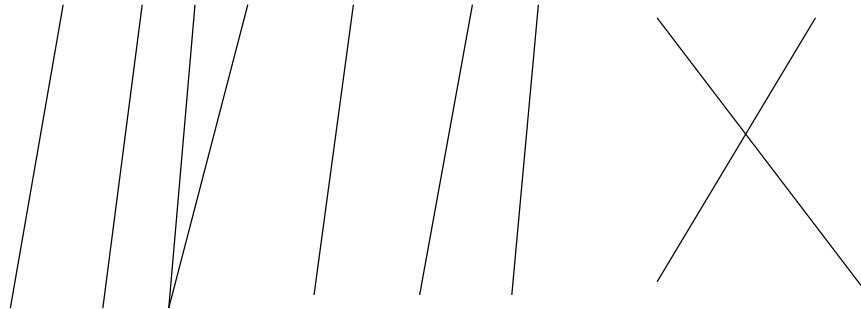Japan is at the junction of four tectonic plates

Le Japon est au confluent de quatre plaques tectoniques

# IBM Model 2

- Alignments tend to the diagonal (broadly at least)

$$P(f, a | e) = \prod_j P(a_j = i | j, I, J) P(f_j | e_i)$$

$$P(dist = i - j\frac{I}{J})$$

$$\frac{1}{Z} e^{-\alpha(i - j\frac{I}{J})}$$

# EM for Models 1/2

- **Model 1 Parameters:**
  Translation probabilities (1+2)   $P(f_j|e_i)$
  Distortion parameters (2 only)   $P(a_j = i|j, I, J)$

- **Start with** $P(f_j|e_i)$ uniform, including $P(f_j|null)$
- **For each sentence:**
  - For each French position j
    - Calculate posterior over English positions

$$P(a_j = i|f, e) = \frac{P(a_j = i|j, I, J)P(f_j|e_i)}{\sum_{i'} P(a_j = i'|j, I, J)P(f_j|e'_i)}$$

    - (or just use best single alignment)
    - Increment count of word $f_j$ with word $e_i$ by these amounts
    - Also re-estimate distortion probabilities for model 2
- **Iterate until convergence**

# HMM Model: Local Monotonicity

# Phrase Movement

On Tuesday Nov. 4, earthquakes rocked Japan once again

Des tremblements de terre ont à nouveau touché le Japon jeudi 4 novembre.

# IBM Models 1/2



**E:** Thank you , I shall do so gladly .

(positions: 1 2 3 4 5 6 7 8 9)

**A:** ①  ③  ⑦  ⑥  ⑧  ⑧  ⑧  ⑧  ⑨

**F:** Gracias , lo haré de muy buen grado .

**Model Parameters**

*Translation:* P( $F_1$ = Gracias | $E_{A_1}$ = Thank )   *Alignment:* P( $A_2$ = 3)

# The HMM Model

**E:**   Thank you   ,   I   shall   do   so   gladly   .

**A:**   ○ → ①  → ③ → ⑦ → ⑥ → ⑧ → ⑧ → ⑧ → ⑧ → ⑨ → ○

**F:**   Gracias   ,   lo   haré   de   muy   buen grado   .

---

**Model Parameters**

*Emissions:* P( F$_1$ = Gracias | E$_{A_1}$ = Thank )      *Transitions:* P( A$_2$ = 3 | A$_1$ = 1)

# The HMM Model

- Model 2 preferred global monotonicity
- We want local monotonicity:
  - Most jumps are small
- HMM model (Vogel 96)

| f | $t(f \mid e)$ |
|---|---|
| nationale | 0.469 |
| national | 0.418 |
| nationaux | 0.054 |
| nationales | 0.029 |

$$P(f, a|e) = \prod_j P(a_j|a_{j-1})P(f_j|e_i)$$

$$P(a_j - a_{j-1})$$

-2 -1 0 1 2 3

  - Re-estimate using the forward-backward algorithm
  - Handling nulls requires some care
- What are we still missing?

# Models 3+: Fertility

# IBM Models 3/4/5

Mary did not slap the green witch

Mary not slap slap slap the green witch

$n(3|slap)$

Mary not slap slap slap NULL the green witch

$P(NULL)$

Mary no daba una botefada a la verde bruja

$t(la|the)$

$d(j|i)$

Mary no daba una botefada a la bruja verde

[from Al-Onaizan and Knight, 1998]

*the*

| f | $t(f \mid e)$ | $\phi$ | $n(\phi \mid e)$ |
|------|------|------|------|
| le | 0.497 | 1 | 0.746 |
| la | 0.207 | 0 | 0.254 |
| les | 0.155 | | |
| l' | 0.086 | | |
| ce | 0.018 | | |
| cette | 0.011 | | |

*not*

| f | $t(f \mid e)$ | $\phi$ | $n(\phi \mid e)$ |
|------|------|------|------|
| ne | 0.497 | 2 | 0.735 |
| pas | 0.442 | 0 | 0.154 |
| non | 0.029 | 1 | 0.107 |
| rien | 0.011 | | |

*farmers*

| f | $t(f \mid e)$ | $\phi$ | $n(\phi \mid e)$ |
|------|------|------|------|
| agriculteurs | 0.442 | 2 | 0.731 |
| les | 0.418 | 1 | 0.228 |
| cultivateurs | 0.046 | 0 | 0.039 |
| producteurs | 0.021 | | |

# Example: Idioms

*nodding*

| f | $t(f \mid e)$ | $\phi$ | $n(\phi \mid e)$ |
|---|---|---|---|
| signe | 0.164 | 4 | 0.342 |
| la | 0.123 | 3 | 0.293 |
| tête | 0.097 | 2 | 0.167 |
| oui | 0.086 | 1 | 0.163 |
| fait | 0.073 | 0 | 0.023 |
| que | 0.073 | | |
| hoche | 0.054 | | |
| hocher | 0.048 | | |
| faire | 0.030 | | |
| me | 0.024 | | |
| approuve | 0.019 | | |
| qui | 0.019 | | |
| un | 0.012 | | |
| faites | 0.011 | | |

he is nodding

il hoche la tête

# Example: Morphology

*should*

| f | $t(f \mid e)$ | $\phi$ | $n(\phi \mid e)$ |
|---|---|---|---|
| devrait | 0.330 | 1 | 0.649 |
| devraient | 0.123 | 0 | 0.336 |
| devrions | 0.109 | 2 | 0.014 |
| faudrait | 0.073 | | |
| faut | 0.058 | | |
| doit | 0.058 | | |
| aurait | 0.041 | | |
| doivent | 0.024 | | |
| devons | 0.017 | | |
| devrais | 0.013 | | |

# Machine Translation

Dan Klein
UC Berkeley

Many slides from John DeNero and Philip Koehn

# Phrase-Based Model

| natuerlich | hat | john | spass am | spiel |
|---|---|---|---|---|

| of course | john | has | fun with the | game |
|---|---|---|---|---|

- Foreign input is segmented in phrases
- Each phrase is translated into English
- Phrases are reordered

# Getting Phrases

# Word Alignment with IBM Models

- IBM Models create a **many-to-one** mapping

  - words are aligned using an alignment function
  - a function may return the same value for different input
    (one-to-many mapping)
  - a function can not return multiple values for one input
    (no many-to-one mapping)

- Real word alignments have **many-to-many** mappings

# Symmetrization

- Run IBM Model training in both directions

$\rightarrow$ two sets of word alignment points

- Intersection: high precision alignment points

- Union: high recall alignment points

- Refinement methods explore the sets between intersection and union

# Example

## english to spanish

|  | Maria | no | daba | una | bofetada | a | la | bruja | verde |
|---|---|---|---|---|---|---|---|---|---|
| Mary | ■ |  |  |  |  |  |  |  |  |
| did |  |  |  |  |  | ■ |  |  |  |
| not |  | ■ |  |  |  |  |  |  |  |
| slap |  |  | ■ | ■ | ■ |  |  |  |  |
| the |  |  |  |  |  |  | ■ |  |  |
| green |  |  |  |  |  |  |  |  | ■ |
| witch |  |  |  |  |  |  |  | ■ |  |

## spanish to english

|  | Maria | no | daba | una | bofetada | a | la | bruja | verde |
|---|---|---|---|---|---|---|---|---|---|
| Mary | ■ |  |  |  |  |  |  |  |  |
| did |  | ■ |  |  |  |  |  |  |  |
| not |  | ■ |  |  |  |  |  |  |  |
| slap |  |  |  | ■ |  |  |  |  |  |
| the |  |  |  |  |  |  | ■ |  |  |
| green |  |  |  |  |  |  |  |  | ■ |
| witch |  |  |  |  |  |  |  | ■ |  |

## intersection

|  | Maria | no | daba | una | bofetada | a | la | bruja | verde |
|---|---|---|---|---|---|---|---|---|---|
| Mary | ■ |  |  |  |  |  |  |  |  |
| did |  |  |  |  |  |  |  |  |  |
| not |  | ■ |  |  |  |  |  |  |  |
| slap |  |  |  | ■ |  |  |  |  |  |
| the |  |  |  |  |  |  | ■ |  |  |
| green |  |  |  |  |  |  |  |  | ■ |
| witch |  |  |  |  |  |  |  | ■ |  |

# Growing Heuristics



**black**: intersection      **grey**: additional points in union

- Add alignment points from union based on heuristics:
  - directly/diagonally neighboring points
  - finally, add alignments that connect unaligned words in source and/or target
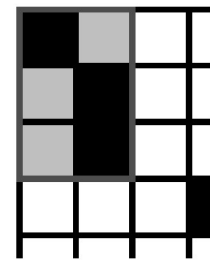- Popular method: grow-diag-final-and

# Extracting Phrase Pairs



extract phrase pair consistent with word alignment:

assumes that / geht davon aus , dass

# Consistent



consistent     inconsistent     consistent
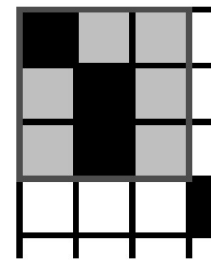
**ok**        **violated**        **ok**

one       unaligned

alignment      word is fine

point outside

All words of the phrase pair have to align to each other.

# Phrase Pair Extraction



## Smallest phrase pairs:

michael — michael

assumes — geht davon aus / geht davon aus ,

that — dass / , dass

he — er

will stay — bleibt

in the — im

house — haus

unaligned words (here: German comma) lead to multiple translations

# Larger Phrase Pairs

|  | michael | geht | davon | aus | , | dass | er | im | haus | bleibt |
|---|---|---|---|---|---|---|---|---|---|---|
| michael | ■ | | | | | | | | | |
| assumes | | ■ | ■ | ■ | | | | | | |
| that | | | | | | ■ | | | | |
| he | | | | | | | ■ | | | |
| will | | | | | | | | | | ■ |
| stay | | | | | | | | | | ■ |
| in | | | | | | | | ■ | | |
| the | | | | | | | | ■ | | |
| house | | | | | | | | | ■ | |

michael assumes — michael geht davon aus / michael geht davon aus ,

assumes that — geht davon aus , dass   ;   assumes that he — geht davon aus , dass er

that he — dass er / , dass er   ;   in the house — im haus

michael assumes that — michael geht davon aus , dass

michael assumes that he — michael geht davon aus , dass er

michael assumes that he will stay in the house  — michael geht davon aus , dass er im haus bleibt

assumes that he will stay in the house — geht davon aus , dass er im haus bleibt

that he will stay in the house — dass er im haus bleibt   ;   dass er im haus bleibt ,

he will stay in the house — er im haus bleibt   ;   will stay in the house — im haus bleibt

# Phrase Translation Table

- Main knowledge source: table with phrase translations and their probabilities

- Example: phrase translations for natuerlich

| Translation | Probability $\phi(\bar{e}|f)$ |
|:---:|:---:|
| of course | 0.5 |
| naturally | 0.3 |
| of course , | 0.15 |
| , of course , | 0.05 |

# Scoring Phrase Translations

- Phrase pair extraction: collect all phrase pairs from the data

- Phrase pair scoring: assign probabilities to phrase translations

- Score by relative frequency:

$$\phi(\bar{f}|\bar{e}) = \frac{\text{count}(\bar{e}, \bar{f})}{\sum_{\bar{f}_i} \text{count}(\bar{e}, \bar{f}_i)}$$

# Real Example

- Phrase translations for den Vorschlag learned from the Europarl corpus:

| English | $\phi(\bar{e}|f)$ | English | $\phi(\bar{e}|f)$ |
|---|---|---|---|
| the proposal | 0.6227 | the suggestions | 0.0114 |
| 's proposal | 0.1068 | the proposed | 0.0114 |
| a proposal | 0.0341 | the motion | 0.0091 |
| the idea | 0.0250 | the idea of | 0.0091 |
| this proposal | 0.0227 | the proposal , | 0.0068 |
| proposal | 0.0205 | its proposal | 0.0068 |
| of the proposal | 0.0159 | it | 0.0068 |
| the proposals | 0.0159 | ... | ... |

  - lexical variation (proposal vs suggestions)
  - morphological variation (proposal vs proposals)
  - included function words (the, a, ...)
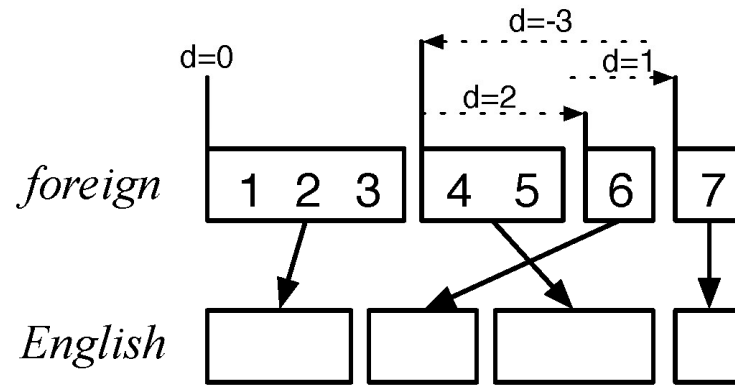  - noise (it)

# Other Scoring Terms

# More Feature Functions

- Bidirectional alignment probabilities: $\phi(\bar{e}|\bar{f})$ and $\phi(\bar{f}|\bar{e})$

- Rare phrase pairs have unreliable phrase translation probability estimates
  $\rightarrow$ lexical weighting with word translation probabilities



$$\text{lex}(\bar{e}|\bar{f}, a) = \prod_{i=1}^{\text{length}(\bar{e})} \frac{1}{|\{j|(i,j) \in a\}|} \sum_{\forall (i,j) \in a} w(e_i|f_j)$$

# Distance-Based Reordering



| phrase | translates | movement | distance |
|:---:|:---:|:---:|:---:|
| 1 | 1–3 | start at beginning | 0 |
| 2 | 6 | skip over 4–5 | +2 |
| 3 | 4–5 | move back over 4–6 | -3 |
| 4 | 7 | skip over 6 | +1 |

Scoring function: $d(x) = \alpha^{|x|}$ — exponential with distance

# Phrase-Based Decoding

# Translation Options

| er | geht | ja | nicht | nach | hause |
|---|---|---|---|---|---|
| he | is | yes | not | after | house |
| it | are | is | do not | to | home |
| , it | goes | , of course | does not | according to | chamber |
| , he | go | , | is not | in | at home |

| | | | |
|---|---|---|---|
| it is | not | | home |
| he will be | is not | | under house |
| it goes | does not | | return home |
| he goes | do not | | do not |

| | |
|---|---|
| is | to |
| are | following |
| is after all | not after |
| does | not to |

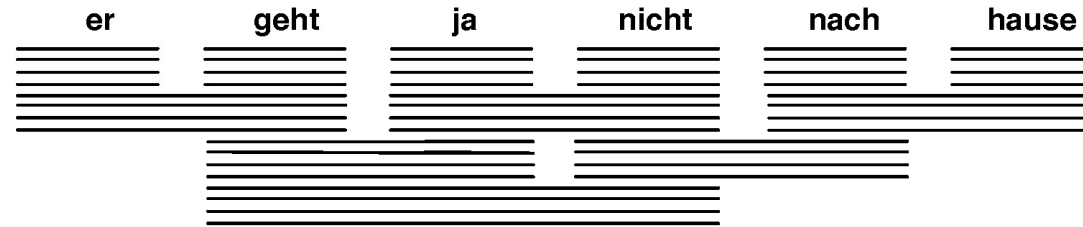| |
|---|
| not |
| is not |
| are not |
| is not a |

- Many translation options to choose from

  - in Europarl phrase table: 2727 matching phrase pairs for this sentence
  - by pruning to the top 20 per phrase, 202 translation options remain
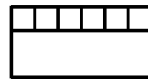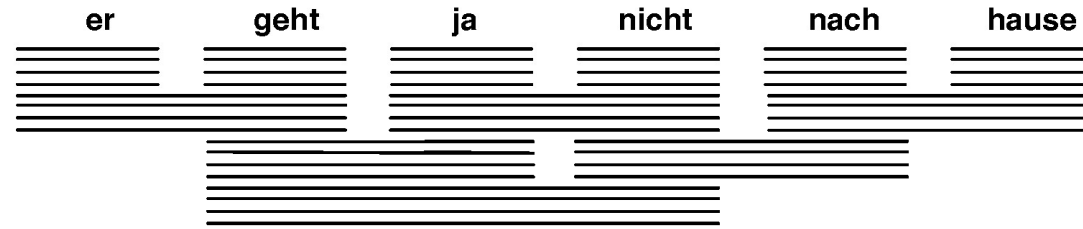
# Translation Options

| er | geht | ja | nicht | nach | hause |
|---|---|---|---|---|---|

| | | | | | |
|---|---|---|---|---|---|
| he | is | yes | not | after | house |
| it | are | is | do not | to | home |
| , it | goes | , of course | does not | according to | chamber |
| , he | go | | is not | in | at home |

| it is | not | home |
|---|---|---|
| he will be | is not | under house |
| it goes | does not | return home |
| he goes | do not | do not |

| is | to |
|---|---|
| are | following |
| is after all | not after |
| does | not to |

| not |
|---|
| is not |
| are not |
| is not a |

- The machine translation decoder does not know the right answer
  - picking the right translation options
  - arranging them in the right order

→ Search problem solved by heuristic beam search

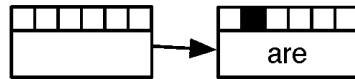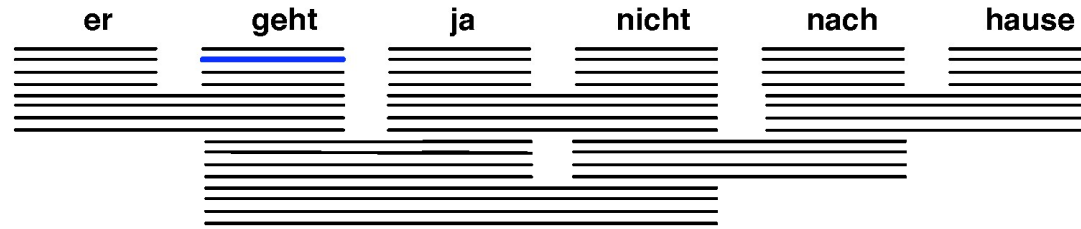# Decoding: Precompute Translation Options

er        geht        ja        nicht        nach        hause

consult phrase translation table for all input phrases

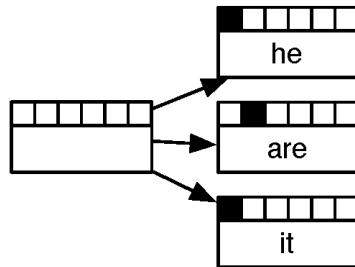# Decoding: Start with Initial Hypothesis

er　geht　ja　nicht　nach　hause

initial hypothesis: no input words covered, no output produced

# Decoding: Hypothesis Expansion

er        geht        ja        nicht        nach        hause

are

pick any translation option, create new hypothesis

# Decoding: Hypothesis Expansion

er      geht      ja      nicht      nach      hause

create hypotheses for all other translation options

# Decoding: Hypothesis Expansion

er       geht       ja       nicht       nach       hause

| | | | |
|---|---|---|---|
| | yes | | |
| he | goes | home | |
| are | does not | go | home |
| it | | to | |

also create hypotheses from created partial hypothesis

# Decoding: Find Best Path



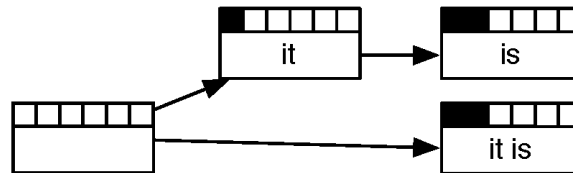backtrack from highest scoring complete hypothesis

# Dynamic Programming

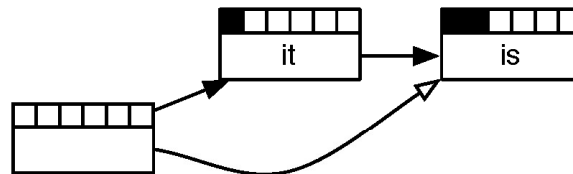# Computational Complexity

- The suggested process creates exponential number of hypothesis

- Machine translation decoding is NP-complete

- Reduction of search space:
  - recombination (risk-free)
  - pruning (risky)

# Recombination

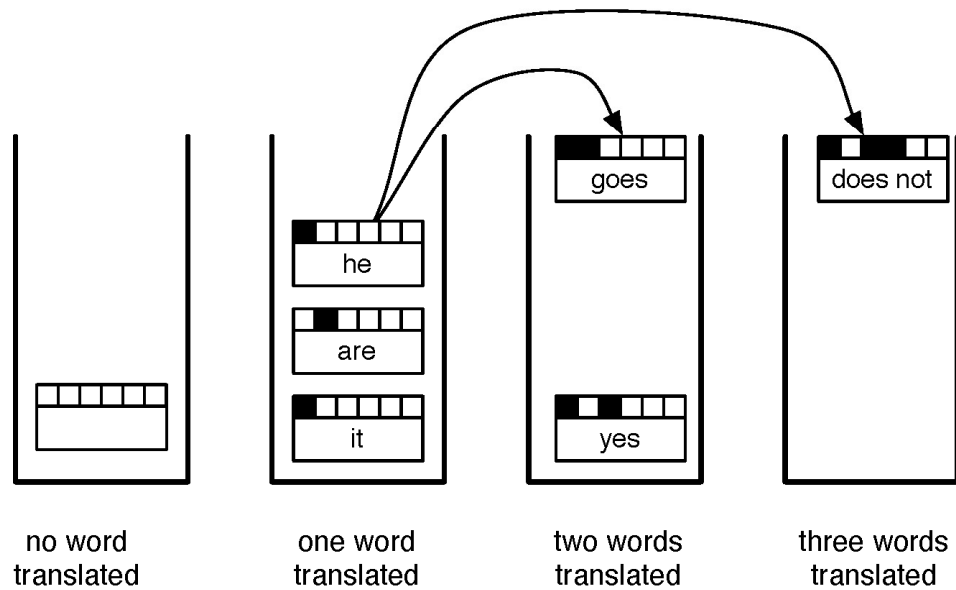- Two hypothesis paths lead to two matching hypotheses

    - same foreign words translated
    - same English words in the output



- Worse hypothesis is dropped

# Stacks



|no word<br>translated|one word<br>translated|two words<br>translated|three words<br>translated|

- Hypothesis expansion in a stack decoder
    - translation option is applied to hypothesis
    - new hypothesis is dropped into a stack further down
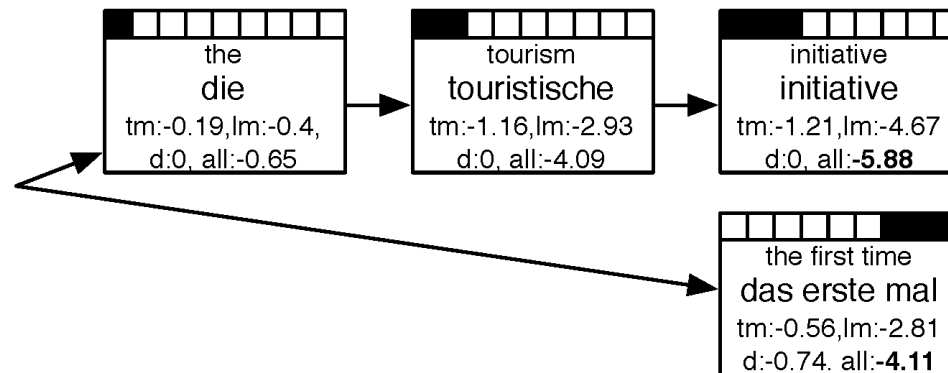
# Stack Decoding Algorithm

1: place empty hypothesis into stack 0
2: **for all** stacks $0...n-1$ **do**
3:     **for all** hypotheses in stack **do**
4:         **for all** translation options **do**
5:             **if** applicable **then**
6:                 create new hypothesis
7:                 place in stack
8:                 recombine with existing hypothesis **if** possible
9:                 prune stack **if** too big
10:             **end if**
11:         **end for**
12:     **end for**
13: **end for**

# Future Costs

# Translating the Easy Part First?

**the tourism initiative addresses this for the first time**



| | | |
|---|---|---|
| the | tourism | initiative |
| **die** | **touristische** | **initiative** |
| tm:-0.19,lm:-0.4, | tm:-1.16,lm:-2.93 | tm:-1.21,lm:-4.67 |
| d:0, all:-0.65 | d:0, all:-4.09 | d:0, all:-**5.88** |

the first time
**das erste mal**
tm:-0.56,lm:-2.81
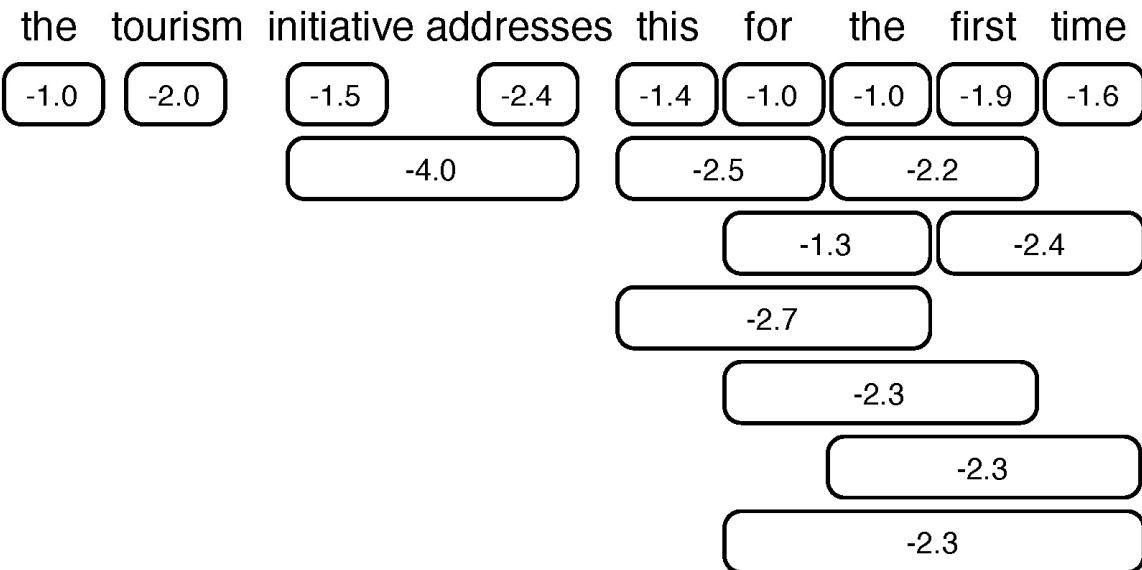d:-0.74. all:-**4.11**

both hypotheses translate 3 words
worse hypothesis has better score

# Estimating Future Cost

- Future cost estimate: how expensive is translation of rest of sentence?

- Optimistic: choose cheapest translation options

- Cost for each translation option
    - **translation model**: cost known
    - **language model:** output words known, but not context
      $\rightarrow$ estimate without context
    - **reordering model:** unknown, ignored for future cost estimation

# Cost Estimates from Translation Options

| the | tourism | initiative | addresses | this | for | the | first | time |
|---|---|---|---|---|---|---|---|---|

| -1.0 | -2.0 | -1.5 | -2.4 | -1.4 | -1.0 | -1.0 | -1.9 | -1.6 |

-4.0  -2.5  -2.2

-1.3  -2.4

-2.7

-2.3

-2.3

-2.3

cost of cheapest translation options for each input span (log-probabilities)
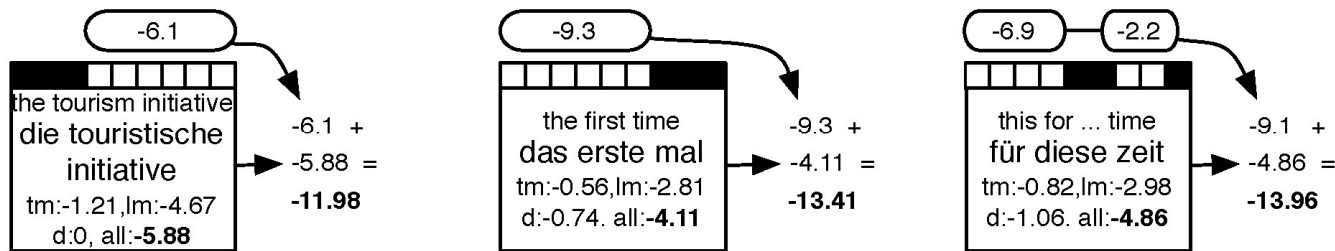
# Cost Estimates for all Spans

- Compute cost estimate for all contiguous spans by combining cheapest options

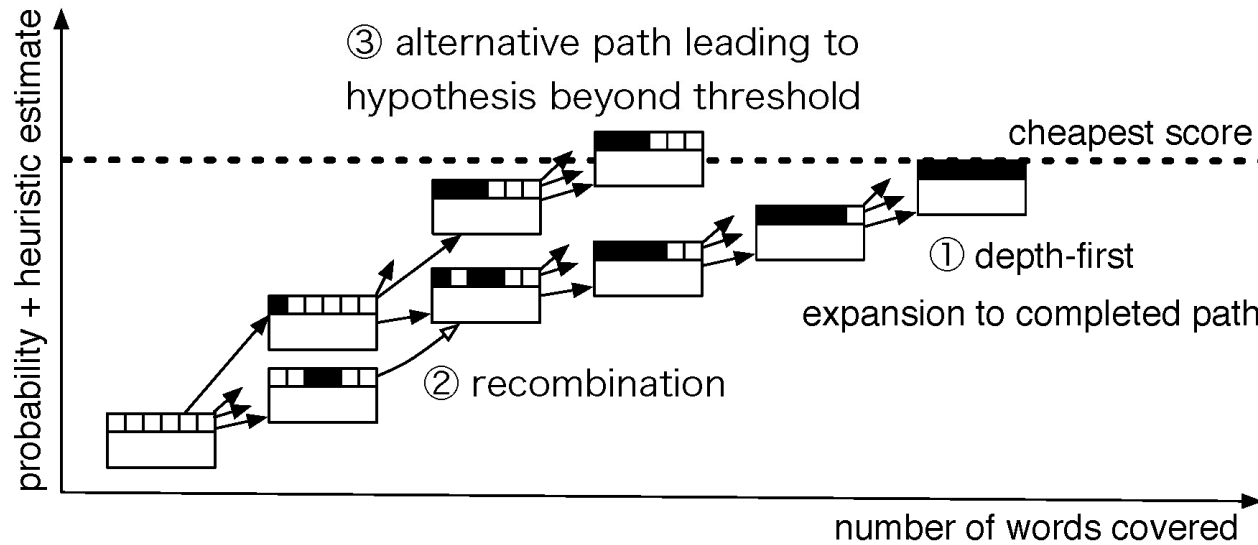| first word | future cost estimate for $n$ words (from first) | | | | | | | | |
|:---:|:---:|:---:|:---:|:---:|:---:|:---:|:---:|:---:|:---:|
| | 1 | 2 | 3 | 4 | 5 | 6 | 7 | 8 | 9 |
| the | -1.0 | -3.0 | -4.5 | -6.9 | -8.3 | -9.3 | -9.6 | -10.6 | -10.6 |
| tourism | -2.0 | -3.5 | -5.9 | -7.3 | -8.3 | -8.6 | -9.6 | -9.6 | |
| initiative | -1.5 | -3.9 | -5.3 | -6.3 | -6.6 | -7.6 | -7.6 | | |
| addresses | -2.4 | -3.8 | -4.8 | -5.1 | -6.1 | -6.1 | | | |
| this | -1.4 | -2.4 | -2.7 | -3.7 | -3.7 | | | | |
| for | -1.0 | -1.3 | -2.3 | -2.3 | | | | | |
| the | -1.0 | -2.2 | -2.3 | | | | | | |
| first | -1.9 | -2.4 | | | | | | | |
| time | -1.6 | | | | | | | | |

- Function words cheaper (the: -1.0) than content words (tourism -2.0)
- Common phrases cheaper (for the first time: -2.3)
  than unusual ones (tourism initiative addresses: -5.9)

# Combining Score and Future Cost



- Hypothesis score and future cost estimate are combined for pruning

  - left hypothesis starts with hard part: the tourism initiative
    score: -5.88, future cost: -6.1 $\rightarrow$ total cost -11.98

  - middle hypothesis starts with easiest part: the first time
    score: -4.11, future cost: -9.3 $\rightarrow$ total cost -13.41

  - right hypothesis picks easy parts: this for ... time
    score: -4.86, future cost: -9.1 $\rightarrow$ total cost -13.96

# A* Search



- Uses *admissible* future cost heuristic: never overestimates cost

- Translation agenda: create hypothesis with lowest score + heuristic cost

- Done, when complete hypothesis created