# Machine Translation
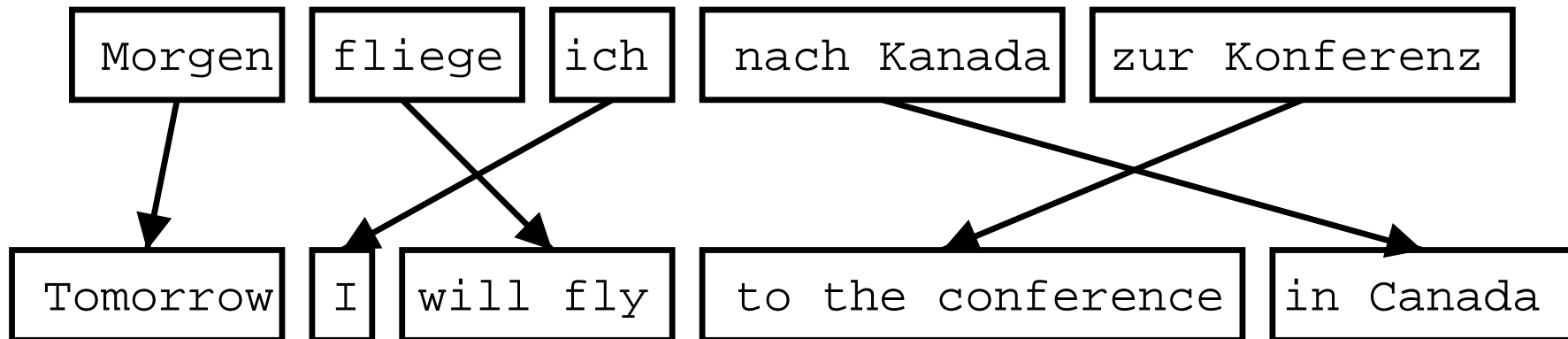# Phrase Models

Philipp Koehn, University of Edinburgh

12 February 2009

# Phrase-based translation

| Morgen | fliege | ich | nach Kanada | zur Konferenz |
|--------|--------|-----|-------------|---------------|

| Tomorrow | I | will fly | to the conference | in Canada |
|----------|---|----------|-------------------|-----------|

- Foreign input is segmented in phrases
  - any sequence of words, not necessarily linguistically motivated

- Each phrase is translated into English

- Phrases are reordered

School of
**informatics**

# Phrase-based translation model

- Major components of phrase-based model
  - **phrase translation model** $\phi(\mathbf{f}|\mathbf{e})$
  - **reordering model** $\Omega(\mathbf{f}|\mathbf{e})$
  - **language model** $p_{\mathrm{LM}}(\mathbf{e})$

- Bayes rule

$$
\begin{aligned}
\mathrm{argmax}_{\mathbf{e}}p(\mathbf{e}|\mathbf{f}) &= \mathrm{argmax}_{\mathbf{e}}p(\mathbf{f}|\mathbf{e})p(\mathbf{e}) \\
&= \mathrm{argmax}_{\mathbf{e}}\phi(\mathbf{f}|\mathbf{e})\ p_{\mathrm{LM}}(\mathbf{e})\ \Omega(\mathbf{f}|\mathbf{e})
\end{aligned}
$$

- Sentence $\mathbf{f}$ is decomposed into $I$ phrases $\bar{f}_1^I = \bar{f}_1, ..., \bar{f}_I$

- Decomposition of $\phi(\mathbf{f}|\mathbf{e})$

$$
\phi(\bar{f}_1^I|\bar{e}_1^I) = \prod_{i=1}^{I} \phi(\bar{f}_i|\bar{e}_i)\ \omega^{d(\mathsf{start}_i - \mathsf{end}_{i-1} - 1)})
$$

School of **informatics**

# Advantages of phrase-based translation

- *Many-to-many* translation can handle non-compositional phrases

- Use of *local context* in translation

- The more data, the *longer phrases* can be learned

School of
**informatics**

# Phrase translation table

- Phrase translations for *den Vorschlag*

| English | $\phi(e|f)$ | English | $\phi(e|f)$ |
|---|---|---|---|
| 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 | ... | ... |

School of **informatics**

# How to learn the phrase translation table?

- Start with the *word alignment*:



- Collect all phrase pairs that are **consistent** with the word alignment

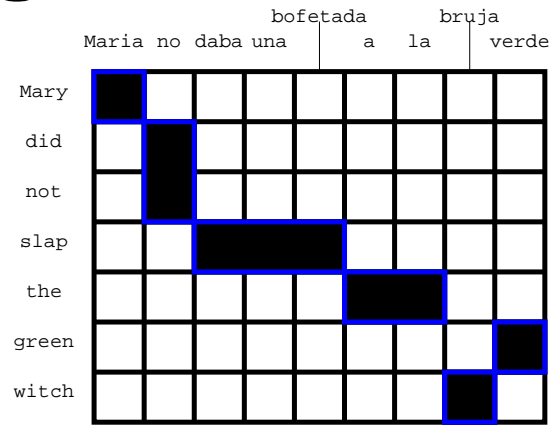# Consistent with word alignment



consistent     inconsistent     inconsistent

- **Consistent with the word alignment** :=

  phrase alignment has to *contain all alignment points* for all covered words
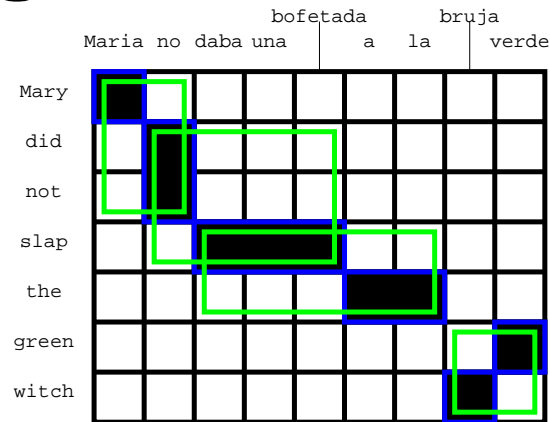
$$(\overline{e}, \overline{f}) \in BP \Leftrightarrow \qquad \forall e_i \in \overline{e} : (e_i, f_j) \in A \rightarrow f_j \in \overline{f}$$

$$\text{AND} \quad \forall f_j \in \overline{f} : (e_i, f_j) \in A \rightarrow e_i \in \overline{e}$$

# Word alignment induced phrases



**(Maria, Mary), (no, did not), (slap, daba una bofetada), (a la, the), (bruja, witch), (verde, green)**
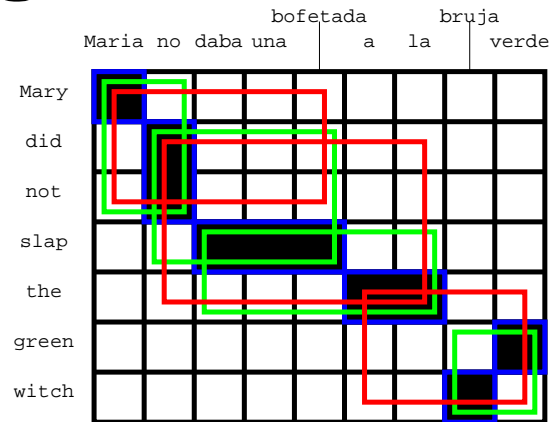
# Word alignment induced phrases



(Maria, Mary), (no, did not), (slap, daba una bofetada), (a la, the), (bruja, witch), (verde, green),

(Maria no, Mary did not), (no daba una bofetada, did not slap), (daba una bofetada a la, slap the),

(bruja verde, green witch)
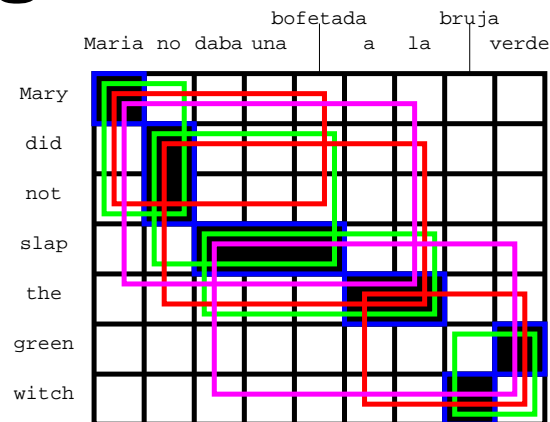
# Word alignment induced phrases



(Maria, Mary), (no, did not), (slap, daba una bofetada), (a la, the), (bruja, witch), (verde, green),

(Maria no, Mary did not), (no daba una bofetada, did not slap), (daba una bofetada a la, slap the),

(bruja verde, green witch),  (Maria no daba una bofetada, Mary did not slap),

(no daba una bofetada a la, did not slap the), (a la bruja verde, the green witch)

School of **informatics**

# Word alignment induced phrases



**(Maria, Mary), (no, did not), (slap, daba una bofetada), (a la, the), (bruja, witch), (verde, green),**

**(Maria no, Mary did not), (no daba una bofetada, did not slap), (daba una bofetada a la, slap the),**

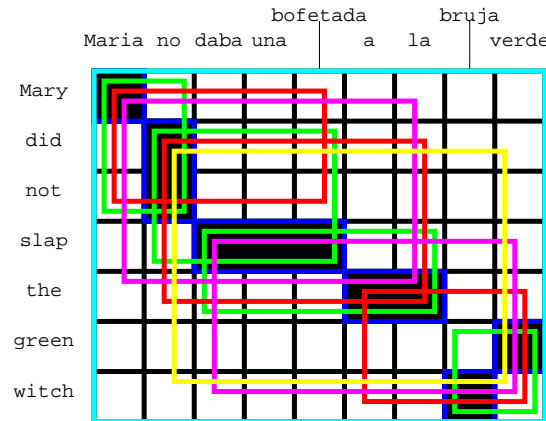**(bruja verde, green witch), (Maria no daba una bofetada, Mary did not slap),**

**(no daba una bofetada a la, did not slap the), (a la bruja verde, the green witch),**

**(Maria no daba una bofetada a la, Mary did not slap the),**

**(daba una bofetada a la bruja verde, slap the green witch)**

School of **informatics**

# Word alignment induced phrases (5)



(Maria, Mary), (no, did not), (slap, daba una bofetada), (a la, the), (bruja, witch), (verde, green),

(Maria no, Mary did not), (no daba una bofetada, did not slap), (daba una bofetada a la, slap the),

(bruja verde, green witch), (Maria no daba una bofetada, Mary did not slap),

(no daba una bofetada a la, did not slap the), (a la bruja verde, the green witch),

(Maria no daba una bofetada a la, Mary did not slap the), (daba una bofetada a la bruja verde,

slap the green witch), (no daba una bofetada a la bruja verde, did not slap the green witch),

(Maria no daba una bofetada a la bruja verde, Mary did not slap the green witch)

School of
**informatics**

# Probability distribution of phrase pairs

- We need a **probability distribution** $\phi(\overline{f}|\overline{e})$ over the collected phrase pairs
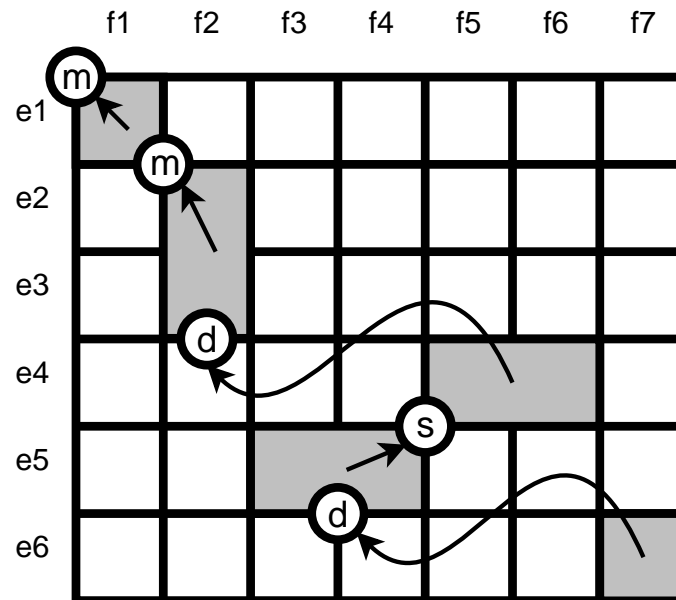
$\Rightarrow$ Possible *choices*

- *relative frequency* of collected phrases: $\phi(\overline{f}|\overline{e}) = \dfrac{\text{count}(\overline{f},\overline{e})}{\sum_{\overline{f}}\text{count}(\overline{f},\overline{e})}$

- or, conversely $\phi(\overline{e}|\overline{f})$
- use *lexical translation probabilities*

School of
**informatics**

# **Reordering**

- *Monotone* translation

  - do not allow any reordering
  → worse translations

- *Limiting* reordering (to movement over max. number of words) helps

- *Distance-based* reordering cost

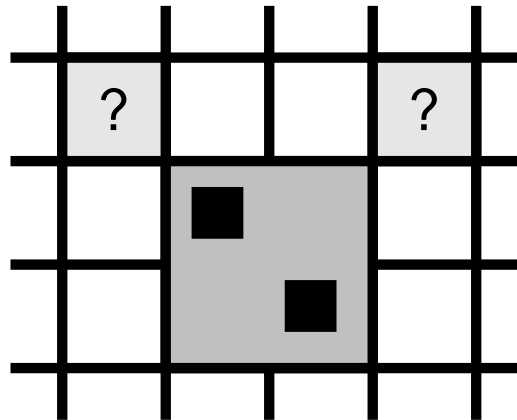  - moving a foreign phrase over $n$ words: cost $\omega^n$

- *Lexicalized* reordering model

School of **informatics**

# Lexicalized reordering models

- Three **orientation** types: **monotone**, **swap**, **discontinuous**

- Probability $p(swap|e, f)$ depends on foreign (and English) *phrase* involved

School of **informatics**

# Learning lexicalized reordering models



[from Koehn et al., 2005, IWSLT]

• Orientation type is *learned during phrase extractions*

• *Alignment point* to the *top left* (monotone) or *top right* (swap)?

• For more, see [Tillmann, 2003] or [Koehn et al., 2005]

# Names and Numbers

- All word tokens are treated the same

- Names and numbers pose special problems

  - there are many different names and numbers
  - if input and output use different scripts, translation is not easy

- Name translation is hard

  - names may not have a properly defined spelling in non-native scripts
  - training data is not always easy to come by
  - treated as special **transliteration** problem

School of
**informatics**

# XML Markup

`Er erzielte <NUMBER english='17.55'>17,55</NUMBER> Punkte .`

- *Add additional translation options*

  – number translation
  – name translation

- Additional options

  – provide multiple translations
  – provide probability distribution along with translations
  – allow bypassing of provided translations