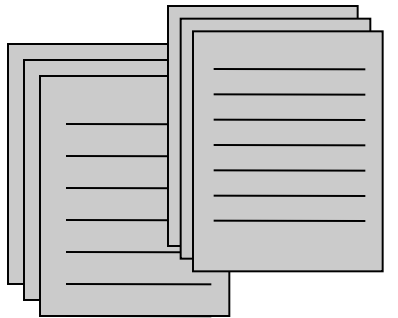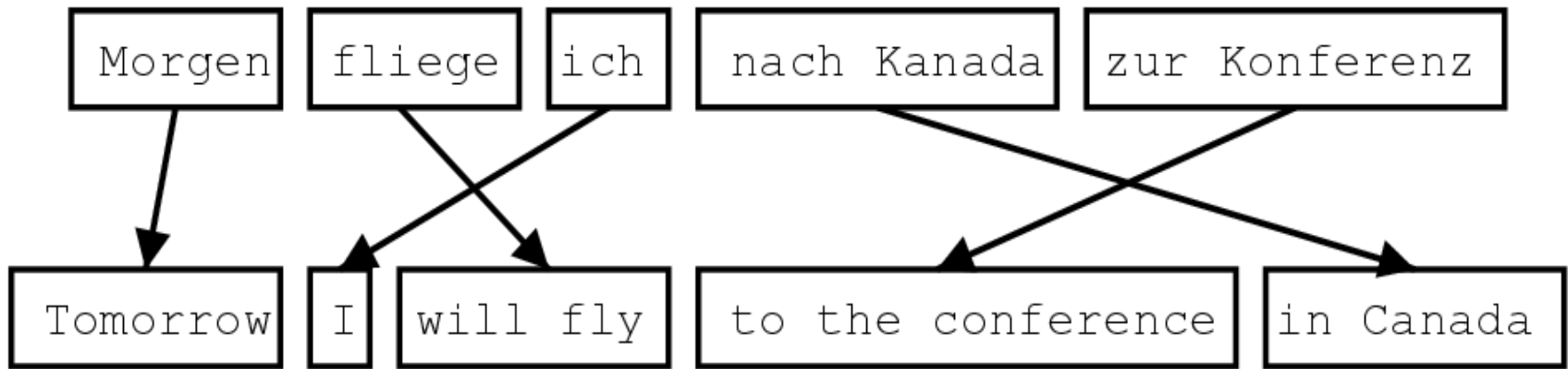# CSE 517
# Natural Language Processing
# Winter 2015

## Phrase Based Translation
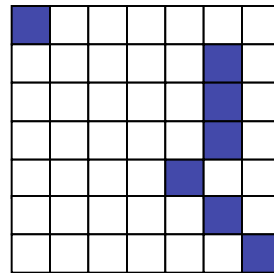
## Yejin Choi

Slides from Philipp Koehn, Dan Klein, Luke Zettlemoyer

# Phrase-Based Systems

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

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

Sentence-aligned corpus

Word alignments

cat ||| chat ||| 0.9
the cat ||| le chat ||| 0.8
dog ||| chien ||| 0.8
house ||| maison ||| 0.6
my house ||| ma maison ||| 0.9
language ||| langue ||| 0.9
…

Phrase table
(translation model)

# Phrase Translation Tables

- Defines the space of possible translations
  - each entry has an associated "probability"
- One learned example, for "den Vorschlag" from Europarl data

  - This table is noisy, has errors, and the entries do not necessarily match our linguistic intuitions about consistency….

# Phrase Translation Model

- Bayes rule

$$\mathbf{e}_{\text{best}} = \text{argmax}_{\mathbf{e}} \, p(\mathbf{e}|\mathbf{f})$$

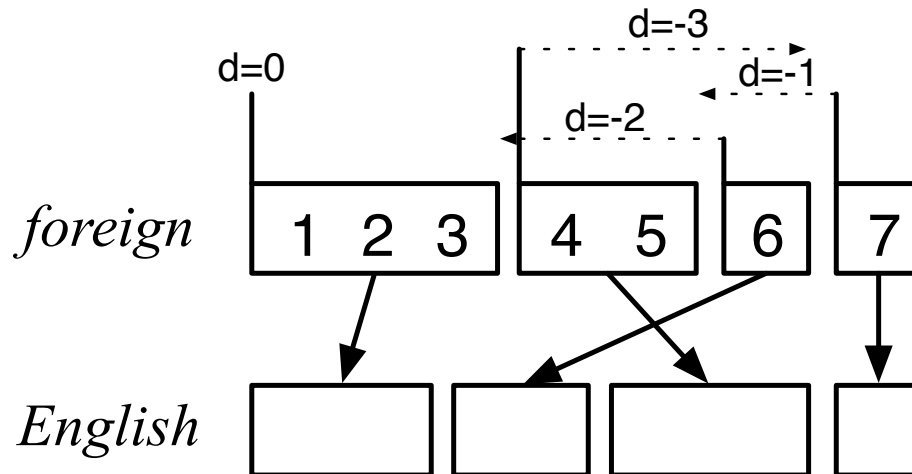$$= \text{argmax}_{\mathbf{e}} \, p(\mathbf{f}|\mathbf{e}) \, p_{\text{LM}}(\mathbf{e})$$

  – translation model $p(\mathbf{e}|\mathbf{f})$
  – language model $p_{\text{LM}}(\mathbf{e})$

- Decomposition of the translation model

$$p(\bar{f}_1^I|\bar{e}_1^I) = \prod_{i=1}^{I} \phi(\bar{f}_i|\bar{e}_i) \, d(start_i - end_{i-1} - 1)$$

  – phrase translation probability $\phi$
  – reordering probability $d$

# Distortion Model



| 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

# Extracting Phrases

- We will use word alignments to find phrases

- Question: what is the best set of phrases?

# Extracting Phrases

- Phrase alignment must
  - Contain at least one alignment edge
  - Contain all alignments for phrase pair

- Extract all such phrase pairs!

# Phrase Pair Extraction Example

(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)

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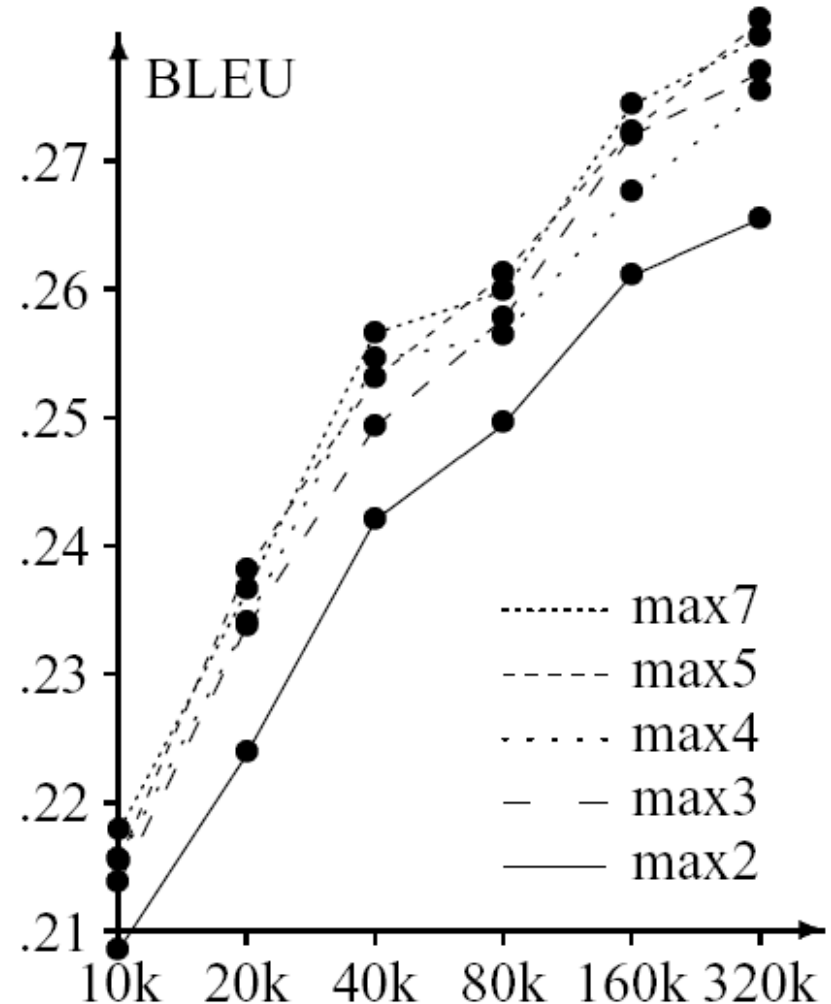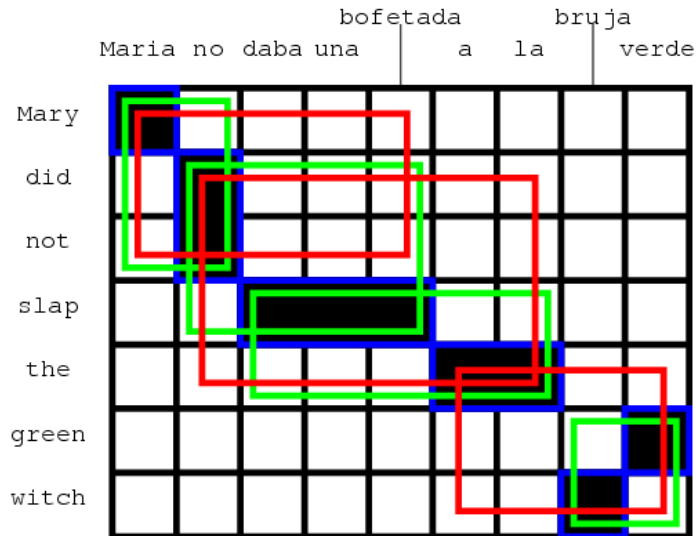
# Linguistic Phrases?

- Model is not limited to linguistic phrases
  (noun phrases, verb phrases, prepositional phrases, ...)

- Example non-linguistic phrase pair

$$\text{spass am} \rightarrow \text{fun with the}$$

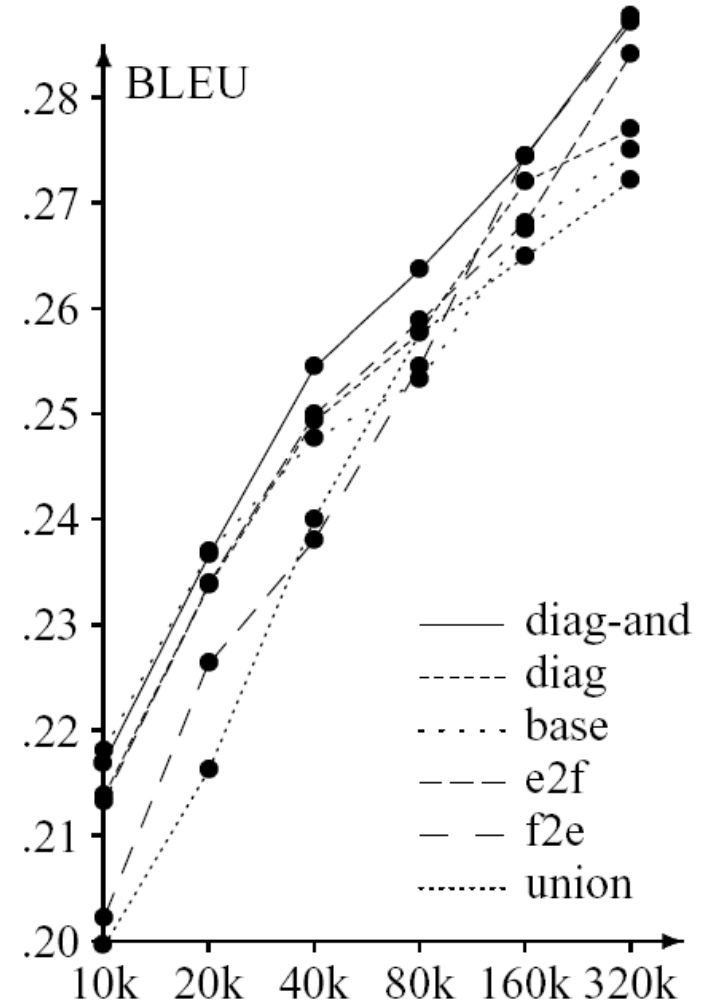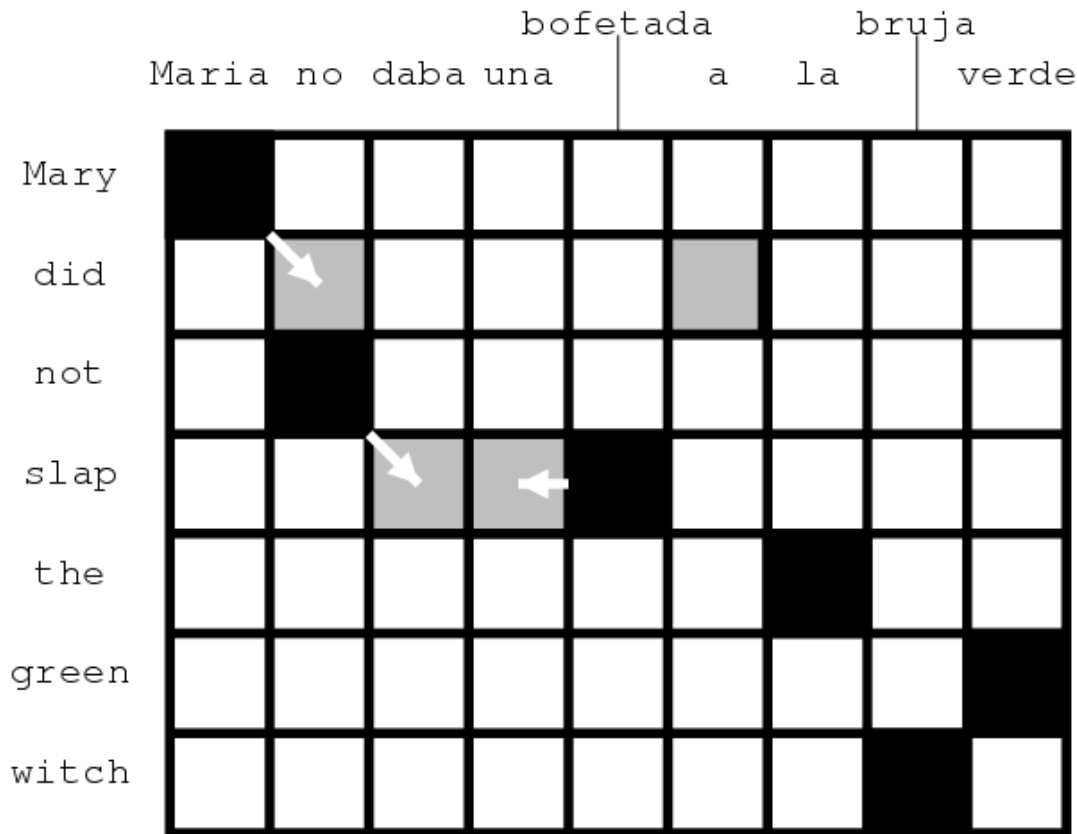- Prior noun often helps with translation of preposition

- Experiments show that limitation to linguistic phrases hurts quality

# Phrase Size

- ## Phrases do help
  - ### But they don't need to be long
  - ### Why should this be?

# Bidirectional Alignment

english to spanish

spanish to english

intersection

# Alignment Heuristics

# Size of Phrase Translation Table

- Phrase translation table typically bigger than corpus

  ... even with limits on phrase lengths (e.g., max 7 words)

$\rightarrow$ Too big to store in memory?

- Solution for training

  – extract to disk, sort, construct for one source phrase at a time

- Solutions for decoding

  – on-disk data structures with index for quick look-ups
  – suffix arrays to create phrase pairs on demand

# Why not Learn Phrases w/ EM?

## EM Training of the Phrase Model

- We presented a heuristic set-up to build phrase translation table
  (word alignment, phrase extraction, phrase scoring)

- Alternative: align phrase pairs directly with EM algorithm

  - initialization: uniform model, all $\phi(\bar{e}, \bar{f})$ are the same
  - expectation step:
    * estimate likelihood of all possible phrase alignments for all sentence pairs
  - maximization step:
    * collect counts for phrase pairs $(\bar{e}, \bar{f})$, weighted by alignment probability
    * update phrase translation probabilties $p(\bar{e}, \bar{f})$

- However: method easily overfits
  (learns very large phrase pairs, spanning entire sentences)

# Phrase Scoring

$$g(f, e) = \log \frac{c(e, f)}{c(e)}$$

$$g(\text{les chats}, \text{cats}) = \log \frac{c(\text{cats}, \text{les chats})}{c(\text{cats})}$$

- Learning weights has been tried, several times:
  - [Marcu and Wong, 02]
  - [DeNero et al, 06]
  - … and others

- Seems not to work well, for a variety of partially understood reasons

- Main issue: big chunks get all the weight, obvious priors don't help
  - Though, [DeNero et al 08]

# Translation: Codebreaking?

*"Also knowing nothing official about, but having guessed and inferred considerable about, the powerful new mechanized methods in cryptography—methods which I believe succeed even when one does not know what language has been coded—one naturally wonders if the problem of translation could conceivably be treated as a problem in cryptography.*

*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 (1955:18, quoting a letter he wrote in 1947)

# Translation is hard!



zi   zhu    zhong  duan
自   助    终    端

self  help terminal device

help oneself terminating machine

(ATM, "self-service terminal")

2

Examples from Liang Huang

# Translation is hard!

Examples from Liang Huang

# Translation is hard!

Examples from Liang Huang

# Translation is hard!

Examples from Liang Huang

# Translation is hard!



Examples from Liang Huang

# or even…

4

Examples from Liang Huang

# Scoring:

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

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

- **Basic approach, sum up phrase translation scores and a language model**
  - Define y = $p_1 p_2 \ldots p_L$ to be a translation with phrase pairs $p_i$
  - Define e(y) be the output English sentence in y
  - Let h() be the log probability under a tri-gram language model
  - Let g() be a phrase pair score (from last slide)
  - Then, the full translation score is:

$$f(y) = h(e(y)) + \sum_{k=1}^{L} g(p_k)$$

- **Goal, compute the best translation**

$$y^*(x) = \arg \max_{y \in \mathcal{Y}(x)} f(y)$$

# Phrase Scoring

$$g(f, e) = \log \frac{c(e, f)}{c(e)}$$

$$g(\text{les chats}, \text{cats}) = \log \frac{c(\text{cats}, \text{les chats})}{c(\text{cats})}$$



- Learning weights has been tried, several times:
  - [Marcu and Wong, 02]
  - [DeNero et al, 06]
  - … and others

- Seems not to work well, for a variety of partially understood reasons

- Main issue: big chunks get all the weight, obvious priors don't help
  - Though, [DeNero et al 08]

# Phrase-Based Translation

| 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 | | |

Table 1: #11# the seven - member crew includes astronauts from france and russia .

Scoring: Try to use phrase pairs that have been frequently observed.
Try to output a sentence with frequent English word sequences.

# Phrase-Based Translation

| 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 | | |

Table 1: #11# the seven - member crew includes astronauts from france and russia .

Scoring: Try to use phrase pairs that have been frequently observed.
        Try to output a sentence with frequent English word sequences.

# Phrase-Based Translation

| 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 | | | |

Table 1: #11# the seven - member crew includes astronauts from france and russia .

Scoring:  Try to use phrase pairs that have been frequently observed.
Try to output a sentence with frequent English word sequences.

# Phrase-Based Translation

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

| 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 | | |

Table 1: #11# the seven - member crew includes astronauts from france and russia .

Scoring: Try to use phrase pairs that have been frequently observed.
Try to output a sentence with frequent English word sequences.

# The Pharaoh Decoder

| Maria | no | dio | una | bofetada | a | la | bruja | verde |
|-------|-----|-----|-----|----------|---|-----|-------|-------|

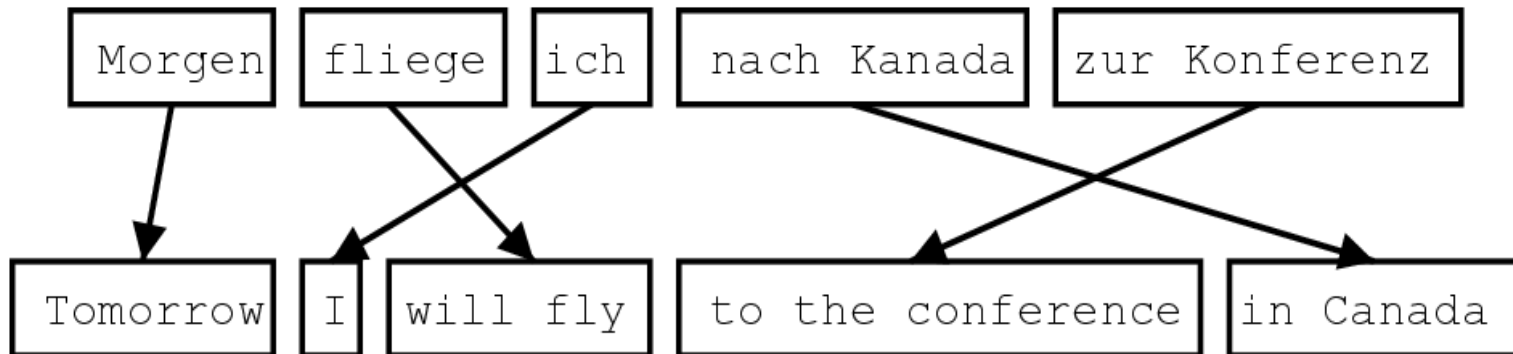Mary    not    give    a    slap    to    the    witch    green
        did not         a slap    by                      green witch
        no            slap        to the
        did not give               to
                                   the
                    slap              the witch

| Maria | no | dio una bofetada | a la | bruja | verde |
|-------|-----|------------------|------|-------|-------|

| Mary | did not | slap | the | green | witch |
|------|---------|------|-----|-------|-------|

- Scores at each step include LM and TM

# The Pharaoh Decoder

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

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

## Space of possible translations

- Phrase table constrains possible translations
- Output sentence is built left to right
  - but source phrases can match any part of sentence
- Each source word can only be translated once
- Each source word must be translated

# Scoring:

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

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

- In practice, much like for alignment models, also include a distortion penalty
  - Define y = $p_1p_2\ldots p_L$ to be a translation with phrase pairs $p_i$
  - Let $s(p_i)$ be the start position of the foreign phrase
  - Let $t(p_i)$ be the end position of the foreign phrase
  - Define η to be the distortion score (usually negative!)
  - Then, we can define a score *with distortion penalty*:

$$f(y) = h(e(y)) + \sum_{k=1}^{L} g(p_k) + \sum_{k=1}^{L-1} \eta \times |t(p_k) + 1 - s(p_{k+1})|$$

- Goal, compute the best translation

$$y^*(x) = \arg \max_{y \in \mathcal{Y}(x)} f(y)$$

# Hypothesis Expansion

# Hypothesis Explosion!

- **Q:** How much time to find the best translation?
  - Exponentially many translations, in length of source sentence
  - NP-hard, just like for word translation models
  - So, we will use approximate search techniques!

# Hypothesis Lattices

| Maria | no | dio | una | bofetada | a | la | bruja | verde |
|-------|-----|-----|-----|----------|---|-----|-------|-------|

Mary — not — give — a — slap — to — the — witch — green
did not — a slap — by — green witch
no — slap — to the
did not give — to
the
slap — the witch



Can recombine if:
- Last two English words match
- Foreign word coverage vectors match

# Decoder Pseudocode

Initialization: Set beam $Q=\{q_0\}$ where $q_0$ is initial state with no words translated

For i=0 … n-1                 [where n in input sentence length]
- For each state $q \in$ beam(Q) and phrase $p \in$ ph(q)
    1. $q'$=next(q,p)          [compute the new state]
    2. Add(Q,q',q,p)          [add the new state to the beam]

Notes:
- ph(q): set of phrases that can be added to partial translation in state q
- next(q,p): updates the translation in q and records which words have been translated from input
- Add(Q,q',q,p): updates beam, q' is added to Q if it is in the top-n overall highest scoring partial translations

# Decoder Pseudocode

Initialization: Set beam $Q=\{q_0\}$ where $q_0$ is initial state with no words translated

For i=0 … n-1                                  [where n in input sentence length]
•For each state q∈beam(Q) and phrase p∈ph(q)
   1.   q'=next(q,p)                  [compute the new state]
   2.   Add(Q,q',q,p)               [add the new state to the beam]

Possible State Representations:
•Full: q = (e, b, α), e.g. ("Joe did not give," 11000000, 0.092)
- e is the partial English sentence
- b is a bit vector recorded which source words are translated
- α is score of translation so far

# Decoder Pseudocode

Initialization: Set beam $Q=\{q_0\}$ where $q_0$ is initial state with no words translated

For i=0 … n-1                     [where n in input sentence length]
•For each state $q \in$ beam(Q) and phrase $p \in$ ph(q)
   1.  q'=next(q,p)               [compute the new state]
   2.  Add(Q,q',q,p)              [add the new state to the beam]

Possible State Representations:
•Full: $q = (e, b, \alpha)$, e.g. ("Joe did not give," 11000000, 0.092)
•Compact: $q = (e_1, e_2, b, r, \alpha)$ ,
   •  e.g. ("not," "give," 11000000, 4, 0.092)
   •  $e_1$ and $e_2$ are the last two words of partial translation
   •  r is the length of the partial translation
•Compact representation is more efficient, but requires back pointers to get the final translation

# Pruning

Maria no      dio una bofetada      a la      bruja verde

```
e: Mary did not
f: **-------
p: 0.154
```
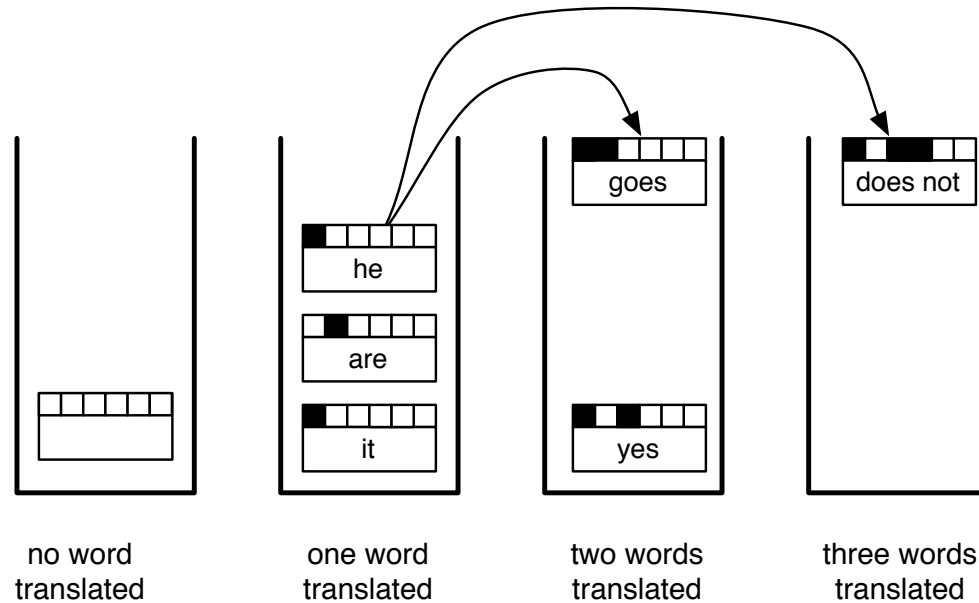**better partial translation**

```
e: the
f: ------**--
p: 0.354
```
**covers easier part --> lower cost**

- Problem: easy partial analyses are cheaper
  - Solution 1: separate bean for each number of foreign words
  - Solution 2: estimate forward costs (A*-like)

# Stack Decoding

## Stacks



| no word translated | one word translated | two words translated | three words translated |

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

# Stack Decoding

## 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**

# Decoder Pseudocode (Multibeam)

Initialization:

- set $Q_0=\{q_0\}$, $Q_i=\{\}$ for I = 1 … n [n is input sent length]

For i=0 … n-1

- For each state $q \in beam(Q_i)$ and phrase $p \in ph(q)$
  1. $q'=next(q,p)$
  2. $Add(Q_j,q',q,p)$ where $j = len(q')$

Notes:

- $Q_i$ is a beam of all partial translations where i input words have been translated
- $len(q)$ is the number of bits equal to one in q (the number of words that have been translated)

# Making it Work (better)

The "Fundamental Equation of Machine Translation" (Brown et al. 1993)

$\hat{e}$ = argmax  $P(e \mid f)$
       e

= argmax  $P(e) \times P(f \mid e) / P(f)$
       e

= argmax  $P(e) \times P(f \mid e)$
       e

# Making it Work (better)

What StatMT people do in the privacy of their own homes

argmax  P(e | f)  =
  e

argmax  P(e) $_x$ P(f | e) / P(f)   =
  e

argmax  P(e)$^{1.9}$ $_x$ P(f | e)        … works better!
  e

Which model are you now paying more attention to?

# Making it Work (better)

## What StatMT people do in the privacy of their own homes

$$\underset{e}{\text{argmax}}\ P(e \mid f)\ =$$

$$\underset{e}{\text{argmax}}\ P(e)\ \times\ P(f \mid e)\ /\ P(f)$$

$$\underset{e}{\text{argmax}}\ P(e)^{1.9}\ \times\ P(f \mid e)\ \times\ 1.1^{\text{length}(e)}$$

Rewards longer hypotheses, since these are 'unfairly' punished by P(e)

# Making it Work (better)

What StatMT people do in the privacy of their own homes

argmax  $P(e)^{1.9}$ x $P(f\,|\,e)$ x $1.1^{length(e)}$ x $KS^{3.7}$ …

e

Lots of knowledge sources vote on any given hypothesis. Each has a weight

"Knowledge source" = "feature function" = "score component".
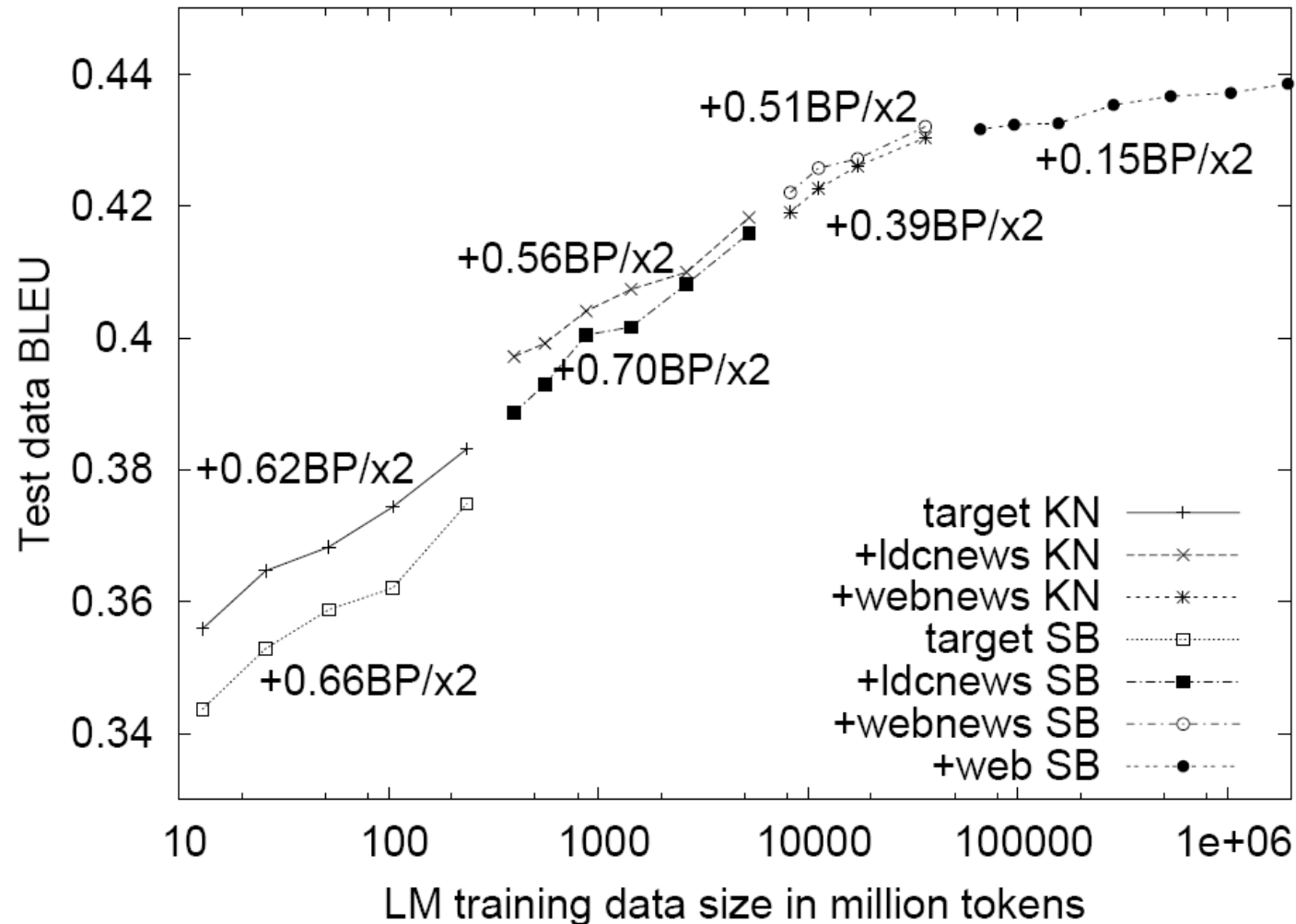
# Making it Work (better)

## Log-linear feature-based MT

argmax$_e$ 1.9×log P(e) + 1.0×log P(f | e) +

1.1× log length(e) + 3.7×KS + …

= argmax$_e$ $\Sigma_i$ $w_i f_i$

So, we have two things:

– "Features" *f*, such as log language model score
– A weight *w* for each feature that indicates how good a job it does at indicating good translations
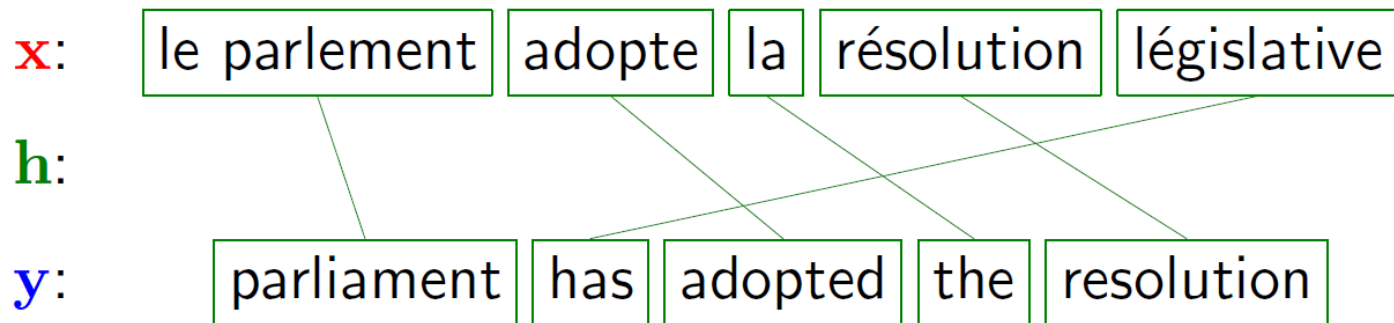
# No Data like More Data!



- Discussed for LMs, but can new understand full model!

# Tuning for MT

- Features encapsulate lots of information
    - Basic MT systems have around 6 features
    - P(e|f), P(f|e), lexical weighting, language model

- How to tune feature weights?

- Idea 1: Use your favorite classifier
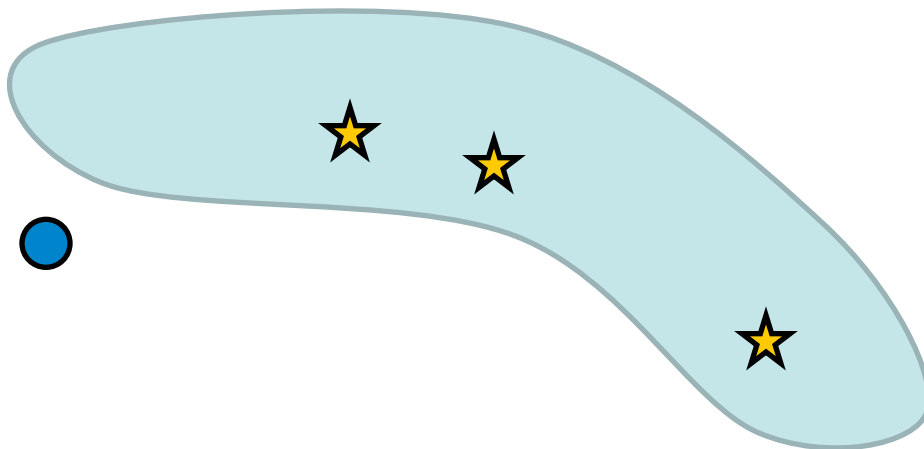
# Why Tuning is Hard

- Problem 1: There are latent variables
  - Alignments and segmentations
  - Possibility: forced decoding (but it can go badly)

$\mathbf{x}$: | le parlement | adopte | la | résolution | législative |

$\mathbf{h}$:

$\mathbf{y}$: | parliament | has | adopted | the | resolution |

# Why Tuning is Hard

- ## Problem 2: There are many right answers
  - The reference or references are just a few options
  - No good characterization of the whole class



  - BLEU isn't perfect, but even if you trust it, it's a corpus-level metric, not sentence-level

# Linear Models: Perceptron

- **The perceptron algorithm**
  - Iteratively processes the training set, reacting to training errors
  - Can be thought of as trying to drive down training error
- **The (online) perceptron algorithm:**
  - Start with zero weights
  - Visit training instances $(x_i, y_i)$ one by one
    - Make a prediction

$$y^* = \arg\max_y w \cdot \phi(x_i, y)$$

    - If correct ($y^* == y_i$): no change, goto next example!
    - If wrong: adjust weights

$$w = w + \phi(x_i, y_i) - \phi(x_i, y^*)$$

# Perceptron training

For each training example $(\mathbf{x}, \mathbf{y})$: [Collins '02]

$$\mathbf{w} \leftarrow \mathbf{w} \ +\Phi(\mathbf{x}, \mathbf{y}_t) \qquad \mathbf{y}_t \quad = \mathbf{y}$$
$$-\Phi(\mathbf{x}, \mathbf{y}_p) \qquad \mathbf{y}_p \quad = \text{DECODE}(\mathbf{x})$$

$$\mathbf{w} \leftarrow \mathbf{w} \ +\Phi(\mathbf{x}, \mathbf{y}_t, \mathbf{h}_t) \quad | \quad \mathbf{y}_t, \mathbf{h}_t \ = \ \textbf{???}$$
$$-\Phi(\mathbf{x}, \mathbf{y}_p, \mathbf{h}_p) \quad | \quad \mathbf{y}_p, \mathbf{h}_p = \text{DECODE}(\mathbf{x})$$

# Update strategies

$$\mathbf{w} \leftarrow \mathbf{w} + \Phi(\mathbf{x}, \boxed{\mathbf{y}_t, \mathbf{h}_t}) - \Phi(\mathbf{x}, \mathbf{y}_p, \mathbf{h}_p)$$

Training example (reference)

$\mathbf{x}$: voté sur demande d ' urgence

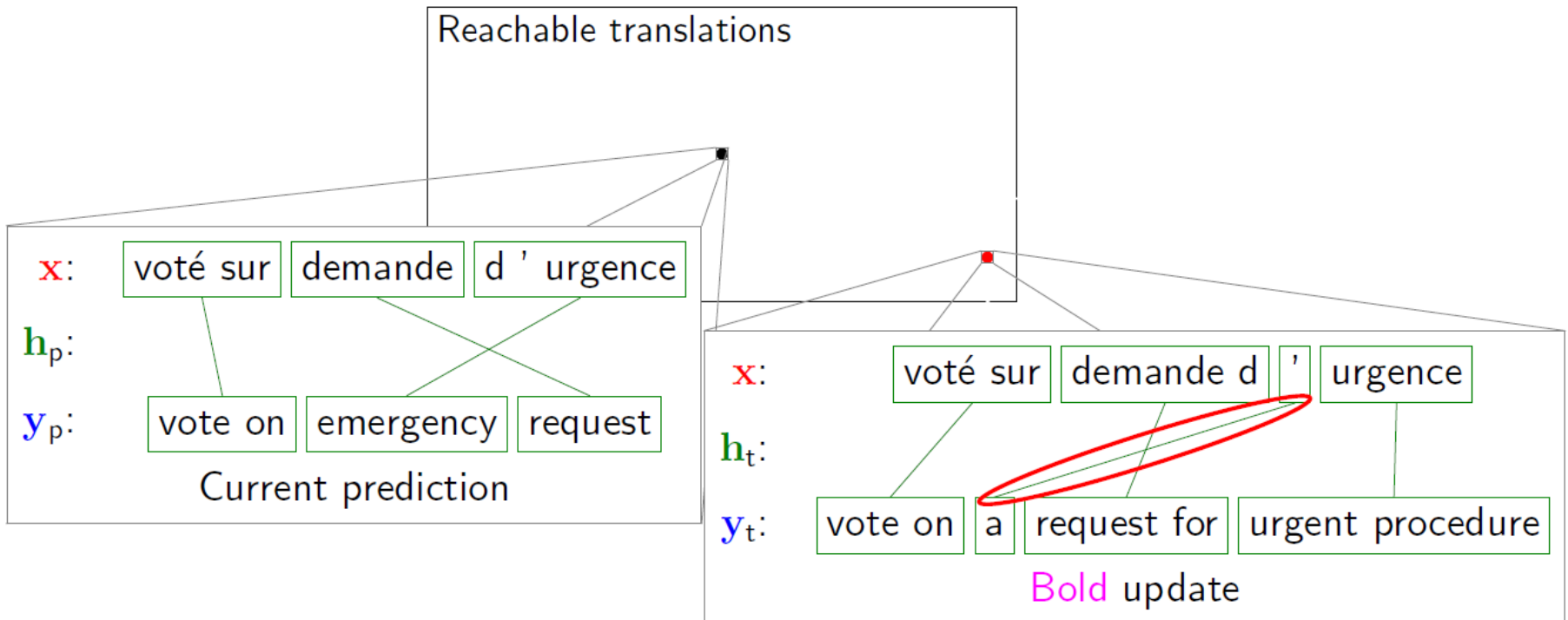$\mathbf{y}$: vote on a request for urgent procedure

# Update strategies

$$\mathbf{w} \leftarrow \mathbf{w} + \Phi(\mathbf{x}, \boxed{\mathbf{y}_t, \mathbf{h}_t}) - \Phi(\mathbf{x}, \mathbf{y}_p, \mathbf{h}_p)$$

Training example (reference)

$\mathbf{x}$: voté sur demande d ' urgence
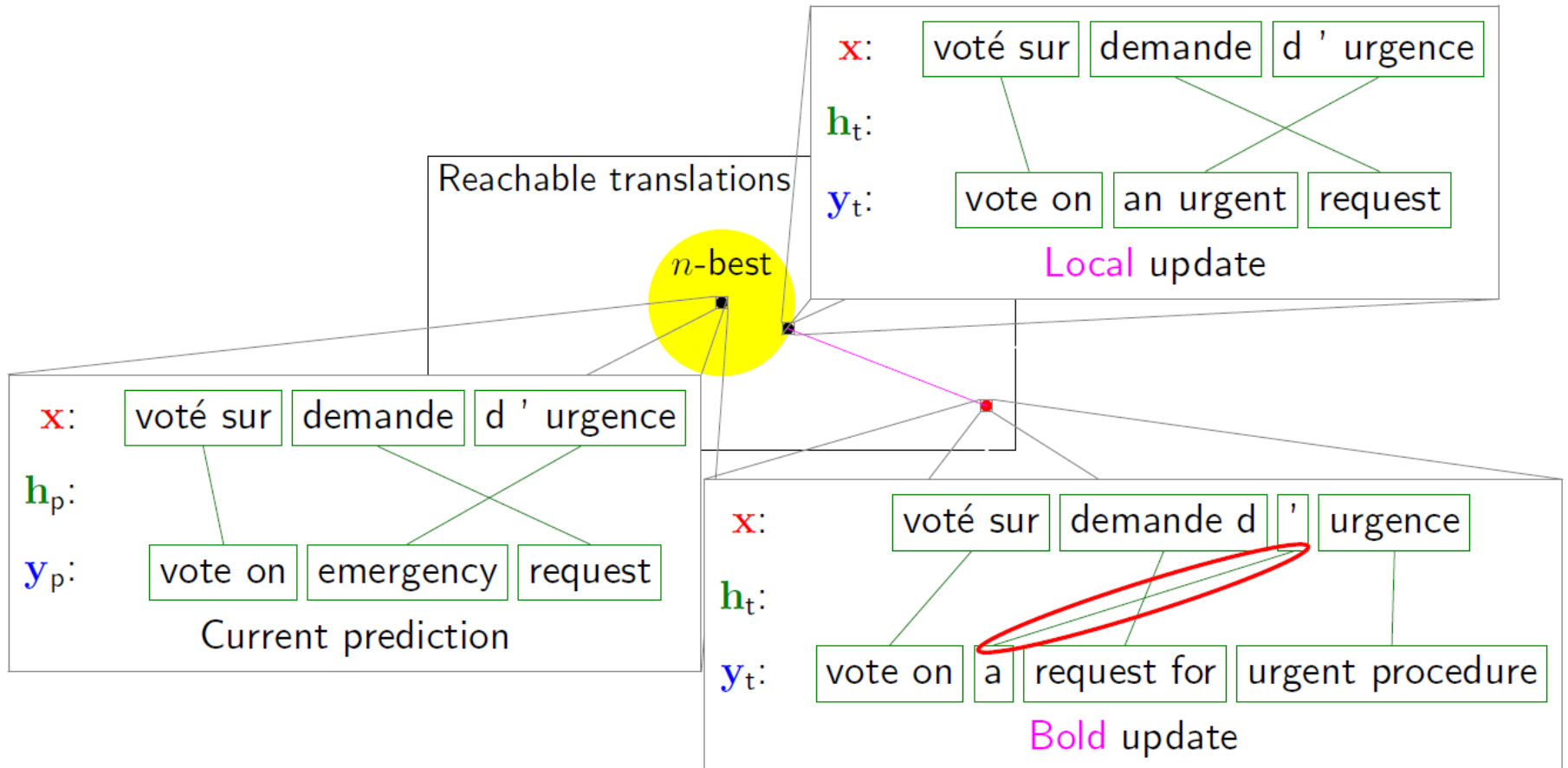
$\mathbf{y}$: vote on a request for urgent procedure

# Update strategies

$$\mathbf{w} \leftarrow \mathbf{w} + \Phi(\mathbf{x}, \boxed{\mathbf{y}_t, \mathbf{h}_t}) - \Phi(\mathbf{x}, \mathbf{y}_p, \mathbf{h}_p)$$

Training example (reference)

$\mathbf{x}$: voté sur demande d ' urgence
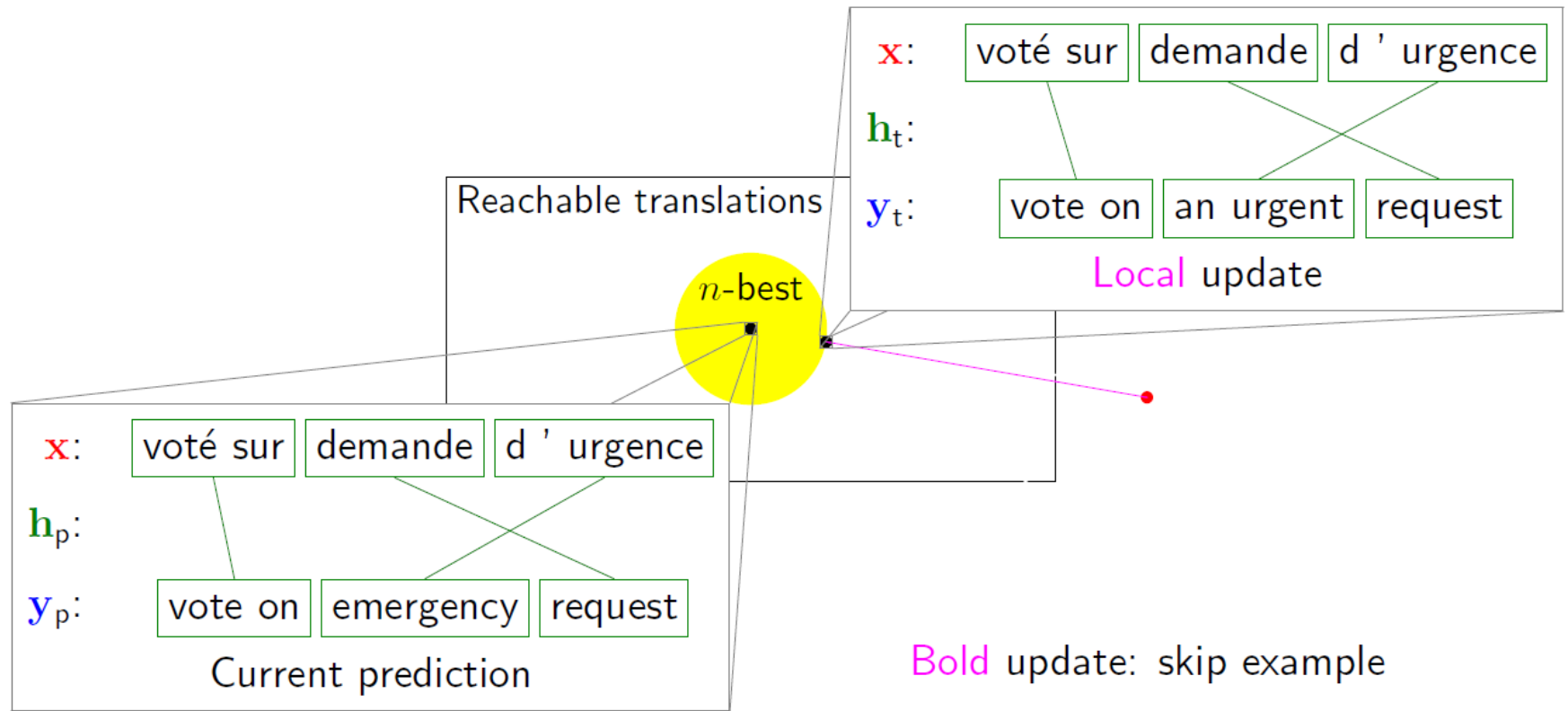
$\mathbf{y}$: vote on a request for urgent procedure



Reachable translations

$n$-best

$\mathbf{x}$: | voté sur | demande | d ' urgence |

$\mathbf{h}_t$:

$\mathbf{y}_t$: | vote on | an urgent | request |

Local update

$\mathbf{x}$: | voté sur | demande | d ' urgence |

$\mathbf{h}_p$:

$\mathbf{y}_p$: | vote on | emergency | request |

Current prediction

$\mathbf{x}$: | voté sur | demande d | ' | urgence |

$\mathbf{h}_t$:

$\mathbf{y}_t$: | vote on | a | request for | urgent procedure |

Bold update

# Update strategies

$$\mathbf{w} \leftarrow \mathbf{w} + \Phi(\mathbf{x}, \boxed{\mathbf{y}_t, \mathbf{h}_t}) - \Phi(\mathbf{x}, \mathbf{y}_p, \mathbf{h}_p)$$

Training example (reference)
$\mathbf{x}$: voté sur demande d ' urgence
$\mathbf{y}$: vote on a request for urgent procedure

# Update strategies

$$\mathbf{w} \leftarrow \mathbf{w} + \Phi(\mathbf{x}, \boxed{\mathbf{y}_t, \mathbf{h}_t}) - \Phi(\mathbf{x}, \mathbf{y}_p, \mathbf{h}_p)$$

Training example (reference)

$\mathbf{x}$: voté sur demande d'urgence

$\mathbf{y}$: vote on a request for urgent procedure

$\mathbf{x}$:

$\mathbf{h}_t$:

| Decoder | Bold | Local |
|---|---|---|
| Monotonic | 34.3 | **34.6** |
| Limited distortion | 33.5 | **34.7** |

$\mathbf{x}$:

$\mathbf{h}_p$:

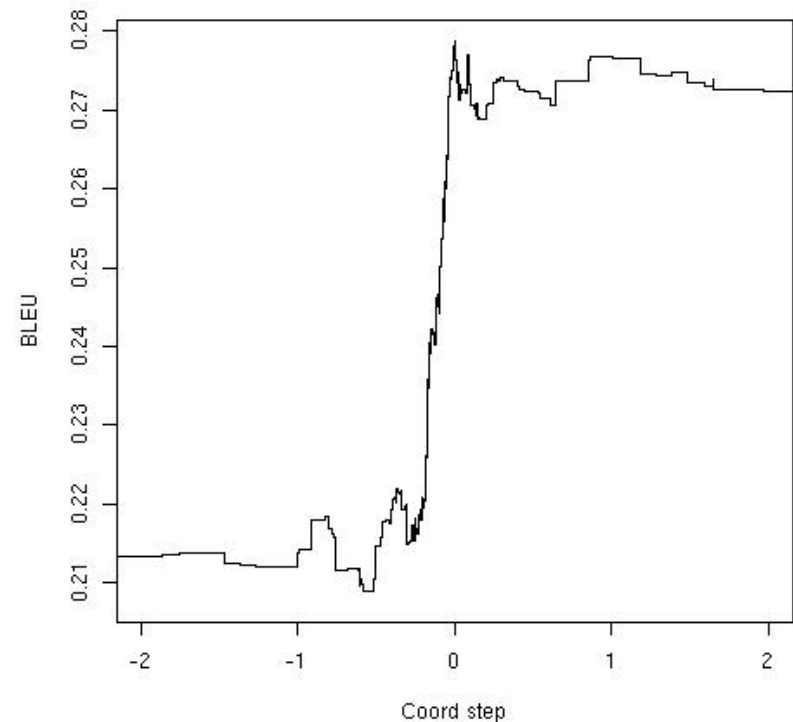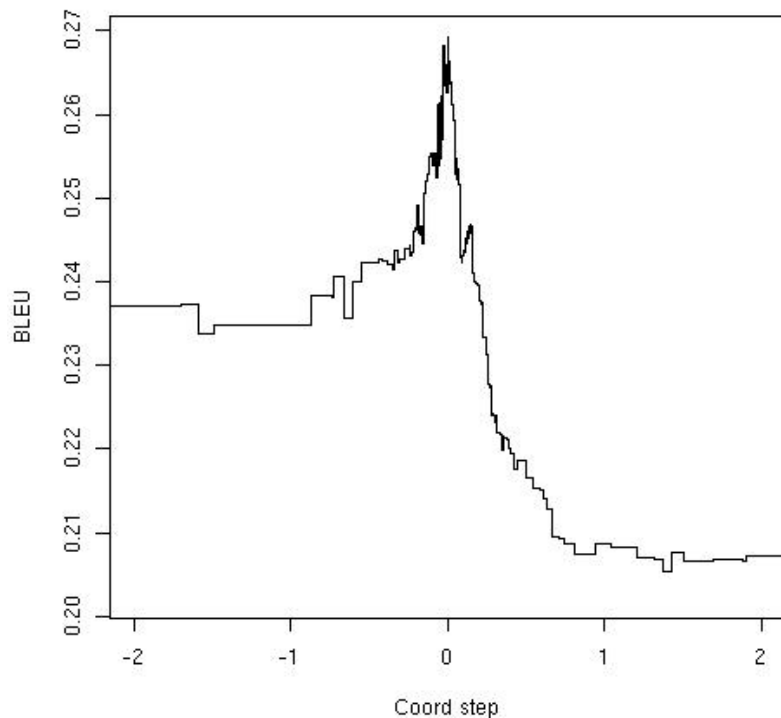$\mathbf{y}_p$:

Current prediction

Bold update: skip example

# Why Tuning is Hard

- Problem 3: Computational constraints
  - Discriminative training involves repeated decoding
  - Very slow!  So people tune on sets much smaller than those used to build phrase tables

# Minimum Error Rate Training

- Standard method: minimize BLEU directly (Och 03)
  - MERT is a discontinuous objective
  - Only works for max ~10 features, but works very well then
  - Here: k-best lists, but forest methods exist (Machery et al 08)

# MERT: Convex Upper Bound of BLEU