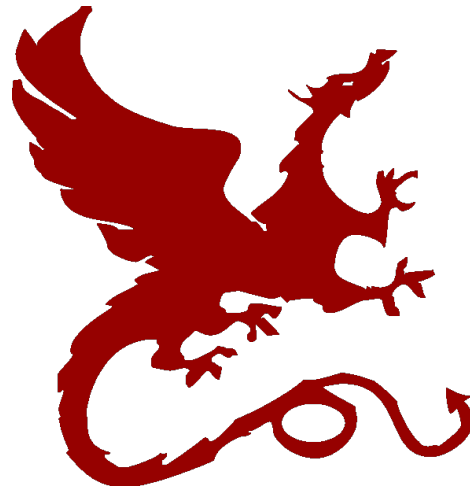# Algorithms for NLP



# Machine Translation II

Taylor Berg-Kirkpatrick – CMU

Slides: Dan Klein – UC Berkeley
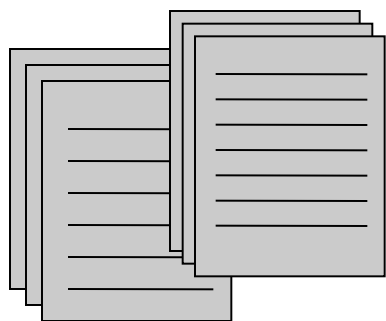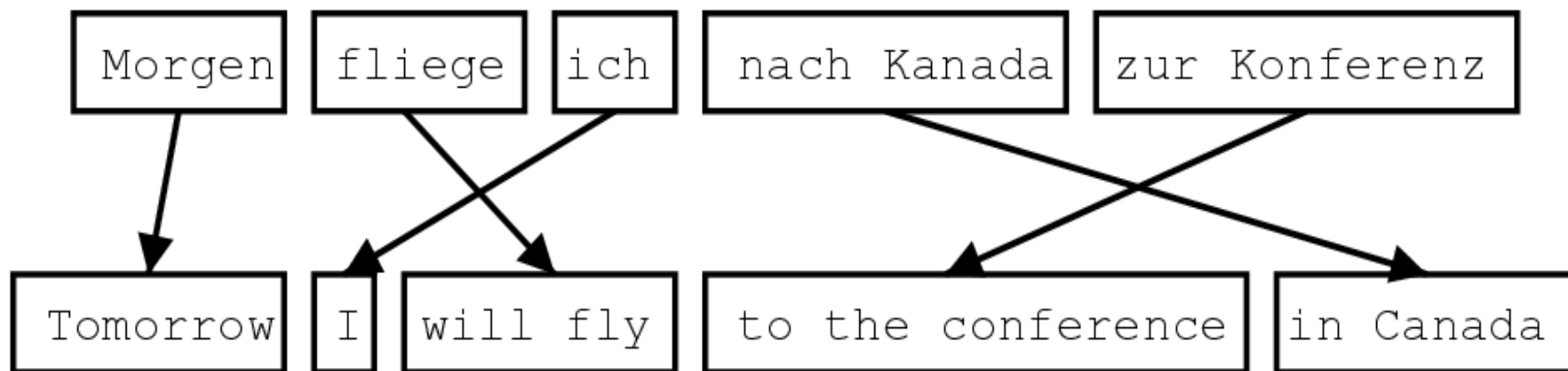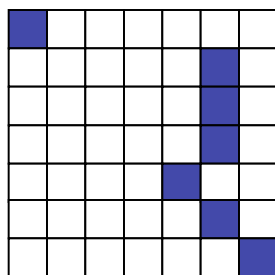
# Announcments

- Project 1 grades out today
  - Requirements (6 points)
  - Write-up (2 points)
  - Clearly written (1 point)
  - Exceeds reqs (1 point)
  - Extra (unbounded)
- Out of 9 points

# Phrase-Based System Overview

Morgen → Tomorrow

fliege → will fly

ich → I

nach Kanada → in Canada

zur Konferenz → to the conference

| | 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)

Many slides and examples from Philipp Koehn or John DeNero

# Word Alignment

# IBM Models 1/2

|  | 1 | 2 | 3 | 4 | 5 | 6 | 7 | 8 | 9 |
|---|---|---|---|---|---|---|---|---|---|
| **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)

Japan shaken by two new quakes

Le Japon secoué par deux nouveaux séismes

Japan is at the junction of four tectonic plates

Le Japon est au confluent de quatre plaques tectoniques

# 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.

# The HMM Model

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

*(word positions: 1 Thank, 2 you, 3 comma, 4 I, 5 shall, 6 do, 7 so, 8 gladly, 9 period)*

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

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

## Model Parameters

*Emissions:* $P(F_1 = \text{Gracias} \mid E_{A_1} = \text{Thank})$     *Transitions:* $P(A_2 = 3 \mid A_1 = 1)$

# 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?

# HMM Examples

# AER for HMMs

| Model | AER |
|---|---|
| Model 1 INT | 19.5 |
| HMM E→F | 11.4 |
| HMM F→E | 10.8 |
| HMM AND | 7.1 |
| HMM INT | 4.7 |
| GIZA M4 AND | 6.9 |

# Models 3, 4, and 5: 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]

# Examples: Translation and Fertility

*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*

he is nodding

il hoche la tête

| 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 | | |

# 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 | | |

# Some Results

- <span style="color:blue">[Och and Ney 03]</span>

| Model | Training scheme | 0.5K | 8K | 128K | 1.47M |
|-------|-----------------|------|-----|------|-------|
| Dice | | 50.9 | 43.4 | 39.6 | 38.9 |
| Dice+C | | 46.3 | 37.6 | 35.0 | 34.0 |
| Model 1 | $1^5$ | 40.6 | 33.6 | 28.6 | 25.9 |
| Model 2 | $1^5 2^5$ | 46.7 | 29.3 | 22.0 | 19.5 |
| HMM | $1^5 H^5$ | 26.3 | 23.3 | 15.0 | 10.8 |
| Model 3 | $1^5 2^5 3^3$ | 43.6 | 27.5 | 20.5 | 18.0 |
| | $1^5 H^5 3^3$ | 27.5 | 22.5 | 16.6 | 13.2 |
| Model 4 | $1^5 2^5 3^3 4^3$ | 41.7 | 25.1 | 17.3 | 14.1 |
| | $1^5 H^5 3^3 4^3$ | 26.1 | 20.2 | 13.1 | 9.4 |
| | $1^5 H^5 4^3$ | 26.3 | 21.8 | 13.3 | 9.3 |
| Model 5 | $1^5 H^5 4^3 5^3$ | 26.5 | 21.5 | 13.7 | 9.6 |
| | $1^5 H^5 3^3 4^3 5^3$ | 26.5 | 20.4 | 13.4 | 9.4 |
| Model 6 | $1^5 H^5 4^3 6^3$ | 26.0 | 21.6 | 12.8 | 8.8 |
| | $1^5 H^5 3^3 4^3 6^3$ | 25.9 | 20.3 | 12.5 | 8.7 |

# Phrase-Based MT

# Phrase-Based Translation Overview

**Input:** lo haré | rápidamente | .

**Translations:** I'll do it | quickly | .

quickly | I'll do it | .

The decoder...

tries different segmentations,

translates phrase by phrase,

and considers reorderings.

**Objective:**

$$\arg\max_{\mathbf{e}}\left[P(\mathbf{f}|\mathbf{e}) \cdot P(\mathbf{e})\right]$$

$$\arg\max_{\mathbf{e}}\left[\prod_{\langle\bar{e},\bar{f}\rangle} P(\bar{f}|\bar{e}) \cdot \prod_{i=1}^{|\mathbf{e}|} P(e_i|e_{i-1}, e_{i-2})\right]$$
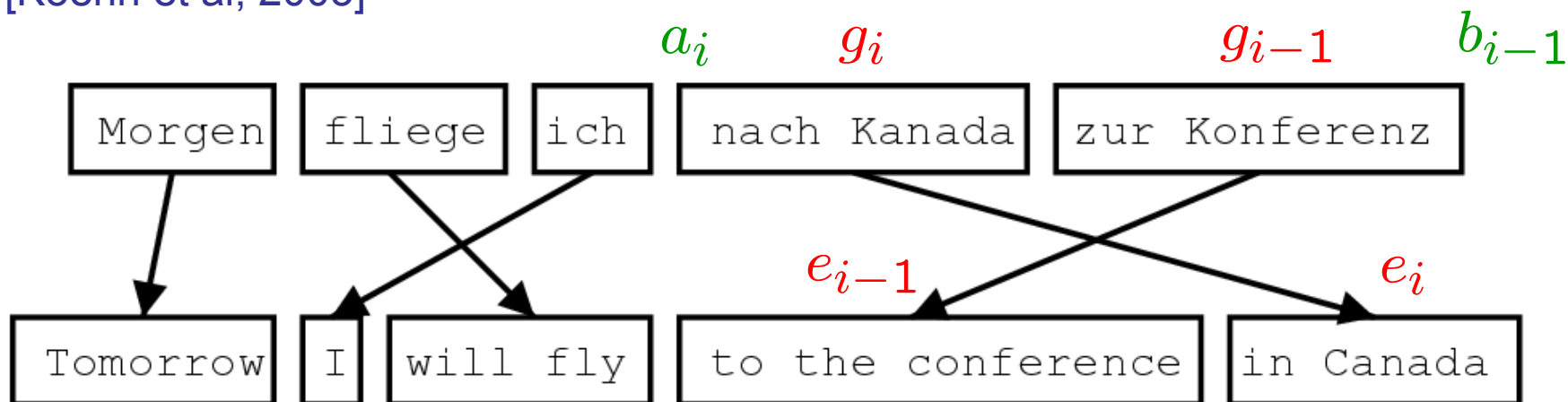
# Phrase-Based Decoding

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 | | | | |

Decoder design is important: [Koehn et al. 03]

# The Pharaoh "Model"

[Koehn et al, 2003]

$$a_i \qquad g_i \qquad\qquad g_{i-1} \qquad b_{i-1}$$

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

$$e_{i-1} \qquad\qquad e_i$$

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

$$P(e|g) = P(\{\bar{g}_i\}|g) \prod_i \phi(\bar{e}_i|\bar{g}_i) d(a_i - b_{i-1})$$

Segmentation  Translation  Distortion

# The Pharaoh "Model"

$$P(f|e) = P(\{\bar{e}_i\}|e) \prod_i \phi(\bar{f}_i|\bar{e}_i) d(a_i - b_{i-1})$$

$$\frac{1}{K}$$

$$\frac{count(\bar{f}_i, \bar{e}_i)}{count(\bar{e}_i)}$$

$$\alpha^{|a_i - b_{i-1}|}$$

*Where do we get these counts?*

# Phrase Weights

How the MT community estimates $P(\bar{f}|\bar{e})$

*Parallel training sentences*        *provide phrase pair counts.*

Gracias , <u>lo haré</u> de muy buen grado .

Thank you , <u>I shall do so</u> gladly .

lo haré ⟺ I shall do so

*44 times in the corpus*

*All phrase pairs are counted,*        *and counts are normalized.*

Gracias , lo haré de muy buen grado .

Thank you , I shall do so gladly .

$$P(\bar{f}|\bar{e}) = \frac{\operatorname{count}(\bar{f},\bar{e})}{\operatorname{count}(\bar{e})}$$

# Phrase-Based Decoding

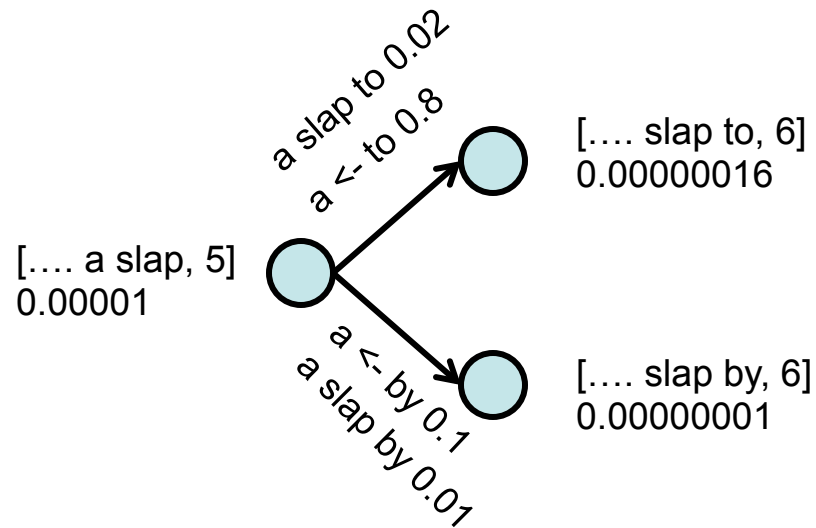| Maria | no | dio | una | bofetada | a | la | bruja | verde |
|-------|-----|-----|-----|----------|---|-----|-------|-------|

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

# Monotonic Word Translation

| Maria | no | dio | una | bofetada | a | la | bruja | verde |
|-------|-----|-----|-----|----------|---|-----|-------|-------|
| Mary | not | give | a | slap | to | the | witch | green |
| | did not | | | | by | | | |
| | no | | | | | | | |

- Cost is LM * TM
- It's an HMM?
  - $P(e|e_{-1}, e_{-2})$
  - $P(f|e)$
- State includes
  - Exposed English
  - Position in foreign
- Dynamic program loop?

[…. a slap, 5]
0.00001

a slap to 0.02
a <- to 0.8

[…. slap to, 6]
0.00000016

a <- by 0.1
a slap by 0.01

[…. slap by, 6]
0.00000001

```
for (fPosition in 1…|f|)
  for (eContext in allEContexts)
    for (eOption in translations[fPosition])
      score = scores[fPosition-1][eContext] * LM(eContext+eOption) * TM(eOption, fWord[fPosition])
      scores[fPosition][eContext[2]+eOption] =max  score
```

# Beam Decoding

- For real MT models, this kind of dynamic program is a disaster (why?)

- Standard solution is beam search: for each position, keep track of only the best k hypotheses

```
for (fPosition in 1…|f|)
  for (eContext in bestEContexts[fPosition])
    for (eOption in translations[fPosition])
      score = scores[fPosition-1][eContext] * LM(eContext+eOption) * TM(eOption, fWord[fPosition])
      bestEContexts.maybeAdd(eContext[2]+eOption, score)
```

- Still pretty slow… why?

- Useful trick: cube pruning (Chiang 2005)



Example from David Chiang

# Phrase Translation

| 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 | | |

- **If monotonic, almost an HMM; technically a semi-HMM**

      for (fPosition in 1…|f|)
        for (lastPosition < fPosition)
          for (eContext in eContexts)
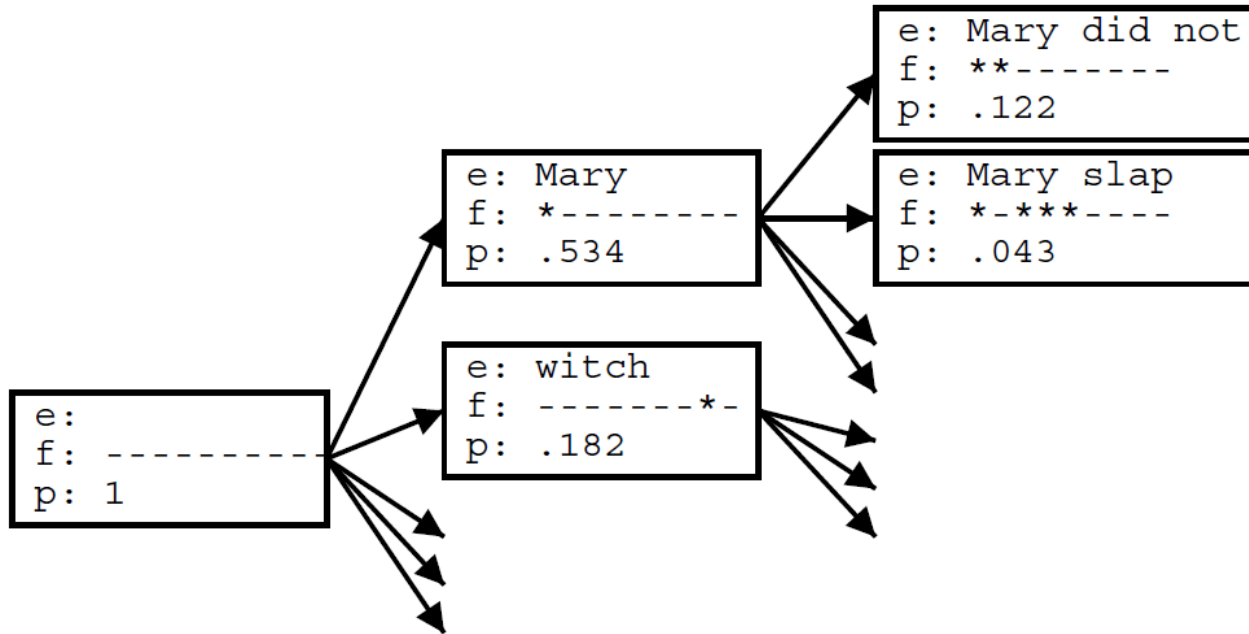            for (eOption in translations[fPosition])
              … combine hypothesis for (lastPosition ending in eContext) with eOption

- **If distortion… now what?**

```
                                              ┌──────────────────┐
                                              │ e: Mary did not  │
                                         ────►│ f: **-------     │
                                              │ p: .122          │
                                              └──────────────────┘
              ┌──────────────────┐
              │ e: Mary          │           ┌──────────────────┐
         ────►│ f: *---------    │           │ e: Mary slap     │
              │ p: .534          │──────────►│ f: *-***----     │
              └──────────────────┘           │ p: .043          │
┌──────────────────┐                         └──────────────────┘
│ e:               │
│ f: ----------    │
│ p: 1             │
└──────────────────┘
              ┌──────────────────┐
              │ e: witch         │
         ────►│ f: -------*-     │
              │ p: .182          │
              └──────────────────┘
```

```
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: use beams per foreign subset
  - Solution 2: estimate forward costs (A*-like)

| 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 |
|------|---------|------|-----|-------|-------|

# Hypotheis 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 |  |  |