#### **Decoding in Phrase-Based Machine** Translation (Building the translation) Not required for the homework



- Noisy channel model:  $P(e|f) \propto P(f|e) P(e)$ (ignore P(f) term) Translation Language model (TM) model (LM)
- Inputs needed
  - $\blacktriangleright$  Language model that scores I
- Phrase table: set of phrase pairs (e, f) with probabilities P(f|e) What we want to find: e produced by a series of phrase-by-phrase
- translations from an input **f**

## Phrase-Based Decoding

$$P(e_i|e_1,\ldots,e_{i-1}) \approx P(e_i|e_{i-n-1},\ldots,e_i)$$

### Phrase Lattice





- translations of all possible spans
- paths in the lattice that don't skip any words
- Looks like Viterbi, but the scoring is more complicated

| ofetada   | a         | la         | bruja | verde                 |
|-----------|-----------|------------|-------|-----------------------|
| slap<br>p | to<br>by  | <u>the</u> | witch | <u>green</u><br>witch |
|           | t.o<br>t. | to the to  |       |                       |
|           | t.ł       | the v      | vitch |                       |

• Given an input sentence, look at our phrase table to find all possible

Monotonic translation: need to translate each word in order, explore

Koehn (2004)



beam?

#### Score

- Where are we in the sentence
- What words have we produced so far (actually only need to remember the last 2 words when using a 3-gram LM)

If we translate with beam search, what state do we need to keep in the

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







| ofetada   | a        | la            | bruja | verde     |
|-----------|----------|---------------|-------|-----------|
| slap<br>D | to<br>by | <u>t.he</u>   | witch | <br>witch |
|           | t.o      | <u>to the</u> |       |           |
|           |          | ne            |       |           |
|           |          | the v         | vitch |           |

- Beam state: where we're at, what -0.1 the current translation so far is, and score of that translation
  - Advancing state consists of trying each possible translation that could get us to this timestep













| ofetada   | a        | la         | bruja | verde                 |
|-----------|----------|------------|-------|-----------------------|
| slap<br>p | to<br>by | <u>the</u> | witch | <u>green</u><br>witch |
|           | t.o      | to the to  |       |                       |
|           | t.ł      | ie         |       |                       |
|           |          | the w      | vitch |                       |

Two ways to get here: Maria + no dio or Maria no + dio

 Beam is filled with options from multiple segmentations of input





| ofetada   | a        | la         | bruja | verde                 |
|-----------|----------|------------|-------|-----------------------|
| slap<br>p | to<br>by | <u>the</u> | witch | <u>green</u><br>witch |
|           | to<br>t  | the<br>o   |       |                       |
|           | t.ł      | ne         |       |                       |
|           |          | the the    | vitch |                       |

With beam size k = infinity, how many possible entries in the beam here?





- State needs to describe which
- Big enough phrases already isn't as important as you think



#### score = $\alpha \log P(t) + \beta \log P(s|t)$ ...and P(s|t) is in fact more complex

- Usually 5-20 feature weights to set, want to optimize for BLEU score which is not differentiable
- MERT (Och 2003): decode to get 1000best translations for each sentence in a small training set (<1000 sentences), do line search on parameters to directly optimize for BLEU

# **Training Decoders**







- Toolkit for machine translation due to Philipp Koehn + Hieu Hoang
  Pharaoh (Koehn, 2004) is the decoder from Koehn's thesis
- Moses implements word alignment, language models, and this decoder, plus a ton more stuff
  - Highly optimized and heavily engineered, could more or less build SOTA translation systems with this from 2007-2013



| SOURCE     | Cela constituerait une solution conduire à terme à une charte                 |
|------------|---|
| HUMAN      | That would be an interim soluti work towards a binding charter                |
| 1x DATA    | [this] [constituerait] [assistance<br>[licences] [to] [terme] [to] [a] [cl    |
| 10x DATA   | [it] [would] [a solution] [transitie<br>[to] [term] [to a] [charter] [to] [va |
| 100x DATA  | [this] [would be] [a transitional<br>charter] [legally binding] [.]           |
| 1000x DATA | [that would be] [a transitional s<br>lead to] [a binding charter] [.]         |

#### Moses

- transitoire qui permettrait de à valeur contraignante.
- ion which would make it possible to r in the long term.
- e] [transitoire] [who] [permettrait] harter] [to] [value] [contraignante] [.]
- onal] [which] [would] [of] [lead] alue] [binding] [.]
- solution] [which would] [lead to] [a

solution] [which would] [eventually

slide credit: Dan Klein

# Evaluating MT

# Evaluating MT



- Fluency: does it sound good in the target language?
- Fidelity/adequacy: does it capture the meaning of the original?
- Automatic evaluation tries to approximate this...
- BLEU score: geometric mean of 1-, 2-, 3-, and 4-gram precision vs. a reference, multiplied by brevity penalty (penalizes short translations) I-gram precision: do you predict words that are in the reference?
- - 4-gram precision: to get this right, you need those words to be in the right order!
- Better metrics: human-in-the-loop variants

# Syntactic MT



 $DI_1$ 

- Rather than use phrases, use a synchronous context-free grammar
- $NP \rightarrow [DT_1 JJ_2 NN_3; DT_1 NN_3 JJ_2]$
- $DT \rightarrow [the, la]$
- $DT \rightarrow [the, le]$
- $NN \rightarrow [car, voiture]$
- $JJ \rightarrow [yellow, jaune]$

yellow the la voiture jaune car

- other half
- Assumes parallel tree structures, but there can be reordering

### Syntactic MT



Translation = parse the input with "half" of the grammar, read off the







- Use lexicalized rules, look like "syntactic phrases"
- Leads to HUGE grammars, parsing is slow

## Syntactic MT

#### Grammar

 $s \rightarrow \langle VP .; | VP . \rangle \circ R s \rightarrow \langle VP .; you VP . \rangle$ VP -> ( lo haré ADV ; will do it ADV ) S → 〈 lo haré ADV . ; I will do it ADV . 〉 ADV  $\rightarrow$  ( de muy buen grado ; gladly ) Slide credit: Dan Klein

