

# Agenda

- Language in 10
- HW5 notes
- Kristina Toutanova speaking tomorrow

# Morphology and Translation

April 24, 2014

# Today's goals

- Have a basic understanding of morphology in languages, along with some of the complexities it introduces for MT
- Cover a few approaches to morphological *analysis* (translation from rich → poor)
- And *generation* (poor → rich)

# Motivation

- To this point, we have treated words as atomic white-space delimited units, with no relationships among them
- This hides a lot of information, since words are related

*house*  $\Leftrightarrow$  *Haus*

*houses*  $\Leftrightarrow$  *Hause*

- ...which information is hidden from the computer

# Example

*Das ist ein kleines Haus .*

*That is a small house .*

# Example

*174 19 182 40626 991 50*

*192 4 19 27 200 49*

# Morphology

- *Morphology: the study of the forms of words*
  - **Inflectional** – words change to reflect grammatical roles
    - e.g., *groß, große, großem, großen, größer, großes*
  - **Derivational** – shared semantics, often across PsOS
    - e.g., *employ (V), employee, employer (N), employable (JJ)*
- **Lemma** – the basic, canonical form of the word (*correr*)
- **Stem** – the shared base form (often a prefix) across inflectional variants
  - e.g., *corr-* (*Spanish*)

# Related problem: tokenization

- Morphology is not the only means by which data are unnecessarily fragmented
- **Tokenization** is largely a task of splitting off punctuation
  - e.g., **house**, becomes **house** ,
  - **“No,” he said.** becomes **“ No , ” he said .**
- A related step, **normalization**, removes case distinctions, standardizes character sets (e.g., quotations, numerals)
- These are largely deterministic processes that are also important for aggregating statistics, but they are largely artifacts of *written* language



# Simple morphology: English

- Words are inflected for
  - case (objective, accusative, genitive)  
*I, me, my/mine, 's*
  - tense (past, present, or future)  
*-ed, -ing, will*
  - person (1st, 2nd, 3rd)  
*I, you, he/she/they*
  - number (singular vs. plural)  
*-s*

# Complex morphology: German

- Inflections of the English definite determiner *the*: *the*
- Inflections of the German definite determiner *der*:

Case	Singular			Plural		
	male	fem.	n.	male	fem.	n.
nominative (subject)	der	die	das	die	die	die
genitive (possessive)	des	der	des	der	der	der
dative (indirect object)	dem	der	dem	den	den	den
accusative (direct object)	den	die	das	die	die	die

**Figure 2.6** Morphology of the definite determiner in German (in English always *the*). It varies depending on count, case, and gender. Each word form is highly ambiguous: *der* is male singular nominative, but also female singular genitive/dative, as well as plural genitive for any gender.

# Worse: Arabic

- Concept defined by three consonants
- Example inflectional morphology:

- concept: ktb (*to write*)

• kataba	he wrote	<i>CaCaCa</i>
katabna	we wrote	<i>CaCaCna</i>
katabuu	they wrote	<i>CaCaCuu</i>
yaktubu	he writes	<i>yaCCuCu</i>
naktubu	we write	<i>naCCuCu</i>
yaktabuuna	they write	<i>yaCCaCuuna</i>
sayaktubu	he will write	<i>sayaCCuCu</i>
sanaktubu	we will write	<i>sanaCCuCu</i>
sayaktabuuna	they will write	<i>sayaCCaCuuna</i>

# MRLs and MT

- Morphologically rich languages (MRLs) cause lots of problems for MT
- *Data sparsity*: alignments to words in the other language are needlessly divided, fracturing statistics

- Common relationships are hidden

houses = plural(house)

was = past-tense(is)

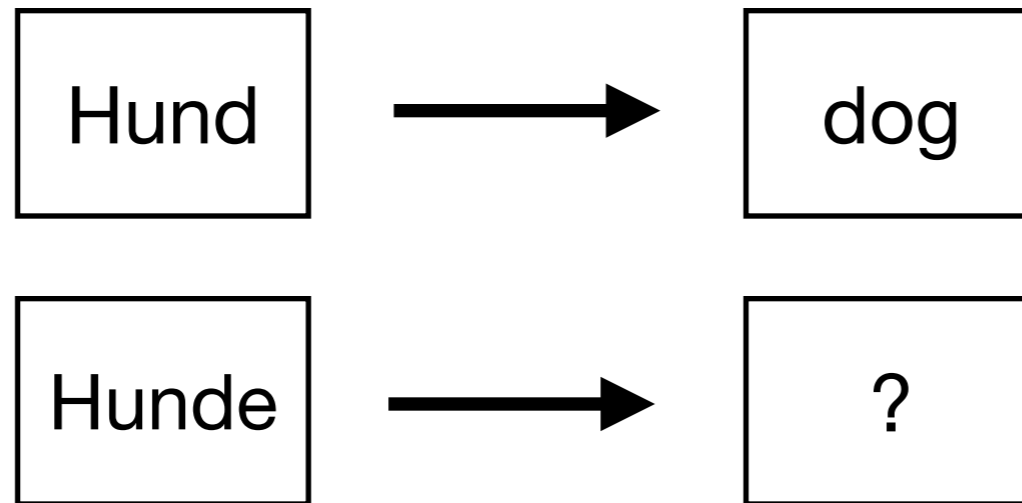
children = plural(child)

- Data is fragmented

	English	German	Finnish
Vocabulary size	65,888	195,290	358,344
Unknown word rate	0.22%	0.78%	1.82%

**Figure 10.4** Vocabulary size and effect on the unknown word rate: Numbers reported for 15 million words of the Europarl corpus for vocabulary collection and unknown word rate on additional 2000 sentences (data from the 2005 ACL workshop shared task).

- *Unseen inflections*: complex inflectional morphology may result in particular versions of a word not being seen



- *Long distance dependencies.* Richer morphology trades off with word order, which can mean fewer local dependencies captureable by ngrams and phrasal translations

# Approaches

- Two focal points for morphology in MT
  - *Analysis*: when translating **from** an MRL to a morphologically poor one (e.g., German to English)
  - *Generation*: translating **into** an MRL (e.g., English to Czech)



# Translation from MRLs

- The “easy” direction
- Common approaches
  - Lemmatization and stemming
  - Compound splitting
  - Lattice translation

# Lemmatization and Stemming

- *Lemmatization* maps inflected forms back to their lemmas

- e.g., am, are, is → be

*The chicken's chicks are different sizes  
→ the chicken chick is differ size*

- This requires a (language-specific) lemmatizer that knows how to remove morph. features

- If you don't have a morphological analyzer, a poor man's approximation is to simply truncate the word
- Czech example:

Words:	Pro někoho by její provedení mělo smysl .
Lemmas:	pro někdo být jeho provedení mít smysl .
Lemmas+Pseudowords:	pro někdo být PER_3 jeho provedení mít PER_X smysl .
Modified Lemmas:	pro někdo být+PER_3 jeho provedení mít+PER_X smysl .

Figure 2: Various transformations of the Czech sentence from Figure 1. The pseudowords and modified lemmas encode the verb person feature, with the values 3 (third person) and X (“any” person).

- Stemming / truncation isn't as effective as a true lemmatizer, but it's better than nothing
- Czech-English results:

	Dev	Test
word-to-word	.311	.270
lemmatize all	.355	.299
except Pro	.350	
except Pro, V, N	.346	
lemmatize $n < 50$	.370	.306
truncate all	.353	.283

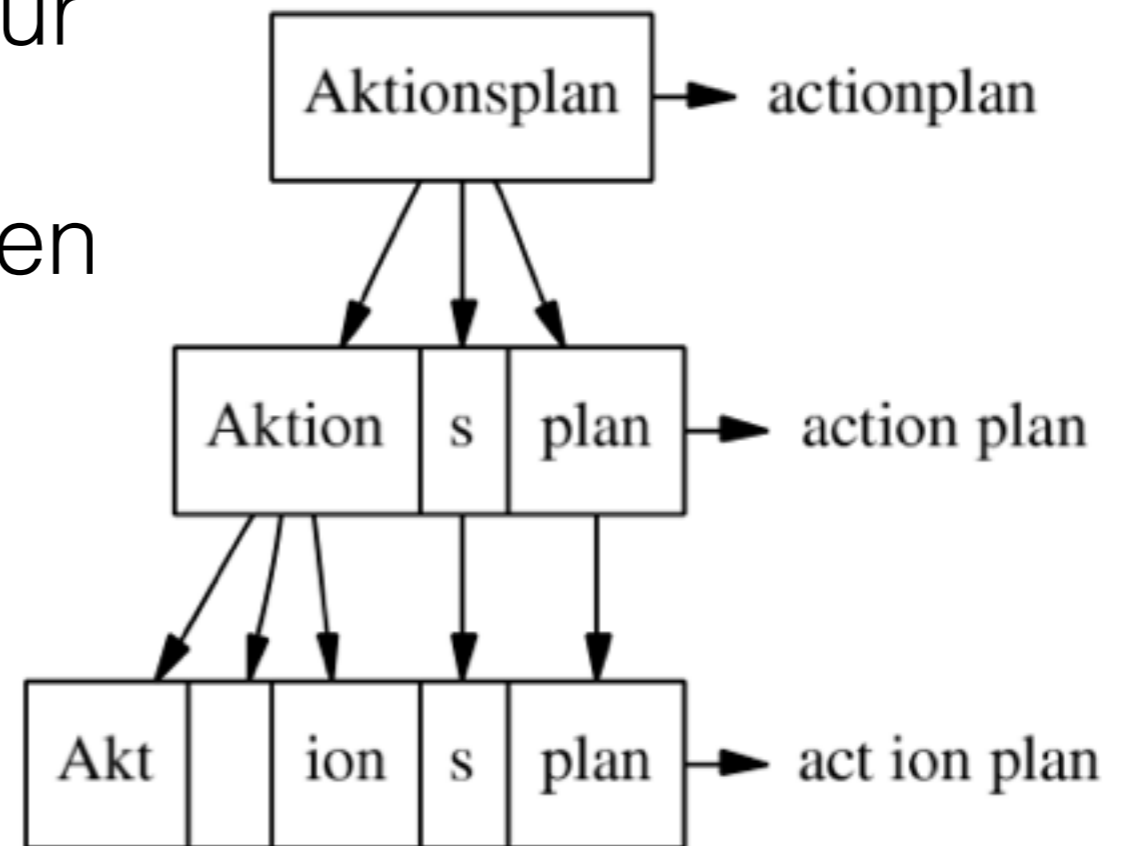
Table 1: BLEU scores for the word-to-word baseline, lemmatization, and word truncation experiments.

- These approaches work for translation from MPLs

# Compound splitting

- German is known for long noun compounds
  - *Großeltern* (grandparents)
  - *Waschmaschine* (washing machine)
  - *Museenverwaltung* (museum management)
- Sometimes this is fine, but sometimes this complicates learning word translations

- What if this wasn't seen in our training data? We might be hosed, but maybe we've seen subparts (especially for German)
- An obvious approach is to split up tokens



- Technique 1: break word into parts that occur elsewhere, and heuristically score them to choose the most likely decomposition
  - aktionsplan (count = 852)  $\rightarrow$  852
  - aktion (960) + plan (710)  $\rightarrow$  825.6
  - aktions (5) + plan (710)  $\rightarrow$  59.6
  - akt (224) + ion (1) + plan (710)  $\rightarrow$  54.2
- Problem:
  - frei (885) + tag (1864)  $>$  freitag (556)

- Technique 2: make sure parts have translations on the English side
  - since *Frei* (*free*) and *Tag* (*day*) are unlikely to exist in the translation of the sentence, *Freitag* (*Friday*) would not be split
- Problem: ambiguity (the word translations might not always appear)
  - *Grundrechte* (*basic rights*)  
*Grund* (*reason/foundation*) + *rechte* (*rights*)



- Technique 3: create a separate translation table from the Method 1 technique, use that as a second-level check
- Further issue: common words result in splits
  - *folgenden (following)*  
*folgen (consequences) + den (the)*
  - solution: POS tag German, limit splitting to certain classes

# Results

- BLEU score: 30.5 (raw), 34.4 (best splitting)
- Lessons
  - Heuristic splitting is complicated and messy: a cascade of exceptions
  - These approaches are also largely specific to German (assuming a particular kind of morphology, and requiring a tagger, for example)

# Translation from lattices

- In the German-English example, we chose a split for the words prior to learning phrase tables and to translation
- This can be problematic if the segmentation had mistakes
- Idea: preserve the ambiguity of splitting and let the decoder efficiently explore *all* splits

- Can be applied to segmentation in non-morphological settings, as well

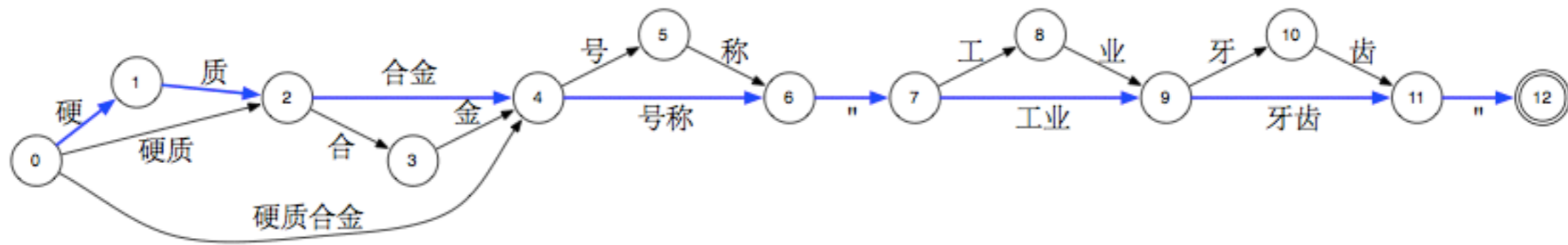


Figure 5: Sample Chinese segmentation lattice using three segmentations.

- Czech–English task
  - By themselves, lemmatization and stemming were not especially helpful
  - A backoff model (in which lower-order models are consulted only when needed) showed some improvement
  - The best model made use of a lemmatized confusion network

Input	BLEU	Sample translation
SURFACE	22.74	From the <b>US side</b> of the Atlantic all such <b>odůvodnění</b> appears to be a totally bizarre.
LEMMA	22.50	From the <b>side</b> of the Atlantic <b>with</b> any such <b>justification</b> seem completely bizarre.
TRUNC ( $l=6$ )	22.07	From the <b>bank</b> of the Atlantic, all such <b>justification</b> appears to be totally bizarre.
backoff (SURFACE+LEMMA)	23.94	From the <b>US bank</b> of the Atlantic, all such <b>justification</b> appears to be totally bizarre.
<b>CN (SURFACE+LEMMA)</b>	<b>25.01</b>	From the <b>US side</b> of the Atlantic all such <b>justification</b> appears to be a totally bizarre.
CN (SURFACE+TRUNC)	23.57	From the <b>US</b> Atlantic any such <b>justification</b> appears to be a totally bizarre.

# Translation into MRLs

- The hard direction
  - Sometimes represents an *increase* in information
- Common approaches
  - Translate stems, inflect afterwards
  - Factored translation

# Translate 'n' inflect

- Apply a lemmatizer or stemmer to target-side training data
- Translate the reduced forms, then apply inflection to 1-best or n-best list
- Reduced setting: just inflect lemmatized words (HW5)
  - Done for Arabic and Hebrew
  - “Generating Complex Morphology for Machine Translation” (Minkov, Toutanova, & Suzuki, 2007)
  - Lots of feature ideas!

- Can also be incorporated into the decoder
- Translate stems (or fully inflected forms)
- Built an inflection sequence model applied in three different ways

*“Applying Morphology Generation Models to Machine Translation”  
(Toutanova et al., 2008)*



- Method 1: Fully-inflected translation model, reinfect output
- Method 2: Align fully inflected forms, stem target side, build system
- Method 3: Stem target side, align, build system

Model	BLEU	Oracle BLEU
Base MT ( $n=1$ )	29.24	-
Method 1 ( $n=1$ )	30.44	36.59
Method 1 ( $n=16$ )	30.61	45.33
Method 2 ( $n=1$ )	30.79	37.38
Method 2 ( $n=16$ )	31.24	48.48
Method 3 ( $n=1$ )	31.42	38.06
Method 3 ( $n=32$ )	31.80	49.19

Table 3: Test set performance for English-to-Russian MT (BLEU) results by model using a treelet MT system.

Model	BLEU	Oracle BLEU
Base MT ( $n=1$ )	35.54	-
Method 1 ( $n=1$ )	37.24	42.29
Method 1 ( $n=2$ )	37.41	52.21
Method 2 ( $n=1$ )	36.53	42.46
Method 2 ( $n=4$ )	36.72	54.74
Method 3 ( $n=1$ )	36.87	42.96
Method 3 ( $n=2$ )	36.92	54.90

Table 5: Test set performance for English-to-Arabic MT (BLEU) results by model using a treelet MT system.

# Factored Translation

- Standard phrase-based model: translate sequences of whitespace-delimited tokens
- An alternative is [factored translation](#) (Koehn & Hoang, 2007), which simultaneously considers multiple sources of evidence

# Factored translation

- Integrates a more complex representation of words directly into the decoder
- Contrast this with some of the other approaches we have considered

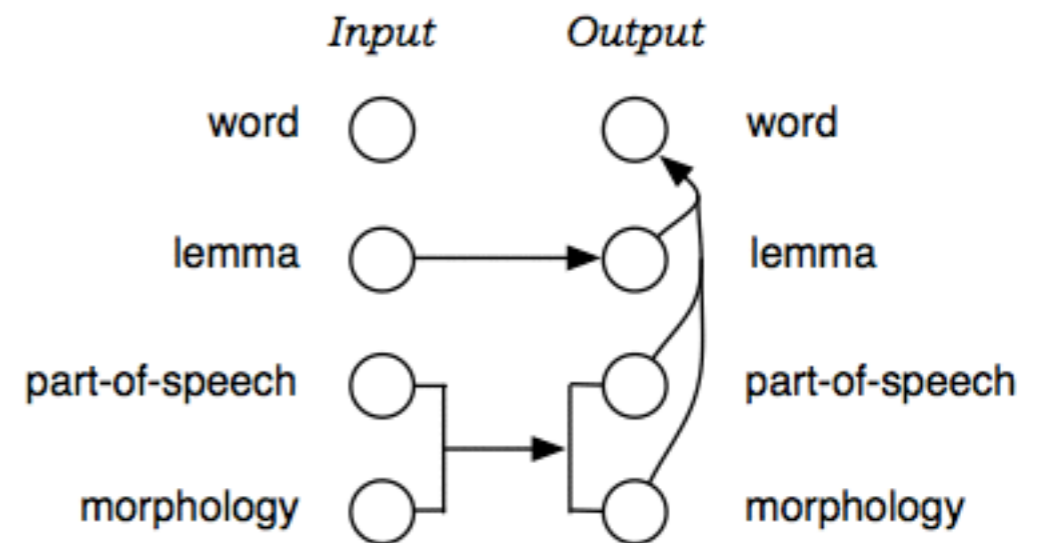


Figure 2: Example factored model: morphological analysis and generation, decomposed into three mapping steps (translation of lemmas, translation of part-of-speech and morphological information, generation of surface forms).

# Factored translation

- Steps
  - **Translate** input factors (phrases) into output factors
  - **Generate** surface forms from the output factors (words)
- Example from paper {surface form | lemma | POS | infl}
  - Map lemma {häuser | haus | NN | pl-nom-neut}  
{*?*|*house*|*?*|*?*, *?*|*home*|*?*|*?*, *?*|*building*|*?*|*?*}
  - Map morphology  
{*?*|*house*|NN|pl, *?*|*home*|NN|pl,  
*?*|*building*|NN|pl, *?*|*house*|NN|sg}
  - Generate surface  
{*houses*|*house*|NN|pl, *homes*|*home*|NN|pl,  
*buildings*|*building*|NN|pl, *house*|*house*|NN|sg}

### English–German

<b>Model</b>	<b>BLEU</b>
best published result	18.15%
baseline (surface)	18.04%
surface + POS	18.15%
surface + POS + morph	18.22%

### English–Spanish

<b>Model</b>	<b>BLEU</b>
baseline (surface)	23.41%
surface + morph	24.66%
surface + POS + morph	24.25%

### English–Czech

<b>Model</b>	<b>BLEU</b>
baseline (surface)	25.82%
surface + all morph	27.04%
surface + case/number/gender	27.45%
surface + CNG/verb/prepositions	27.62%

# Summary

- Morphology is a real problem in translation, especially for low-resource languages
- Delving into morphology requires us to abandon straightforward “language agnostic” approaches
- Linguistics are useful (e.g., lemmatization), and even linguistic approximations (e.g., truncating) can do well
- Morphology is far from a solved problem

# References

- *Empirical Methods for Compound Splitting* (Koehn & Knight, 2005)
- *The ‘noisier channel’: translation from morphologically complex languages* (Dyer, WMT 2007)
- *Factored Translation Models* (Koehn & Hoang, 2007)
- *Improving Statistical MT through Morphological Analysis* (Goldwater & McClosky, 2005)