

# Natural language processing, morphology, and finite-state transducers

Informatics 2A: Lecture 15

Adam Lopez

School of Informatics  
University of Edinburgh  
alopez@inf.ed.ac.uk

21 October 2016

- 1 Natural language processing
- 2 Morphology
- 3 Finite-state transducers
- 4 FSTs for morphology parsing and generation

This lecture is (mostly) based on Jurafsky & Martin chapter 3, sections 1–7.

# Natural language processing



# Formal languages are like this



(image by flickr user Ref54)

# Natural languages are like this



(image by flickr user seliaymiwell)

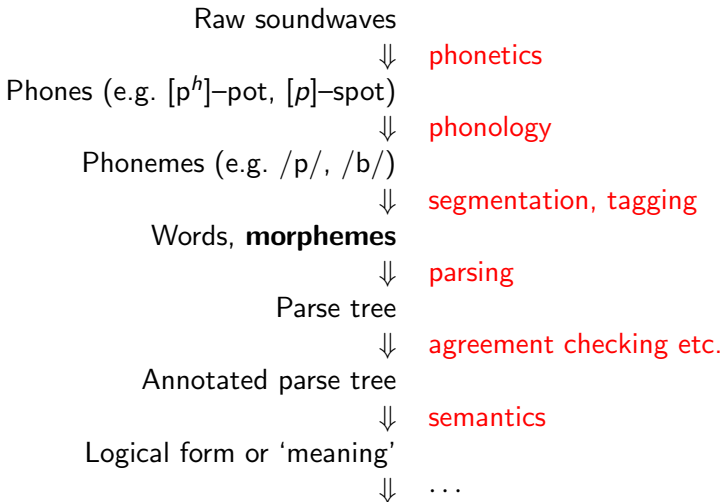
	<i>Exact sciences</i>	<i>Empirical sciences</i>	<i>Engineering</i>
<i>Deals in...</i>	Axioms and theorems	Facts and theories	Artifacts
<i>Truth is...</i>	Forever	Temporary	It works!
<i>Examples...</i>	Maths, CS theory	Physics, Biology, <b>Linguistics</b>	Many, inc. Applied CS and <b>NLP</b>

*Essentially, all models are wrong, but some are useful.*

— *George Box*

# The natural language pipeline

For spoken language:



# What are morphemes? What is morphology?

**Morphology** is the study of the **structure of words**.

A **morpheme** is the smallest meaningful unit of language.

*Morpheme+s are the small+est mean+ing+ful unit+s of language.*

... So, to build systems that understand the meaning of a language, we must build systems that can identify its morphemes, and understand how these morphemes are *composed* to form complete words.



# Why bother?

- Any NLP tasks involving **grammatical parsing** will typically involve morphological parsing as a prerequisite.
- **Search engines**: e.g. a search for 'fox' should return documents containing 'foxes', and vice versa.
- Even a humble task like **spell checking** can benefit: e.g. is 'walking' a possible word form?

But why not just list all derived forms separately in our wordlist (e.g. walk, walks, walked, walking)?

How hard is this problem?

## Example: English morphology

English is primarily **concatenative**: stems and affixes (prefixes, suffixes, infixes and circumfixes) are concatenated to form words.

Four concatenative processes:

- Inflection (stem + grammar affix): word does not change its grammatical class (*walk* → *walking*)
- Derivation (stem + grammar affix): word changes its grammatical form (*computerize* → *computerization*)
- Compounding (stems together): *doghouse*
- Cliticization: *I've*, *we're*, *he's*

These processes are **productive**: *fleecking*, *baeless*, *selfies*

n.b. Morphology can be non-concatenative (e.g. templatic morphology as in Arabic, reduplication in many languages). English is easy!

# Examples: Inflection in different languages

In English, nouns are inflected for **number**, while verbs are inflected for **person** and **tense**.

- Number: *book / books*
- Person: *you read / she reads*
- Tense: *we talk / we talked*

In German, nouns are inflected for number and **case**, e.g. *house* becomes:

	<i>Singular</i>	<i>Plural</i>
<i>Nominative</i>	das Haus	die Häuser
<i>Genitive</i>	des Hauses	der Häuser
<i>Dative</i>	dem Haus / dem Hause	den Häusern
<i>Accusative</i>	das Haus	die Häuser

In Spanish, inflection depends on **gender**, e.g. *el sol / la luna*.

In Luganda, nouns have ten genders, each with different inflections for singular/ plural!

## Examples: Agglutination and compounding

ostoskeskuksessa  
ostos#keskus+N+Sg+Loc:in  
shopping#center+N+Sg+Loc:in  
'in the shopping center' (Finnish)

qangatasuukkuvimmuuriaqalaaqtunga  
'I'll have to go to the airport' (Inuktitut)

Avrupallatramadklarmzdanmsnzcasna  
'as if you are reportedly of those of ours that we were unable to  
Europeanize' (Turkish)

In the most extreme examples, the meaning of the word is the meaning of a sentence!

# English morphological parsing: the problem

Simplest case: assume words consist of a **stem** (carrying the basic dictionary meaning) plus one or more **affixes** carrying grammatical information. E.g.:

<b>Surface form:</b>	cats	walking	smoothest
<b>Lexical form:</b>	cat+N+PL	walk+V+PresPart	smooth+Adj+Sup

**Morphological parsing** is the problem of extracting the lexical form from the surface form. (For speech processing, this includes identifying the word boundaries.)

We should take account of:

- Irregular forms (e.g. goose → geese)
- Systematic rules (e.g. 'e' inserted before suffix 's' after s,x,z,ch,sh: fox → foxes, watch → watches)

# How expressive is morphology?

Morphemes are concatenated in a rather “regular” way. Their order is highly prescribed, even in morphologically rich languages.

This means that finite-state machines are a good way to model morphology. There is no need for “unbounded memory” to model it (there are no long range dependencies).

This is as opposed to syntax, the study of the order of words in a sentence, which we will learn about in about a week.

# Parsing and generation

**Parsing** here means going from the surface to the lexical form.  
E.g. foxes  $\rightarrow$  fox +N +PL.

**Generation** is the opposite process: fox +N +PL  $\rightarrow$  foxes. It's helpful to consider these two processes together.

Either way, it's often useful to proceed via an intermediate form, corresponding to an analysis in terms of **morphemes** (= minimal meaningful units) before **orthographic rules** are applied.

Surface form:	foxes
Intermediate form:	fox ^ s #
Lexical form:	fox +N +PL

(^ means morpheme boundary, # means word boundary.)

N.B. The translation between surface and intermediate form is exactly the same if 'foxes' is a 3rd person singular verb!

We require a relation that maps surface forms to intermediate forms, and another relation that maps intermediate forms to lexical forms, i.e. relations on sets of strings.

How do we represent a (possibly infinite) set of strings?

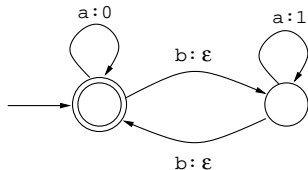
How do we represent a **relation** on two (possibly infinite) sets of strings?



# Finite-state transducers

We can consider  $\epsilon$ -NFAs (over an alphabet  $\Sigma$ ) in which transitions may also (optionally) produce *output* symbols (over a possibly different alphabet  $\Pi$ ).

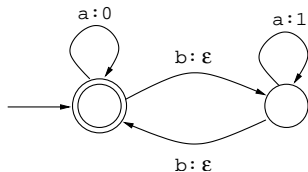
E.g. consider the following machine with input alphabet  $\{a, b\}$  and output alphabet  $\{0, 1\}$ :



Such a thing is called a **finite state transducer**.

In effect, it specifies a (possibly multi-valued) translation from one regular language to another.

## Quick exercise



What output will this produce, given the input *abaaabbab*?

- ① 001110
- ② 001111
- ③ 0011101
- ④ More than one output is possible.

Formally, a **finite state transducer**  $T$  with inputs from  $\Sigma$  and outputs from  $\Pi$  consists of:

- sets  $Q$ ,  $S$ ,  $F$  as in ordinary NFAs,
- a transition relation  $\Delta \subseteq Q \times (\Sigma \cup \{\epsilon\}) \times (\Pi \cup \{\epsilon\}) \times Q$

From this, one can define a many-step transition relation  $\hat{\Delta} \subseteq Q \times \Sigma^* \times \Pi^* \times Q$ , where  $(q, x, y, q') \in \hat{\Delta}$  means “starting from state  $q$ , the input string  $x$  can be translated into the output string  $y$ , ending up in state  $q'$ .” (Details omitted.)

Note that a finite state transducer can be run in either direction! From  $T$  as above, we can obtain another transducer  $\bar{T}$  just by swapping the roles of inputs and outputs.

Formally, a **finite state transducer**  $T$  with inputs from  $\Sigma$  and outputs from  $\Pi$  consists of:

- sets  $Q$ ,  $S$ ,  $F$  as in ordinary NFAs,
- a transition relation  $\Delta \subseteq Q \times (\Sigma \cup \{\epsilon\}) \times (\Pi \cup \{\epsilon\}) \times Q$

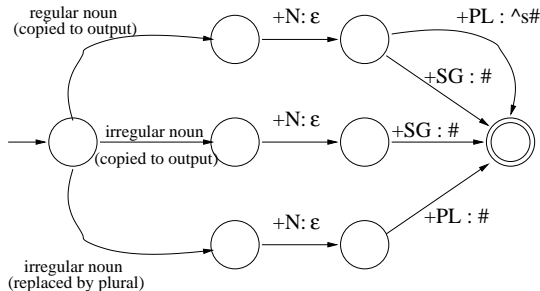
**Reminder:** Formally, an NFA with alphabet  $\Sigma$  consists of:

- A finite set  $Q$  of states.
- A transition relation  $\Delta \subseteq Q \times (\Sigma \cup \{\epsilon\}) \times Q$ ,
- A set  $S \subseteq Q$  of possible starting states,
- A set  $F \subseteq Q$  of accepting states.

## Stage 1: From lexical to intermediate form

Consider the problem of translating a lexical form like 'fox+N+PL' into an intermediate form like 'fox ^ s # ', taking account of irregular forms like goose/geese.

We can do this with a transducer of the following schematic form:

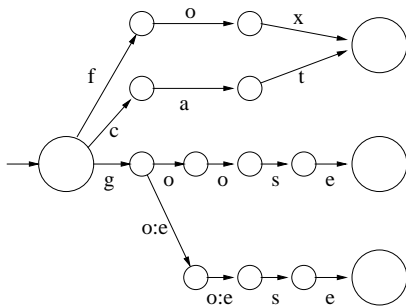


We treat each of +N, +SG, +PL as a single symbol.

The 'transition' labelled +PL : ^s# abbreviates three transitions:  
+PL : ^, ε : s, ε : #.

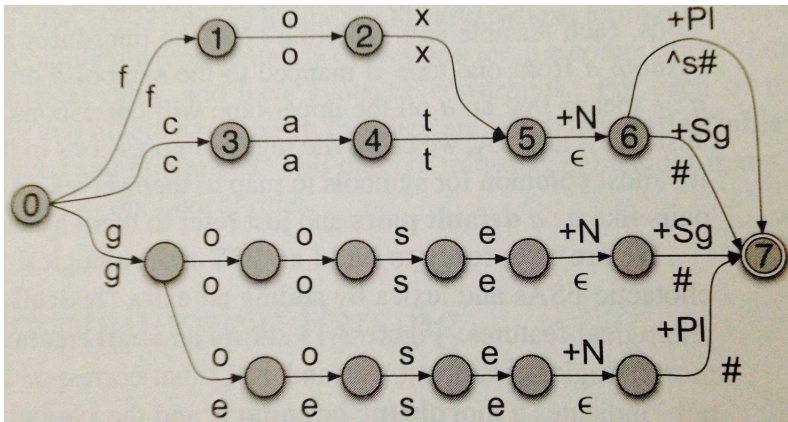
# The Stage 1 transducer fleshed out

The left hand part of the preceding diagram is an abbreviation for something like this (only a small sample shown):



Here, for simplicity, a single label  $u$  abbreviates  $u : u$ .

# Stage 1 in full



## Stage 2: From intermediate to surface form

To convert a sequence of morphemes to surface form, we apply a number of **orthographic rules** such as the following.

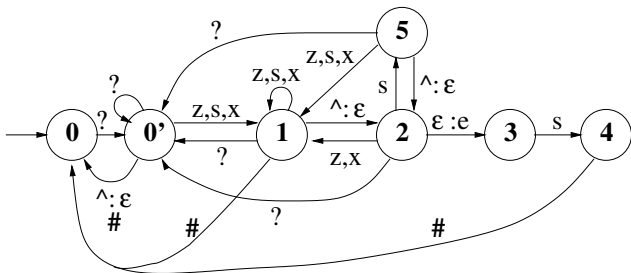
- **E-insertion:** Insert e after s,z,x,ch,sh before a word-final morpheme -s. (fox → foxes)
- **E-deletion:** Delete e before a suffix beginning with e,i. (love → loving)
- **Consonant doubling:** Single consonants b,s,g,k,l,m,n,p,r,s,t,v are doubled before suffix -ed or -ing. (beg → begged)

We shall consider a simplified form of E-insertion, ignoring ch,sh.

(Note that this rule is oblivious to whether -s is a plural noun suffix or a 3rd person verb suffix.)



# A transducer for E-insertion (adapted from J+M)



Here ? may stand for any symbol except  $z,s,x,\hat{\ },\#$ .  
(Treat # as a 'visible space character'.)

At a morpheme boundary following  $z,s,x$ , we arrive in State 2.  
If the ensuing input sequence is  $s\#$ , our only option is to go via states 3 and 4. **Note that there's no #-transition out of State 5.**

State 5 allows e.g. 'ex $\hat{\ }$ service $\hat{\ }$ men $\#$ ' to be translated to 'exservicemen'.

# Putting it all together

FSTs can be **composed**: output from one can be input to another.

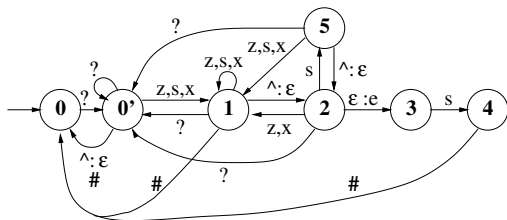
To go from lexical to surface form, use 'Stage 1' transducer followed by a bunch of orthographic rule transducers like the above. (Made more efficient by back-end compilation into one single transducer. To be really efficient we must determinise and minimise the result!)

The results of this **generation** process are typically **deterministic** (each lexical form gives a unique surface form), even though our transducers make use of non-determinism along the way.

Running the same cascade **backwards** lets us do **parsing** (surface to lexical form). Because of ambiguity, this process is frequently **non-deterministic**: e.g. 'foxes' might be analysed as fox+N+PL or fox+V+Pres+3SG.

Such ambiguities are not resolved by morphological parsing itself: left to a later processing stage.

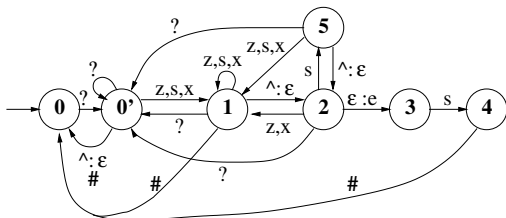
## Quick exercise 2



Apply this **backwards** to translate from surface to int. form.

Starting from state 0, how many **sequences of transitions** are compatible with the input string 'asses' ?

- ① 1
- ② 2
- ③ 3
- ④ 4
- ⑤ More than 4



On the input string 'asses', 10 transition sequences are possible!

- $0 \xrightarrow{a} 0' \xrightarrow{s} 1 \xrightarrow{s} 1 \xrightarrow{\epsilon} 2 \xrightarrow{e} 3 \xrightarrow{s} 4$ , output  $ass^{\wedge}s$
- $0 \xrightarrow{a} 0' \xrightarrow{s} 1 \xrightarrow{s} 1 \xrightarrow{\epsilon} 2 \xrightarrow{e} 0' \xrightarrow{s} 1$ , output  $ass^{\wedge}es$
- $0 \xrightarrow{a} 0' \xrightarrow{s} 1 \xrightarrow{s} 1 \xrightarrow{e} 0' \xrightarrow{s} 1$ , output  $asses$
- $0 \xrightarrow{a} 0' \xrightarrow{s} 1 \xrightarrow{\epsilon} 2 \xrightarrow{s} 5 \xrightarrow{\epsilon} 2 \xrightarrow{e} 3 \xrightarrow{s} 4$ , output  $as^{\wedge}s^{\wedge}s$
- $0 \xrightarrow{a} 0' \xrightarrow{s} 1 \xrightarrow{\epsilon} 2 \xrightarrow{s} 5 \xrightarrow{\epsilon} 2 \xrightarrow{e} 0' \xrightarrow{s} 1$ , output  $as^{\wedge}s^{\wedge}es$
- $0 \xrightarrow{a} 0' \xrightarrow{s} 1 \xrightarrow{\epsilon} 2 \xrightarrow{s} 5 \xrightarrow{e} 0' \xrightarrow{s} 1$ , output  $as^{\wedge}ses$
- Four of these can also be followed by  $1 \xrightarrow{\epsilon} 2$  (output  $\wedge$ ).

Lexicon can be quite large with finite state transducers

Sometimes need to extract the stem in a very efficient fashion  
(such as in IR)

The Porter stemmer: a lexicon-free method for getting the stem of  
a given word

ATIONAL → ATE (e.g., relation → relate)

ING →  $\epsilon$  if stem contains a vowel (e.g. motoring → motor)

SSES → SS (e.g., grasses → grass)

Makes errors:

organization → organ

doing → doe

numerical → numerous

policy → police

A vibrant area of study

Mostly done by *learning from data*, just like many other NLP problems. Two main paradigms: unsupervised morphological parsing and supervised one.

NLP solvers are not perfect! They can make mistakes. Sometimes ambiguity can't even be resolved. BUT, for English, morphological analysis is highly accurate. With other languages, there is still a long way to go.

One of the basic tools that is used for many applications

- Speech recognition
- Machine translation
- Part-of-speech tagging
- ... and many more

## Part-of-speech tagging:

- What are parts of speech?
- What are they useful for?
- Zipf's law and the ambiguity of POS tagging
- One problem NLP solves really well (... for English)