

Linguistic and Computational Morphology¹

Agata Savary

December 19, 2008

Morphology

Linguistic discipline interested in the internal structure of (written) words.

What is a word ?

In **linguistics** a **word** has two senses:

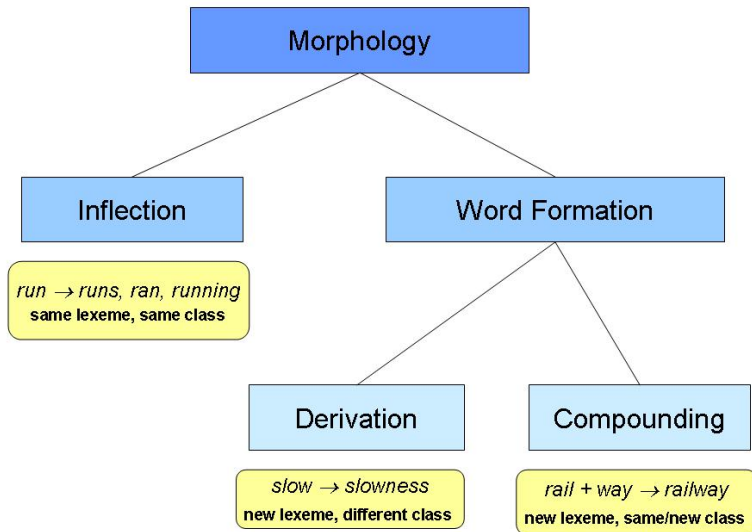
- ▶ **Lexeme** = abstract (“deep”) unit having a certain **meaning**, and belonging to a certain **class**; **lexicon** = set of lexemes
- ▶ **Word form** = different textual (“surface”) realizations of a lexeme

The **inflection paradigm** = all word forms of a lexeme.

A **lemma** = a canonical word form chosen to represent the lexeme.

	Word forms	Lemma
French	{porte, portes}	porte
	{porter, porte, portes, portiez, ... }	porter
	{à}	à
Your language		

Types of morphological rules



Inflectional categories and values (1/4)

	English (Germanic)	French (Latin)	Serbian (Slavic)	Your language (...)
Number (<i>Nb</i>)	singular (<i>s</i>) plural (<i>p</i>)	singular (<i>s</i>) plural (<i>p</i>)	singular (<i>s</i>) plural (<i>p</i>) paukal (<i>w</i>)	
Gender (<i>Gen</i>)		masculine (<i>m</i>) feminine (<i>f</i>)	masculine (<i>m</i>) feminine (<i>f</i>) neuter (<i>n</i>)	
Case			nominative (<i>1</i>) genitive (<i>2</i>) dative (<i>3</i>) accusative (<i>4</i>) instrumental (<i>5</i>) locative (<i>6</i>) vocative (<i>7</i>)	

Inflectional categories and values (2/4)

	English	French	Serbian	Your language
Degree (<i>Deg</i>)	positive (<E> comparative (C) superlative (S)		positive (a) comparative (b) superlative (c)	
Person (<i>Pers</i>)	first (1) second (2) third (3)	first (1) second (2) third (3)	first (x) second (y) third (z)	
Animate-ness (<i>Anim</i>)			animate (v) inanimate (q) no-care (g)	

Inflectional categories and values (3/4)

	English	French
Tense and mood (TM)	infinitive (<i>W</i>): <i>do</i> present indicative (<i>P</i>): <i>does</i> imperfect indicative (<i>I</i>): <i>did</i> past participle (<i>K</i>): <i>done</i> gerund (<i>G</i>): <i>doing</i>	infinitive (<i>W</i>): <i>faire</i> present indicative (<i>P</i>): <i>faisons</i> imperfect indicative (<i>I</i>): <i>faisait</i> present subjunctive (<i>S</i>): <i>fasse</i> imperfect subjunctive (<i>T</i>): <i>fisse</i> present imperative (<i>Y</i>): <i>faites</i> present conditional (<i>C</i>): <i>ferait</i> simple past (<i>J</i>): <i>fit</i> past participle (<i>K</i>): <i>faite</i> gerund (<i>G</i>): <i>faisant</i> future (<i>F</i>): <i>fera</i>

Inflectional categories and values (4/4)

	Your language
Tense and mood (<i>TM</i>)	

Inflectional classes \approx parts of speech (POS) (1/6)

	Noun	Max. forms
English	<u>inflects in</u> number: <i>dog, dogs</i>	2
French	<u>inflects in</u> number <i>toile, toiles</i> has gender <i>toile</i> OR <u>inflects in</u> gender <i>cousin, cousine</i>	4
Serbian	<u>inflects in</u> number has gender OR <u>inflects in</u> gender <u>inflects in</u> case has animateness	28
Your language		

Inflectional classes \approx parts of speech (POS) (2/6)

	Adjective	Max forms
English	<u>uninflected</u> <i>famous</i> OR <u>inflects in</u> <i>big, bigger</i>	3
French	<u>inflects in</u> <i>bleu, bleus</i> <u>inflects in</u> <i>bleue, bleues</i>	4
Serbian	<u>inflects in</u> number <u>has</u> gender OR <u>inflects in</u> gender <u>inflects in</u> case <u>inflects in</u> animateness <u>inflects in</u> degree <u>inflects in</u> determinedness	77
Your language		

Inflectional classes \approx parts of speech (POS) (3/6)

	Verb	Max. forms
English	<u>inflects in</u> <i>go, went, going</i> <u>inflects in</u> <i>go, goes</i> <u>inflects in</u> <i>am, are</i>	9
French	<u>inflects in</u> <i>être, suis, été</i> <u>inflects in</u> <i>suis, es, est</i> <u>inflects in</u> <i>suis, sommes</i> <u>inflects in</u> <i>aimés, aimées</i>	51
Serbian	<u>inflects in</u> tense-mood <u>inflects in</u> person <u>inflects in</u> number <u>inflects in</u> gender	dozens
Your language		

Inflectional classes \approx parts of speech (POS) (4/6)

	Pronoun	Max forms
English	<u>inflects in</u> <i>I, you, he</i> <u>inflects in</u> <i>I, we</i> <u>inflects in</u> <i>he, she</i>	8
French	<u>inflects in</u> <i>je, tu, il</i> <u>inflects in</u> <i>tu, vous</i> <u>inflects in</u> <i>il, elle</i>	8
Serbian	<u>inflects in</u> person <u>inflects in</u> number <u>inflects in</u> gender ...	10
Your language		

Inflectional classes \approx parts of speech (POS) (5/6)

	Adverb	Max. forms
English	<u>uninflected</u> <i>yesterday</i> OR <u>inflects in</u> <i>early, earlier</i>	3
French	<u>uninflected</u> <i>hier, facilement</i>	1
Serbian	<u>uninflected</u>	1
Your language		

Inflectional classes \approx parts of speech (POS) (6/7)

	Determiner	Max. forms
English	<u>has</u> <i>a, this, those, the</i>	1
French	<u>inflects in</u> <i>le, les</i>	4
	<u>inflects in</u> <i>le, la</i>	
Serbian	<u>inexistent</u>	0
Your language		

Inflectional classes \approx parts of speech (POS) (6/6)

	Preposition	Conjunction	Interjection
English	<u>uninflected</u> : <i>to</i>	<u>uninflected</u> : <i>and</i>	<u>uninflected</u> : <i>hurray</i>
French	<u>uninflected</u> : <i>de</i>	<u>uninflected</u> : <i>mais</i>	<u>uninflected</u> : <i>adieu</i>
Serbian	<u>uninflected</u>	<u>uninflected</u>	<u>uninflected</u>
Your language			

Inflectional paradigm (verb lemma *aimer*)

Word form	Features	Word form	Features	Word form	Features
aimer	W	aimais	I2s	aimais	I1s
aimait	I3s	aimions	I1p	aimiez	I2p
aimaient	I3p	aimassent	T3p	aimassiez	T2p
aimassions	T1p	aimât	T3s	aimasses	T2s
aimasse	T1s	aimai	J1s	aima	J3s
aimâmes	J1p	aimâtes	J2p	aimèrent	J3p
aimas	J2s	aimant	G	aimés	Kmp
aimé	Kms	aimées	Kfp	aimée	Kfs
aimons	Y1p	aimons	P1p	aimions	S1p
aimiez	S2p	aimerais	C2s	aimerais	C1s
aimerait	C3s	aimerions	C1p	aimeriez	C2p
aimeraient	C3p	aimerai	F1s	aimeras	F2s
aimera	F3s	aimerons	F1p	aimerez	F2p
aimeront	F3p	aime	Y2s	aime	S3s
aime	S1s	aime	P3s	aime	P1s
aiment	S3p	aiment	P3p	aises	S2s
aises	P2s	aimez	Y2p	aisez	P2p

Derivational morphology (1/2)

- ▶ Source word:
 - a lemma: *small* → *smallness*
 - an inflected form (in French): *normale* → *normalement*
- ▶ Derivational affix:
 - prefix: *ir* + *regular*(adj.) → *irregular*
 - infix: e.g. in Arabic
 - suffix: *small*+*ness* → *smallness*
 - no affix: *to enter* → *an enter*
- ▶ Target word: different lexeme and/or different class (inflects differently)
 - small*(adj.) → *smallness*(noun)
 - astonish*(verb) → *astonishment*(noun)
 - count*(verb) → *countable*(adj.)
 - courage*(noun) → *encourage*(verb)
 - forest*(noun) → *forestry*(noun)
- ▶ Stem modification:
 - regulate* → *regulation*
- ▶ Multiple affixes:
 - un* + *forget* + *able* → *unforgettable*

Compounding

- ▶ Several lexemes form a **new lexeme**.
- ▶ The new lexeme shows some degree of **non-compositionality**
 - ▶ morphological: *un peau rouge*(masc.), unlike *peau*(fem.)
 - ▶ syntactic: *un moulin à vent*, but not **un moulin à brise*
 - ▶ distributional: *un cordon bleu*(human), unlike *cordon*(inanimate)
 - ▶ semantic: *pomme de terre* is not an apple from earth

Headword

- ▶ **Headword**: component from which the compound inherits its features

fireman - noun in singular like *man*

cheval à bascules - noun in singular masculine, like *cheval*

- ▶ Types of compounds:

- ▶ **endocentric** (has a headword): *fireman*
- ▶ **exocentric** (no headword): (EN) *forget-me-not*, (FR) *porte-serviettes*
- ▶ **apposition** (two heads): *man servant* → *men-servants*

Examples of compounds

	Noun	Adjective	Verb
English	<i>air brake</i> <i>forget-me-not</i> <i>man-of-war</i>	<i>bittersweet</i> <i>easy-going</i> <i>as busy as a bee</i>	<i>cut off</i> <i>co-occur</i> <i>make up for</i>
French	<i>rouge-gorge</i> <i>stylo à bille</i> <i>porte-monnaie</i>	<i>à pied</i> <i>anglo-saxon</i> <i>sans domicile fixe</i>	<i>sous-entendre</i> <i>faire avec</i> <i>contre-attaquer</i>
Your language			

Examples of compounds

	Adverb	Preposition	Conjunction
English	<i>all of a sudden</i> <i>as soon as possible</i> <i>on and on</i>	<i>instead of</i> <i>contrary to</i> <i>in front of</i>	<i>as well as</i> <i>if and only if</i> <i>neither . . . nor</i>
French	<i>trop bien</i> <i>un peu</i> <i>à l'envers</i>	<i>à propos de</i> <i>de façon à</i> <i>en cas de</i>	<i>alors que</i> <i>parce que</i> <i>au moment où</i>
Your language			

Ambiguity of compounds

- ▶ Non-ambiguous compound: each occurrence of its components is always a compound.

Je suis venu parce que je le voulais.

- ▶ Ambiguous compound: an occurrence of its components may or may not be a compound.

Je suis venu alors que je ne le voulais pas.

Il m'a dit alors que l'affaire était close.

Natural language \neq formal language

- ▶ Linguistic definitions are circular (.....)
- ▶ Basic elements are not clearly defined (.....)
- ▶ Many notions are based on human intuition, and remain formally undescribed

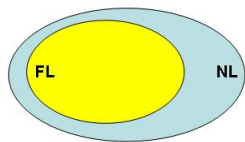
But:

- ▶ Computer programs cannot deal with implicit knowledge
- ▶ They can only treat formal languages

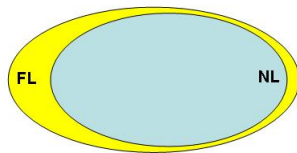
Solution:

- ▶ Define a formal language as close as possible to the natural language

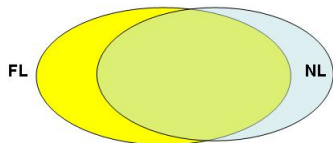
Natural language vs. formal language



Overlooking of exceptions



Overgeneralization



Both overgeneralization and overlooking of exceptions

What is a word ? What is an alphabet ?

In a **formal language**:

- ▶ **alphabet** Σ is a finite set of symbols
- ▶ a **word** over Σ is a (finite or infinite) sequence of elements in Σ : $\omega \in \Sigma^*$
- ▶ a **language** is a (finite or infinite) subset of Σ^* given by a **grammar**

Example:

- ▶ $\Sigma = \{a, b\}$
- ▶ $L = \{a, aba, aabaa, aaabaaa, \dots\}$
- ▶ Grammar = ...

What is a word ? What is an alphabet ?

In a **natural language** - on the **morphological level**

- ▶ an **alphabet** = list of (lowercase, uppercase, accented,...) letters of the language

In English: {A, a, B, b, C, c, ...}

In French: {A, a, Â, â, À, à, B, b, C, c, ...}

In your language:

- ▶ a **language** = list of all correct (*grammatical*) words of the language

In English: {a, the, dog, dogs, make, making, example, ...}

In French: {un, à, cher, chères, exemple ...}

In your language:

- ▶ a **grammar**

- ▶ A set of correct word forms
- ▶ Grammar rules (over)generating sets of words

In English: NOUN → ADJ *ness*

*madness, emptiness, *irregularness, ...*

What is a word ? What is an alphabet ?

In a **natural language** - on the **syntactic level**

- ▶ an **alphabet** = list of valid *morphological words*

In English: {a, the, dog, dogs, make, making, example, ... }

In French: {un, à, cher, chères, exemple ... }

In your language:

- ▶ a **language** = list of all correct (*grammatical*) sentences of the language:

In English: {Dogs like cats., Do cats like dogs?, We will see example 5., ... }

In French: {Ces maisons sont-elles chères?, Tais-toi!, ... }

In your language:

- ▶ a **grammar**

Many formalisms were proposed (DCG, TAG, ..., see lecture on syntax)

An **complete** and **efficient** grammar remains a challenge

Non-alphabet characters

- ▶ They help to separate morphological words in a sentence
- ▶ They separate sentences
- ▶ They may be parts of words (*aujourd'hui*)
- ▶ They may miss between words (*Schul/errinnerung*)
- ▶ They may have a semantic content : λ -calculus, γ -rays

The English paradox

- ▶ It is the dominating language in the Natural Language Processing (NLP) community
- ▶ It is one of the least inflected occidental languages

Computational morphology

- ▶ **Tokenization** = dividing text into elementary graphical units (word forms, separators, ...)
- ▶ **Morphological analysis** = assigning all possible morphological interpretations to a word form (out of context)
- ▶ **Morphological disambiguation (tagging)** = choosing the correct interpretation in the given context
- ▶ **Morphological generation** = for a given lemma and annotation, produce the corresponding word form(s)

Morphological analysis and generation

Morphological analysis: from “surface” form to an (several) annotation(s)

avions → {⟨ Lemma=**avion**, Class=**N**, Nb=**p**⟩,
⟨ Lemma=**avoir**, Class=**V**, Nb=**p**, TM=**I**, Pers=**1**⟩}

Morphological generation: from an annotation to a (several) surface forms

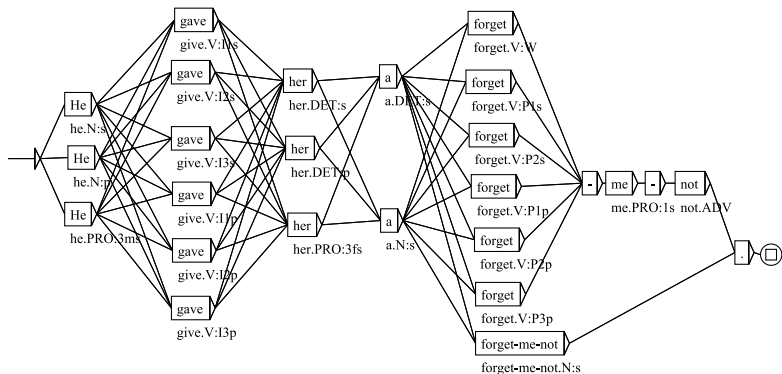
⟨ Lemma=**avoir**, Class=**V**, Nb=**p**, TM=**I**, Pers=**1**⟩ → **avions**

Tokenization and morphological analysis of a sentence

He gave her a forget-me-not.

He	gave	her	a	forget	-	me	-	not
he.N:s	give.V:l1s	her.DET:s	a.DET:s	forget.V:P1s	-	me.PRO:1s	-	not,.A
he.N:p	give.V:l2s	her.DET:p		forget.V:P2s	-		-	
he.PRO:3ms	give.V:l3s	her.PRO:3fs		forget.V:P1p	-		-	
	give.V:l1p			forget.V:P2p	-		-	
	give.V:l2p			forget.V:P3p	-		-	
	give.V:l3p			forget.V:W	-		-	
				forget-me-not.N:s				

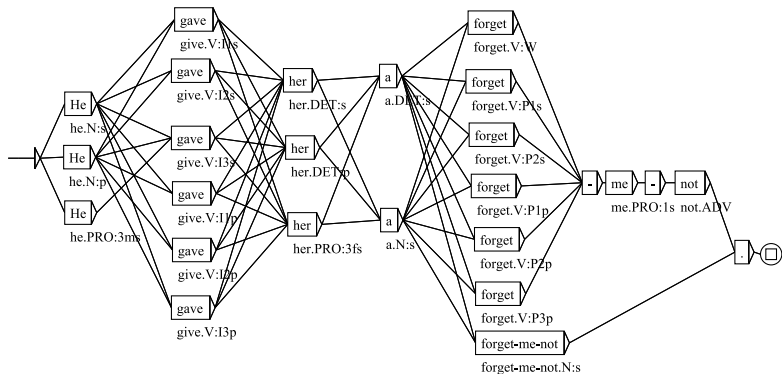
A sentence becomes a graph



How many possible interpretations of the sentence ?

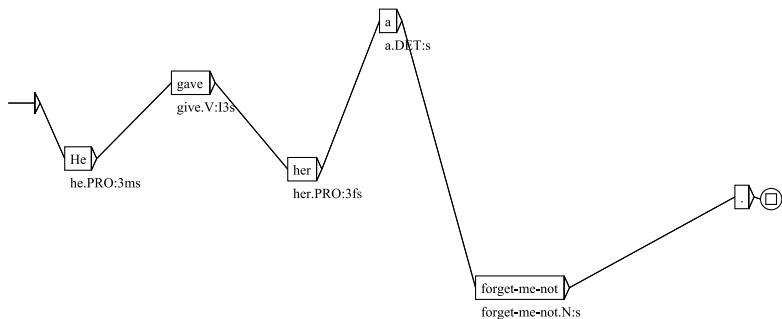
Disambiguation: cutting off forbidden paths

Disambiguating rule - example: *If a personal pronoun is followed by a verb, both must agree in number and person.*



How many possible interpretations of the sentence were eliminated?
?

Tagging = choosing the correct interpretation of the sentence



A perfect tagging is not always possible

Truly ambiguous sentences exist:

La petite brise la glace.

