

Computational Linguistics I: Introduction and Machine Translation

Linguistics 201
Spring 2004

Overview

1. What is it?
2. Where is it used?
3. Machine Translation

2

What is it?

Computational linguistics – study of how to process natural language with computers.

Purpose:

- practical: computers can help us with many tasks involving language
- theoretical: linguistics can test its theories

Methods used and developed in CL are used in other areas, too (DNA decoding, recognition of faces on photographs, etc.)

3

Where is it used? Applications

Some of the following applications are already available (to a certain extent) some are still waiting to be made real:

- Machine translation from one language to another
 - MT can be fully or partly automatic (assisted, i.e. Computer assisting you or you assisting computer, with pre-emption or post-emption)
 - important especially for multilingual countries (Canada, India, Switzerland), international institutions (UN, IMF), multinational companies, exporters
- For example, Canada Meteo – translate Eng weather reports into French
- The European Union used to have 11 official languages, since May 1 it has 20. All federal laws and other documents have to be translated into all languages.

4

- Searching the internet or a database for relevant documents
Google, a library system, searching for law precedents
- Spell-checking & Grammar checking
e.g., **goed, *sentance, *speach; *this books, *a man are; to, too, two*
- Discovering plagiarism
 - analyzing text to find whether it is composed from blocks written by others,
 - automatically check whether one text is not a rewritten version of another text
- Text summarization
e.g., creating a 3 page summary of a 1000 page book
- Dictation system, automatic closed captioning

5

- Beeping out four letter-words
- Reading a written text aloud, e.g., for blind people
- Automatic customer service via phone or web in natural languages (instead of infinite menus)
e.g., You would say: *I need some help with the form I-765.*
- Adding diacritics to a text without it
e.g., *Munchen* → *München*, *Ceske Budejovice* → *Ceské Budejovice*
- Getting specific information from many sources
e.g., *Give me a one-page summary of what do the media write about Iraq.*
- Control of machinery by voice.

6

- Machines reporting their status by voice.
- Simultaneous generation of texts in several languages
e.g., User guides for multinational companies
- Written text recognition (optical character recognition, OCR).
- Searching documents for a word regardless of its morphology
e.g., A Search for *go* returns sentences with *go, goes, going, gone, went.*
More important and complicated for languages with complex morphology
e.g., In Czech *jít (go)* returns *jdu, jdeme, el, la, jdoucí, priel, odejdi, . . .*
- etc.

7

Components used in such applications

- Speech recognition/synthesis
- Morphological analysis
was is a past tense of *be*, 1st or 3rd person sg.
- Tagging – determining the word classes of words
What are the word classes in *Can he can me for kicking a can?*
- Understanding dates:
12th of April 2002; April 12, 2002; 04/12/2002; 04/12/02; 04/12; April 12; 2002-04-12

8

- Parsing – determining the syntactic structure(s) of a sentence

- Word sense disambiguation

What does *pen* mean in *Put some ink in the pen*, and what in *Put the pig in the pen*?

- What does a pronoun refer to.

What does *they* refer to in the following sentences:

The officials forbade the celebrations, because they were afraid of riots.

The officials forbade the celebrations, because they tend to be violent.

- Determining language of a text (after that you can run the appropriate spelling/grammar checker, translator, etc.)
- etc.

9

Taking a closer look at machine translation

Translation is the process of:

- moving texts from one (human) language (**source language**) to another (**target language**),
- in a way that preserves meaning.

Machine translation (MT) automates the process, or part of the process.

- Fully automatic translation
- Computer-aided (human) translation

10

What is MT good for?

- When you need the gist of something and there are no human translators around:
 - translating e-mails & webpages
 - obtaining information from sources in multiple languages (e.g., search engines)
- If you have a limited vocabulary and a small range of sentence types:
 - translating weather reports
 - translating technical manuals
 - translating terms in scientific meetings
 - determining if certain words or ideas appear in suspected terrorist documents
 - help pin down which documents need to be looked at closely
- If you want your human translators to focus on interesting/difficult sentences while avoiding lookup of unknown words and translation of mundane sentences.

11

What is MT not good for?

- Things that require subtle knowledge of the world and/or a high degree of (literary) skill:
 - translating Shakespeare into Navaho
 - diplomatic negotiations
 - court proceedings
 - actually knowing what the terrorists are going to do (chances are, their messages are rather cryptic)
- Things that may be a life or death situation:
 - Pharmaceutical business
 - Automatically translating frantic 911 calls for a dispatcher who speaks only Spanish

12

Translation examples

It will help to look at a few examples of real translation before talking about how a machine does it.

Spanish and English:

- (1) *Yo hablo español.*
I speak_{1st,sg} Spanish
'I speak Spanish.'

⇒ Words pretty much translate one-for-one, but we have to make sure *hablo* matches with *Yo*. (verb *conjugation*)

13

Translation examples (cont.)

- (2) a. *Tu hablas español?*
You speak_{2nd,sg} Spanish
'Do you speak Spanish?'
- b. *Hablas español?*
Speak_{2nd,sg} Spanish
'Do you speak Spanish?'

⇒ Now there's either a word missing ('you') or in a different order in the Spanish.

14

More translation examples

Russian and English:

- (3) *Ya vac lyublyu.*
I_{nominative} you_{accusative} love.
'I love you.'

⇒ Word order in Russian is often different than in English.

- (4) *Ya dal vam podarok.*
I gave you_{dative} gift
'I gave a gift to you.'

⇒ The word for 'you' changes depending on how it's used in the sentence.

15

What goes into a translation

Some things to note about these examples and thus what we might need to know to translate:

- Languages are rule-based (we saw this with syntax).
- Words have to be translated. (Dictionary)
- Word order can change.
- Verbs have to be conjugated; nouns might look different depending on where they are at in the sentence.

16

Different approaches to MT

- Transformer systems
- Linguistic knowledge systems
 - Direct transfer systems
 - Interlinguas
- Machine learning approaches

All of these (except machine learning) use dictionaries in one form or another, so we will start by looking at dictionaries.

17

Dictionaries

An MT **dictionary** is different than a “paper” dictionary.

- must be computer-usable
 - contain the inherent properties (meaning) of a word
 - need to be able to handle various word inflections
- have* is the dictionary entry, but we want the entry to specify how to conjugate this verb.

18

Dictionaries (cont.)

- contain (syntactic and semantic) restrictions it places on other words
 - e.g., Subcategorization information: *give* needs a giver, a person given to, and an object that is given
 - e.g., Selectional restrictions: if X is *eating*, then X must be animate.
- may also contain frequency information
- can be hierarchically organized: all nouns have person, number, and gender.
 - e.g., Verbs (unless irregular) conjugate in the past tense by adding *ed*.

19

What dictionary entries might look like

WORD: *button*
PART OF SPEECH: noun
HUMAN: no
CONCRETE: yes
GERMAN: Knopf

WORD: *knowledge*
PART OF SPEECH: noun
HUMAN: no
CONCRETE: no
GERMAN: Wissen, Kenntnisse

⇒ Can have separate rules which tell you whether to choose *Wissen* or *Kenntnisse*.

20

A dictionary entry with frequency

WORD: *knowledge*

PART OF SPEECH: noun

HUMAN: no

CONCRETE: no

GERMAN: Wissen: 80%, Kenntnisse: 20%

Probabilities derived from various machine learning techniques → to be discussed later.

21

System 1: Transformers

Transformer architectures transform one language into another. They have:

- a grammar for the source/input language
- a source-to-target language dictionary
- source-to-target language rules

Note that there is no grammar for the target language, only mappings from the source language.

22

Transformer steps

We'll work through a German-to-English example.

- Drehen Sie den Knopf eine Position zurück.
- Turn the button back a position.

1. Using the grammar, give parts of speech to the input words

(5) *Drehen Sie den Knopf eine Position zurück.*
verb pronoun article noun article noun preposition

23

Transformer steps (cont.)

2. Using the grammar, give the sentence a (basic) structure

(6) *Drehen Sie (den Knopf) (eine Position) zurück.*

3. Using the dictionary, find the target language words

(7) *Drehen Sie den Knopf eine Position zurück.*
turn you the button a position back

4. Using the source-to-target rules, reorder/combine/eliminate/add target language words, as needed

- 'turn' and 'back' form one unit.
- Because 'Drehen ... zurück' is a command, 'you' is unnecessary

⇒ End result: *Turn back the button a position.*

24

Transformers: Less than meets the eye

- By their very nature, transformer systems are non-**reversible** because they lack a target language grammar.

If we have a German to English translation system, for example, we are incapable of translating from English to German.

- However, as these systems do not require sophisticated knowledge of the target language, they are usually very **robust** = they will return a result for nearly any input sentence.

25

System 2: Linguistic knowledge architectures

These systems include knowledge of both the source and the target languages. We will look at direct transfer systems and then the more specific instance of interlinguas.

- Direct transfer systems
- Interlinguas

26

Transfer systems

These systems have:

- A source language grammar
- A target language grammar
- Rules relating source language underlying representation to target language underlying representation

27

Transfer systems

A direct transfer system has a **transfer component** which relates a source language representation with a target language representation. This can also be called a **comparative grammar**.

We'll walk through the following French to English example:

- (8) *Londres plaît à Sam.*
London is pleasing to Sam
'Sam like London.'

28

Transfer system steps

1. source language grammar analyzes the input and puts it into an **underlying representation** (UR).

Londres plaît à Sam → Londres plaire Sam (source UR)

2. The transfer component relates this source language UR (French UR) to a target language UR (English UR).

French UR English UR
X plaire Y ↔ Eng(Y) like Eng(X)

(where Eng(X) means the English translation of X)

Londres plaire Sam (source UR) → Sam like London (target UR)

3. target language grammar translates the target language UR into an actual target language sentence.

Sam like London → Sam likes London.

29

Things to note about transfer systems

- The transfer mechanism is essentially reversible; e.g., the *plaire* rule works in both directions (at least in theory)
- Because we have a separate target language grammar, we are able to ensure that the rules of English apply; *like* → *likes*.
- Word order is handled differently than with transformers: the URs are essentially unordered.
- The underlying representation can be of various levels of abstraction – words, syntactic trees, meaning representations, etc.; we will talk about this with the **translation triangle**.

30

Caveat about reversibility

It seems like reversible rules are highly desirable—and in general they are—but we may not always want reversible rules.

e.g., Dutch *aanvangen* should be translated into English as *begin*, but English *begin* should be translated into Dutch as *beginnen*.

31

Levels of abstraction

There are differing levels of abstraction which you can transfer – i.e. your URs can be in different forms. We have looked so far at URs that represent only word information.

We can do a full syntactic analysis, which helps us to know how the words in a sentence relate.

Or we can do only a partial syntactic analysis, such as representing the dependencies between words.

32

Czech-English example

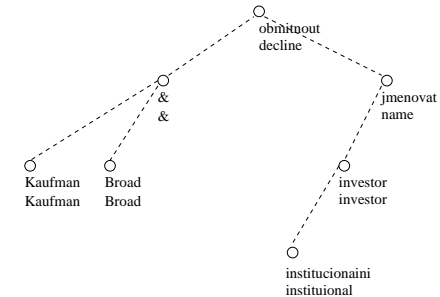
- (9) *Kaufman & Broad odmítla institucionální investory jmenovat.*
Kaufman & Broad declined institutional investors to name/identify
 'Kaufman & Broad refused to name the institutional investors.'

Example taken from Čmejrek, Cuřín, and Havelka (2003).

- They find the baseforms of words (e.g., *obmítout* 'to decline' instead of *odmítla* 'declined')
- They find which words depend on which other words and represent this in a tree (the noun *investory* depends on the verb *odmítla*)
- This dependency tree is then converted to English (comparative grammar) and re-ordered as appropriate.

33

Czech-English dependency tree



34

Interlinguas

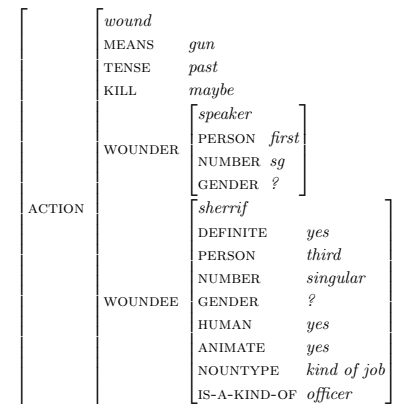
Ideally, we could use an **interlingua** = a language-independent representation of meaning.

Benefit: To add new languages to your MT system, you merely have to provide mapping rules between your language and the interlingua, and then you can translate into any other language in your system.

What your interlingua looks like depends on your goals, but an example might be the following for *I shot the sherrif.*:

35

Interlingua example



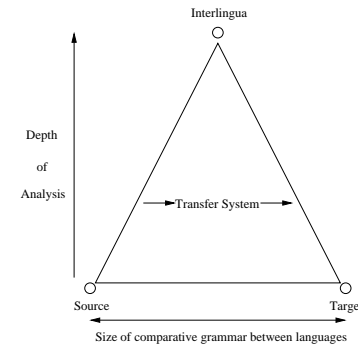
36

Interlingual problems

- What exactly should be represented in the interlingua?
⇒ English *corner* = Spanish *rincón* = 'inside corner' or *esquina* = 'outside corner'
- Can possibly lead to unnecessary work
Because Japanese distinguishes *older brother* from *younger brother*, we have to disambiguate English *brother* to put it into the interlingua. Then, if we translate into French, we have to ignore the disambiguation and simply translate it as *frère*, which simply means 'brother'.

37

The translation triangle



38

System 3: Machine learning

Instead of trying to tell the MT system how we're going to translate, we might try a **machine learning** approach = the computer will learn how to translate based on what it sees.

For this, we need

- **training data**
- a way of learning from that data

39

Let's use frequency (statistical methods)

We can look at how often a source language word is translated as a target language word – i.e. the **frequency** of a given translation, and choose the most frequent translation.

But how can we tell what a word is being translated as? There are two different cases:

- We are told what each word is translated as: **text alignment**
- We are not told what each word is translated as: use a **bag of words**

40

Text alignment

Sometimes humans have provided informative training data for us.

- sentence alignment
- word alignment

(The process of text alignment can also be automated and then fed into an MT system.)

41

Sentence alignment

sentence alignment = determine which source language sentences align with which target language ones (what we assumed in the bag of words example).

Intuitively easy, but can be difficult in practice since different languages have different punctuation conventions.

42

Word alignment

word alignment = determine which source language words align with which target language ones

- Much harder than sentence alignment to do automatically.
- But if it has already been done for us, it gives us good information about what a word's translation equivalent is.

Note that one word can map to one word or to multiple words. Likewise, sometimes it is best for multiple words to align with multiple words.

43

Different word alignments

English-Russian examples

- one-to-one: *khorosho* = *well*
- one-to-many: *kniga* = *the book*
- many-to-one: *to take a walk* = *gulyat'*
- many-to-many: *at least* = *khotya by* ('although if/would')

44

Calculating probabilities

With word alignments, it is relatively easy to calculate probabilities.

e.g., What is the probability that *run* translates as *correr* in Spanish?

1. Count up how many times *run* appears in the English part of your bi-text. e.g., 500 times
2. Out of all those times, count up how many times it was translated as (i.e. aligns with) *correr*. e.g., 275 (out of 500) times.
3. Divide to get a probability: $275/500 = 0.55$, or 55%

45

Word alignment difficulties

But knowing how words align in the training data will not tell us how to handle the new data we see.

- we may have many cases where *fool* is aligned with the Spanish *engañar* = 'to fool'
- but we may then encounter *a fool*, where the translation should be *tonto* (male) or *tonta* (female)

So, word alignment only helps us get some frequency numbers; we still have to do something intelligent with them.

46

Word alignment difficulties

Sometimes it is not even clear that word alignment is possible.

- (10) *Ivan aspirant.*
Ivan graduate student
'Ivan is a graduate student.'

What does *is* align with?

⇒ In cases like this, a word can be mapped to a "null" element in the other language.

47

The "bag of words" method

What if we're not given word alignments?

How can we tell which English words are translated as which German words if we are only given an English text and a corresponding German text?

- We can treat each sentence as a **bag of words** = unordered collection of words.
- If word A appears in a sentence, then we will record all of the words in the corresponding sentence in the other language as appearing with it.

48

Bag of words example

- English *He speaks Russian well.*
- Russian *On khorosho govorit po-russki.*

Eng	Rus	Eng	Rus
He	On	speaks	On
He	khorosho	speaks	khorosho
He	govorit
He	po-russki	well	po-russki

The idea is that, in the long run – over thousands, or even millions, of sentences, *He* will tend to appear more often with *On*, *speaks* will appear with *govorit*, and so on.

49

Calculating probabilities: sentence 1

So, for *He* in *He speaks Russian well/On khorosho govorit po-russki*, we do the following:

1. Count up the number of Russian words: 4.
2. Assign each word equal probability of translation: $1/4 = 0/25$, or 25%.

50

Calculating probabilities: sentence 2

If we also have *He is nice./On simpatich'nyi.*, then for *He*, we do the following:

1. Count up the number of possible translation words: 4 from the first sentence, 2 from the second = 6 total.
2. Count up the number of times *On* is the translation = 2 times out of 6 = $1/3 = 0.33$, or 33%.

Every other word has the probability $1/6 = 0.17$, or 17%, so *On* is clearly the best translation for *He*.

51

What makes MT hard?

We've seen how MT systems can work, but MT is a very difficult task because languages are vastly different. They differ:

- Lexically: In the words they use
- Syntactically: In the constructions they allow
- Semantically: In the way meanings work
- Pragmatically: In what readers take from a sentence.

In addition, there is a good deal of real-world knowledge that goes into a translation.

52

Lexical ambiguity

Words can be **lexically ambiguous** = have multiple meanings.

- *bank* can be a financial institution or a place along a river.
- *can* can be a cylindrical object, as well as the act of putting something into that cylinder (e.g., *John cans tuna.*), as well as being a word like *must*, *might*, or *should*.

⇒ We have to know which meaning before we translate.

53

How words divide up the world (lexical issues)

Words don't line up exactly between languages.

Within a language, we have synonyms, hyponyms, and hypernyms.

- *sofa* and *couch* are synonyms (mean the same thing)
- *sofa* is a hyponym (more specific term) of *furniture*
- *furniture* is a hypernym (more general term) of *sofa*

54

Synonyms

Often we find **synonyms** between two languages (as much as there are synonyms within a language):

- English *book* = Russian *kniga*
- English *music* = Spanish *música*

But words don't always line up exactly between languages.

55

Hypernyms and Hyponyms

- English **hypernyms** = words that are more general in English than in their counterparts in other languages
 - English *know* is rendered by the French *savoir* ('to know a fact') and *connaître* ('to know a thing')
 - English *library* is German *Bucheri* if it is open to the public, but *Bibliothek* if it is intended for scholarly work.
- English **hyponyms** = words that are more specific in English than in their foreign language counterparts.
 - The German word *berg* can mean either *hill* or *mountain* in English.
 - The Russian word *ruka* can mean either *hand* or *arm*.

56

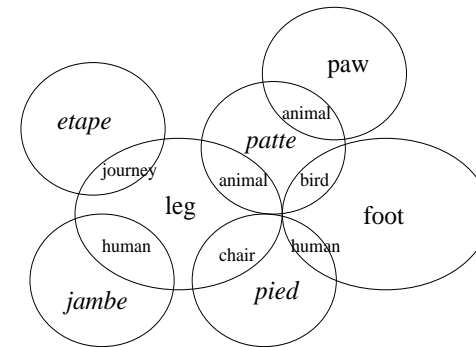
Semantic overlap

And then there's just fuzziness, as in the following English and French correspondences

- *leg* = *etape* (journey), *jambe* (human), *pied* (chair), *patte* (animal)
- *foot* = *pied* (human), *patte* (bird)
- *paw* = *patte* (animal)

57

Venn diagram of semantic overlap



58

Lexical gaps

Sometimes there is no simple equivalent for a word in a language, and the word has to be translated with a more complex phrase. We call this a **lexical gap** or **lexical hole**.

- French *gratiner* means something like 'to cook with a cheese coating'
- Hebrew *stam* means something like 'I'm just kidding' or 'Nothing special.'

59

Light verbs

Sometimes a verb carries little meaning, so we call it a **light verb**

- French *faire une promenade* is literally 'make a walk,' but it has the meaning of the English *take a walk*
- Dutch *een poging doen* 'do an attempt' means the same as the English *make an attempt*

60

Idioms

And we often face **idioms** = expressions whose meaning is not made up of the meanings of the individual words.

English *kick the bucket*

- approximately equivalent to the French *casser sa pipe* 'break his/her pipe'
- but we might want to translate it as *mourir* 'die')
- and we want to treat it differently than *kick the table*

61

Idiosyncracies

Along with that, there are idiosyncratic choices among languages

- English *heavy smoker*
- French *grand fumeur* = 'large smoker'
- German *stark Raucher* = 'strong smoker'

62

Taboo words

And then there are **taboo words** = words which are "forbidden" in some way or in some circumstances (i.e. swear/curse words)

- You of course know several English examples Note that the literal meanings of these words lack the emotive impact of the actual words.
- Other languages/cultures have different taboos: often revolving around death, body parts, bodily functions, disease, and religion.
e.g., The word 'skin' is taboo in a Western Australian (Aboriginal) language (<http://www.aija.org.au/online/ICABenchbook/BenchbookChapter5.pdf>)

Imagine encountering the word 'skin' in English and translating it without knowing this.

63

What sentences look like (syntax)

Different languages have different structures. In English, we have what is called a subject-verb-object (SVO) order, as in (11).

(11) *John punched Bill.*
SUBJECT VERB OBJECT

Japanese is SOV. Arabic is VSO. Dyrbal (Australian aboriginal language) has free word order.

MT systems have to account for these differences.

64

Syntax, pt. 2

Sometimes there are just different ways of saying things.

English *His name is Jerome.*

- (12) a. *Er heißt Jerome. (German)*
He untranslatable verb Jerome
b. *Il s'appelle Jerome. (French)*
He calls himself Jerome.

⇒ Words don't really align here.

65

Syntax & Semantics: what things mean

Even within a language, there are syntactic complications. We can have **structural ambiguities** = sentences where there are multiple ways of interpreting it.

(13) *John saw the boy (with the binoculars).*

with the binoculars can refer to either *the boy* or to how John saw the boy.

- This difference in structure corresponds to a difference in what we think the sentence means.
- Do we attempt to translate only one interpretation? Or do we try to preserve the ambiguity in the target language?

66

How language is used (pragmatics)

Translation becomes even more difficult when we try to translate something in context.

- *Thank you* is usually translated as *merci* in French, but it is translated as *s'il vous plaît* 'please' when responding to an offer.
- *Can you drive a stick-shift?* could be a request for you to drive my manual transmission automobile, or it could simply be a request for information about your driving abilities.

67

Real-world knowledge

Sometimes we have to use **real-world knowledge** to figure out what a sentence means.

(14) *Put the paper in the printer. Then switch it on.*

We know what *it* refers to only because we know that printers, not paper, can be switched on.

68

Ambiguity resolution

- If the source language involves ambiguous words/phrases, but the target language does not have the same ambiguity, we have to resolve ambiguity before translation.
e.g., the hyponyms/hypernyms we saw before.
- But sometimes we might want to preserve the ambiguity, or note that there was ambiguity or that there are a whole range of meanings available.
⇒ In the Bible, the Greek word *hyper* is used in 1 Corinthians 15:29; it can mean 'over', 'for', 'on behalf of', and so on. How you treat it affects how you treat the theological issue of salvation of the already dead.i.e. people care deeply about how you translate this word, yet it is not entirely clear what English meaning it has.

69

Evaluating MT systems

So, we've seen some translation systems and we know that translation is hard. The question now is: how do we evaluate MT systems?

The biggest users of MT systems are large corporations who need a lot of translation.

- How much change in the current setup will the MT system force?
Translator tasks will change from translation to updating the MT dictionaries and post-editing the results.
- How will it fit in with word processors and other software?
- Will the company selling the MT system be around in the next few years for support and updates?
- How fast is the MT system?
- How good is the MT system (quality)?

70

Evaluating quality

- **Intelligibility** = how understandable the output is
- **Accuracy** = how faithful the output is to the input
- **Error analysis** = how many errors we have to sort through (and how do the errors affect intelligibility & accuracy)
- **Test suite** = a set of sentences that our system should be able to handle

71

Intelligibility

Intelligibility Scale (taken directly from Arnold et al)

1. The sentence is perfectly clear and intelligible. It is grammatical and reads like ordinary text.
2. The sentence is generally clear and intelligible. Despite some inaccuracies or infelicities of the sentence, one can understand (almost) immediately what it means.
3. The general idea of the sentence is intelligible only after considerable study. The sentence contains grammatical errors and/or poor word choices.
4. The sentence is unintelligible. Studying the meaning of the sentence is hopeless; even allowing for context, one feels that guessing would be too unreliable.

72

Further reading

Many of the examples are adapted from the following books:

- Doug J. Arnold, Lorna Balkan, Siety Meijer, R. Lee Humphreys and Louisa Sadler (1994). *Machine Translation: an Introductory Guide*. Blackwells-NCC, London. 1994. Available from <http://www.essex.ac.uk/linguistics/clmt/MTbook/>
- Jurafsky, Daniel, and James H. Martin (2000). *Speech and Language Processing: An Introduction to Natural Language Processing, Speech Recognition, and Computational Linguistics*. Prentice-Hall. More info at <http://www.cs.colorado.edu/~martin/slp.html>.