

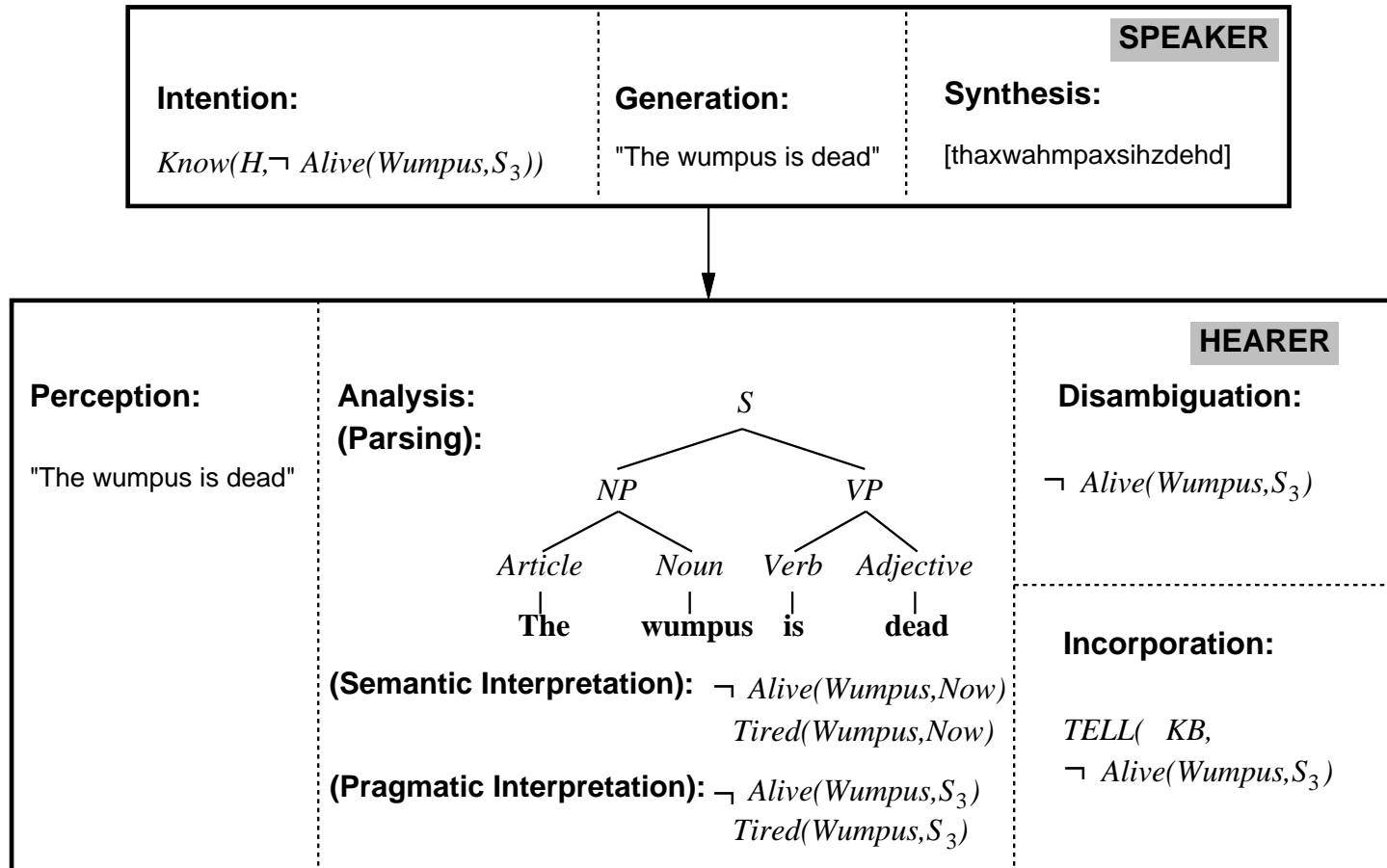
LANGUAGE

CHAPTER 22

Outline

- ◇ Communication
- ◇ Grammar
- ◇ Syntactic analysis
- ◇ Problems
- ◇ Grammar induction

Stages in communication: example



Grammar

Grammar specifies the compositional structure of complex messages

A formal language is a set of strings of terminal symbols

Each string in the language can be analyzed/generated by the grammar

The grammar is a set of rewrite rules, e.g.,

$$S \rightarrow NP VP$$

$$Article \rightarrow \mathbf{the} \mid \mathbf{a} \mid \mathbf{an} \mid \dots$$

Here S is the sentence symbol, NP and VP are nonterminals

Grammar types

Regular: *nonterminal* \rightarrow *terminal*[*nonterminal*]

$$S \rightarrow aS$$

$$S \rightarrow \Lambda$$

Context-free: *nonterminal* \rightarrow *anything*

$$S \rightarrow aSb$$

Context-sensitive: a single nonterminal in a string is replaced by a string

$$ASB \rightarrow AAaBB$$

Recursively enumerable: no constraints

Natural languages probably context-free, parsable in real time!

Wumpus lexicon

- Noun* → *stench* | *breeze* | *glitter* | *nothing*
| *wumpus* | *pit* | *pits* | *gold* | *east* | ...
- Verb* → *is* | *see* | *smell* | *shoot* | *feel* | *stinks*
| *go* | *grab* | *carry* | *kill* | *turn* | ...
- Adjective* → *right* | *left* | *east* | *south* | *back* | *smelly* | ...
- Adverb* → *here* | *there* | *nearby* | *ahead*
| *right* | *left* | *east* | *south* | *back* | ...
- Pronoun* → *me* | *you* | *I* | *it* | ...
- Name* → *John* | *Mary* | *Waikato* | *Otago* | ...
- Article* → *the* | *a* | *an* | ...
- Preposition* → *to* | *in* | *on* | *near* | ...
- Conjunction* → *and* | *or* | *but* | ...
- Digit* → **0** | **1** | **2** | **3** | **4** | **5** | **6** | **7** | **8** | **9**

Divided into **closed** and **open** classes

Wumpus grammar

S	$\rightarrow NP VP$	I + feel a breeze
	S <i>Conjunction</i> S	I feel a breeze + and + I smell a wumpus
NP	\rightarrow <i>Pronoun</i>	I
	<i>Noun</i>	pits
	<i>Article Noun</i>	the + wumpus
	<i>Digit Digit</i>	3 4
	$NP PP$	the wumpus + to the east
	NP <i>RelClause</i>	the wumpus + that is smelly
VP	\rightarrow <i>Verb</i>	stinks
	$VP NP$	feel + a breeze
	VP <i>Adjective</i>	is + smelly
	$VP PP$	turn + to the east
	VP <i>Adverb</i>	go + ahead
PP	\rightarrow <i>Preposition NP</i>	to + the east
<i>RelClause</i>	\rightarrow that VP	that + is smelly

Parse trees

Exhibit the grammatical structure of a sentence

I **shoot** **the** **wumpus**

Parse trees

Exhibit the grammatical structure of a sentence

Pronoun

I

Verb

shoot

Article

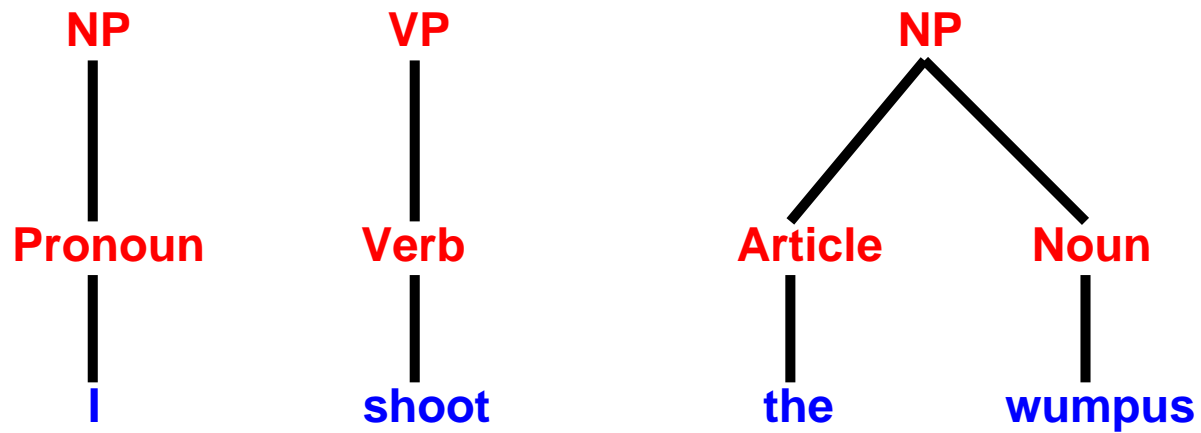
the

Noun

wumpus

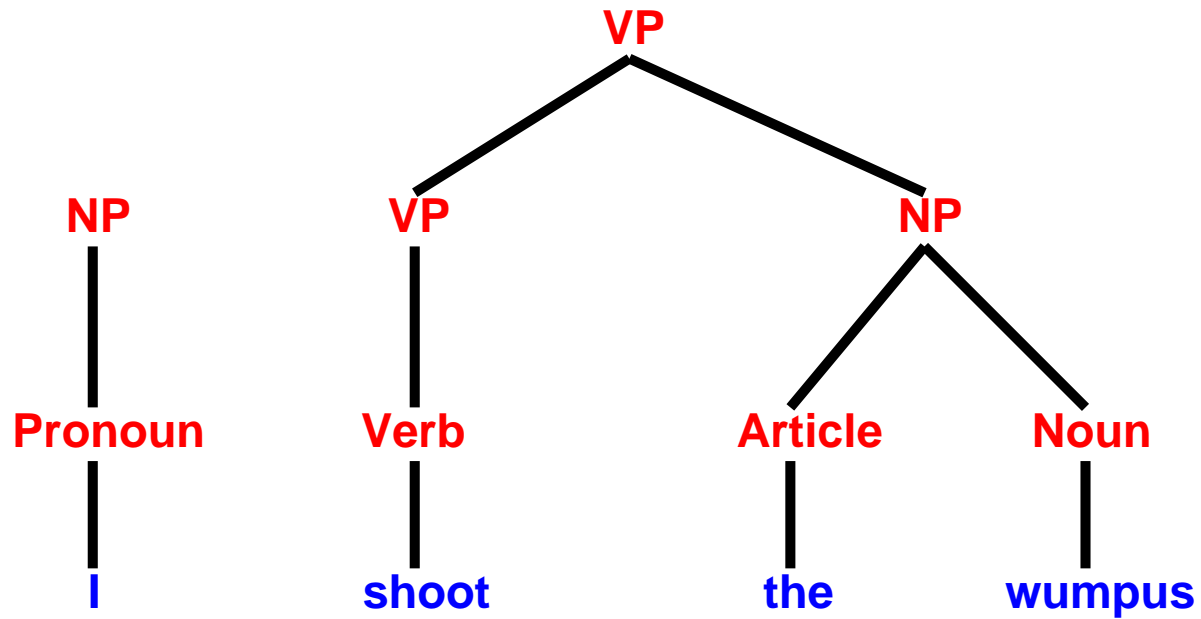
Parse trees

Exhibit the grammatical structure of a sentence



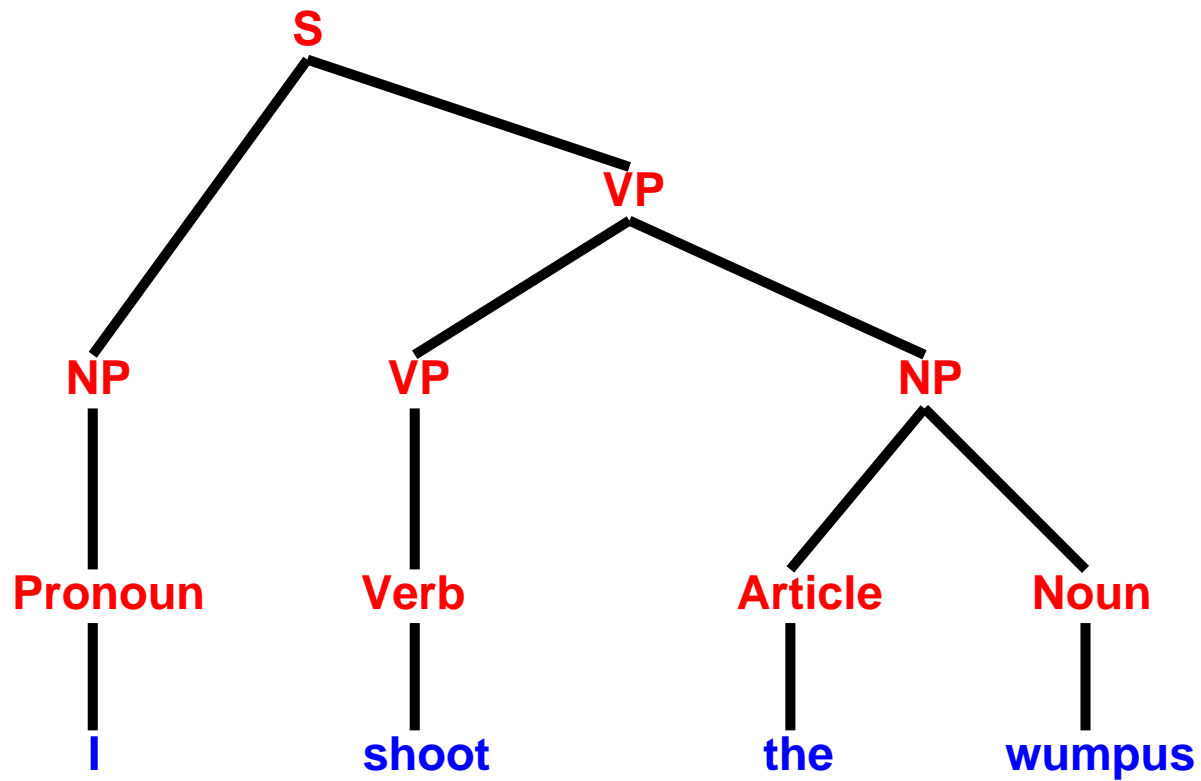
Parse trees

Exhibit the grammatical structure of a sentence



Parse trees

Exhibit the grammatical structure of a sentence



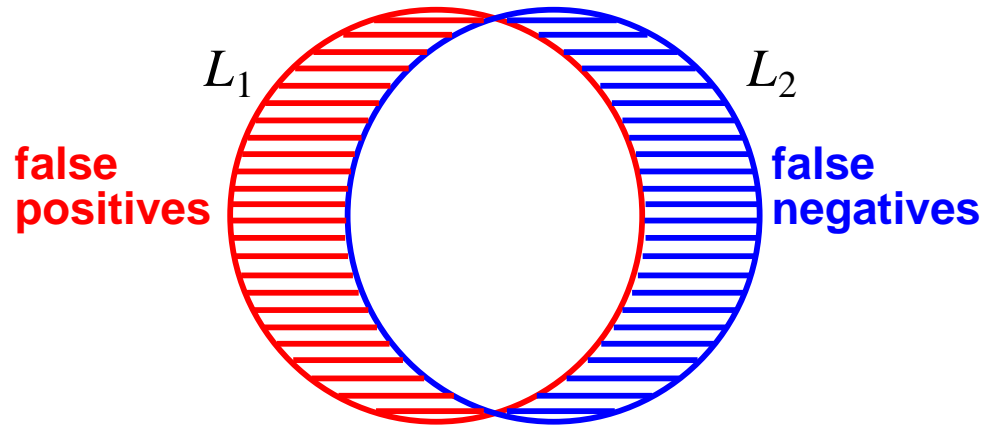
Context-free parsing

Bottom-up parsing works by replacing any substring that matches RHS of a rule with the rule's LHS

Efficient algorithms (e.g., chart parsing, Section 22.3) $O(n^3)$ for context-free, run at several thousand words/sec for real grammars

Language learning

Formal language L_1 may differ from natural language L_2



Adjusting L_1 to agree with L_2 is a learning problem

Real grammars 10–500 pages, insufficient even for “proper” English

Probabilistic grammars are an alternative

Real language

Real human languages provide many problems for NLP:

- ◇ ambiguity
- ◇ anaphora
- ◇ indexicality
- ◇ metonymy
- ◇ metaphor

etc.

Ambiguity

Squad helps dog bite victim

Helicopter powered by human flies

American pushes bottle up Germans

Anaphora

Using pronouns to refer back to entities already introduced in the text

After Mary proposed to John, **they** found a preacher and got married.

For the honeymoon, **they** went to Hawaii

Mary saw a ring through the window and asked John for **it**

Mary threw a rock at the window and broke **it**

Indexicality

Indexical sentences refer to utterance situation (place, time, etc.)

I am over **here**

Why did **you** do **that**?

Metonymy

Using one noun phrase to stand for another

I've read **Shakespeare**

Chrysler announced record profits

The **ham sandwich** on Table 4 wants another beer

Metaphor

“Non-literal” usage of words and phrases, often systematic:

I’ve tried killing the process but it won’t die. Its parent keeps it alive.

A linear-time algorithm for grammar learning

SEQUITUR algorithm developed by Nevill-Manning and Witten at Waikato: learns grammar that generates only one sentence, namely the input text

Basic idea: good grammar is compact grammar (i.e., grammar should compress the input text)

Scans input from left to right, one symbol at a time, building up a grammar for the text seen so far.

Grammar is built by enforcing the following two constraints:

No pair of adjacent symbols must appear more than once in the grammar.

Every rule in the grammar must be used at least twice (apart from the start rule).

SEQUITUR example

	Input	Grammar	Comments
1	<i>a</i>	$S \rightarrow a$	
2	<i>ab</i>	$S \rightarrow ab$	
3	<i>abc</i>	$S \rightarrow abc$	
4	<i>abcd</i>	$S \rightarrow abcd$	
5	<i>abcdb</i>	$S \rightarrow abcdb$	
6	<i>abcdbc</i>	$S \rightarrow abcdbc$	<i>bc</i> twice
		$S \rightarrow aAdA; A \rightarrow bc$	
7	<i>abcdbca</i>	$S \rightarrow aAdAa; A \rightarrow bc$	
8	<i>abcdbcab</i>	$S \rightarrow aAdAab; A \rightarrow bc$	
9	<i>abcdbcabc</i>	$S \rightarrow aAdAabc; A \rightarrow bc$	<i>bc</i> twice
		$S \rightarrow aAdAaA; A \rightarrow bc$	<i>aA</i> twice
		$S \rightarrow BdAB; A \rightarrow bc; B \rightarrow aA$	
10	<i>abcdbcabcd</i>	$S \rightarrow BdABd; A \rightarrow bc; B \rightarrow aA$	<i>Bd</i> twice
		$S \rightarrow CAC; A \rightarrow bc; B \rightarrow aA; C \rightarrow Bd$	<i>B</i> only once
		$S \rightarrow CAC; A \rightarrow bc; C \rightarrow aAd$	