

Complexity and Character of Human Languages

Informatics 2A: Lecture 21

Mirella Lapata

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

04 November 2011

- 1 Human Language Complexity
 - Chomsky Hierarchy
 - The Faculty of Language
 - Strong and Weak Adequacy

- 2 Linear Indexed Grammars

Reading: J&M. Chapter 16.3–16.4; Hauser, Chomsky, and Fitch (2002)

1 / 24

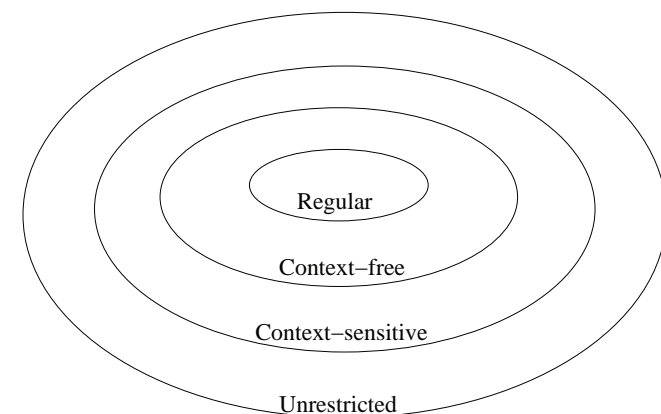
2 / 24

Review

Review

Chomsky Hierarchy: classifies languages on scale of complexity:

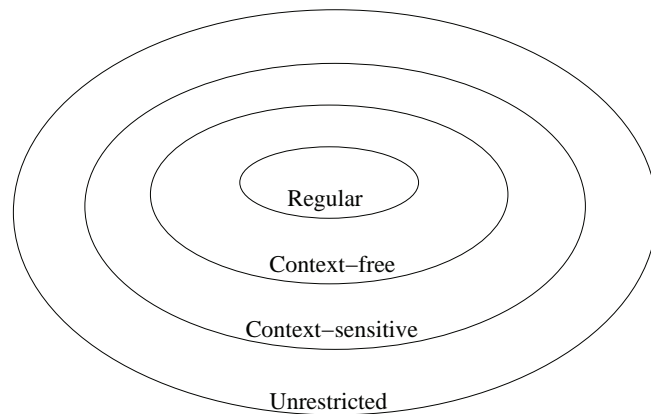
- **Regular** languages: those whose phrases can be 'recognized' by a finite state machine.
- **Context-free** languages: most programming languages, and many aspects of natural languages can be described at this level; the set of languages accepted by pushdown automata.
- **Context-sensitive** languages: equivalent with a linear bounded nondeterministic Turing machine, also called a linear bounded automaton.
- **Unrestricted** languages: *all* languages that can in principle be defined via mechanical rules.



3 / 24

4 / 24

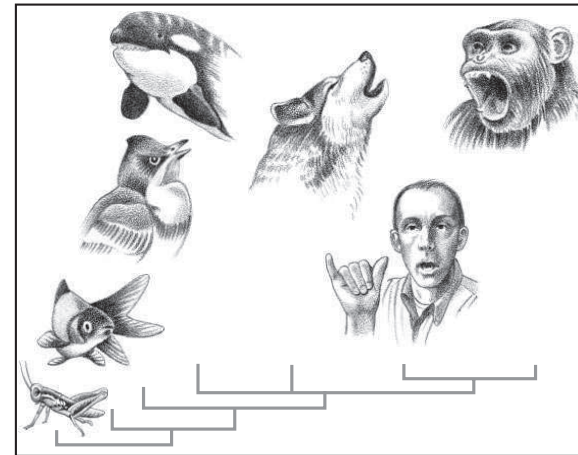
Review



Where do human languages fit within this complexity hierarchy?

4 / 24

The Faculty of Language



The “language faculty” has a **broad sense** and a **narrow sense** (Hauser, Chomsky, and Fitch 2002).

5 / 24

The Faculty of Language (Broad Sense)

Sensory-motor system

- for producing and perceiving linguistic communication
- spoken language: vocal track, auditory system
- sign language: gestural system, visual system
- written language: writing system, visual or tactile system

Conceptual-intentional system

- who to communicate with and what to communicate about
- generating mental states and attributing them to others;
- acquiring conceptual representations that are non-linguistic;
- referring to entities and events.

6 / 24

The Faculty of Language (Narrow Sense)

Abstract computational system

- one part of which is **narrow syntax** which generates representations internal to the mind/brain and maps them to:
- sensory-motor interface through phonological, gestural system;
- conceptual-intentional system through semantic (and pragmatic) systems.

A core property of narrow syntax is **recursion**: takes a finite set of elements and yields a potentially infinite array of discrete expressions.

7 / 24

Recursion

The potential infiniteness of the language faculty has been recognized by Galileo, Descartes, von Humboldt.

Discrete Infinity

- Sentences are built up by discrete units
- There are 6-word sentences, and 7-word sentences, but no 6.5 word sentences
- **There is no longest sentence!**
- **There is no non-arbitrary upper bound to sentence length!**

Mary thinks that John thinks that George thinks that Mary thinks that this course is boring!
I ate lunch and slept and watched tv and went to the bathroom and had a coffee and got dressed ...

8 / 24

Strong and Weak Adequacy

Questions about the formal complexity of language are about the computational power of syntax, as represented by a grammar that's **adequate** for it.

A strongly adequate grammar

- generates all and only the strings of the language;
- assigns them the "right" structures — ones that support a correct representation of meaning.

A weakly adequate grammar

generates all and only the strings of a language but assigns them "wrong" structures.

9 / 24

Is Natural Language Regular?

It is generally agreed that natural languages are not regular!

Center-embedding

[The cat₁ likes tuna fish₁].
[The cat₁ [the dog₂ chased₂] likes tuna fish₁].
[The cat₁ [the dog₂ [the rat₃ bit₃] chased₂] likes tuna fish₁].

Idea of proof

$(\text{the} + \text{noun})^n (\text{transitive verb})^{n-1}$ likes tuna fish.
 $A = \{ \text{the cat, the dog, the rat, the elephant, the kangaroo} \dots \}$
 $B = \{ \text{chased, bit, admired, ate, befriended} \dots \}$
 Intersect $/A^* B^* \text{ likes tuna fish}/$ with English
 $L = x^n y^{n-1}$ likes tuna fish, $x \in A, y \in B$
 Use pumping lemma to show L is not regular

10 / 24

Is Natural Language Context Free?

It doesn't look like it is context free either! Evidence comes from a Swiss German dialect and Bambara, a language spoken in Mali.

Crossing dependencies

omdat Wim₁ Jan₂ Henk₃ de kinderen₄ zag₁ helpen₂ leren₃ zwemmen₄
 because Wim₁ Jan₂ Henk₃ the children₄ saw₁ help learn swim₄
 because Wim saw Jan help Henk teach the children to learn to swim
 $|zag|$ depends on $|Wim|$, and $|helpen|$ depends on $|Jan|$, etc.

Idea of Proof

Languages $\{xx \mid x \in \{a, b\}^*\}$ are not context-free.
 Related $a^n b^m c^n d^m$ language also not context-free.
 Swiss German crossing dependencies equivalent to $a^n b^m c^n d^m$

11 / 24

Is Natural Language Context Free?

It doesn't look like it is context free either! Evidence comes from a Swiss German dialect and Bambara, a language spoken in Mali.

Crossing dependencies

omdat Wim₁ Jan₂ Henk₃ de kinderen₄ zag₁ helpen₂ leren₃ zwemmen₄
 because Wim₁ Jan₂ Henk₃ the children₄ saw₁ help learn swim₄
 because Wim saw Jan help Henk teach the children to learn to swim
 | zag | depends on | Wim |, and | helpen | depends on | Jan |, etc.

Idea of Proof

Languages $\{xx|x \in \{a, b\}^*\}$ are not context-free.
 Related $a^n b^m c^n d^m$ language also not context -free.
 Swiss German crossing dependencies equivalent to $a^n b^m c^n d^m$

11 / 24

Is Natural Language Context Free?

It doesn't look like it is context free either! Evidence comes from a Swiss German dialect and Bambara, a language spoken in Mali.

Crossing dependencies

omdat Wim₁ Jan₂ Henk₃ de kinderen₄ zag₁ helpen₂ leren₃ zwemmen₄
 because Wim₁ Jan₂ Henk₃ the children₄ saw₁ help learn swim₄
 because Wim saw Jan help Henk teach the children to learn to swim
 | zag | depends on | Wim |, and | helpen | depends on | Jan |, etc.

Idea of Proof

Languages $\{xx|x \in \{a, b\}^*\}$ are not context-free.
 Related $a^n b^m c^n d^m$ language also not context -free.
 Swiss German crossing dependencies equivalent to $a^n b^m c^n d^m$

11 / 24

Is Natural Language Context Free?

It doesn't look like it is context free either! Evidence comes from a Swiss German dialect and Bambara, a language spoken in Mali.

Crossing dependencies

omdat Wim₁ Jan₂ Henk₃ de kinderen₄ zag₁ helpen₂ leren₃ zwemmen₄
 because Wim₁ Jan₂ Henk₃ the children₄ saw₁ help learn swim₄
 because Wim saw Jan help Henk teach the children to learn to swim
 | zag | depends on | Wim |, and | helpen | depends on | Jan |, etc.

Idea of Proof

Languages $\{xx|x \in \{a, b\}^*\}$ are not context-free.
 Related $a^n b^m c^n d^m$ language also not context -free.
 Swiss German crossing dependencies equivalent to $a^n b^m c^n d^m$

11 / 24

Is Natural Language Context Free?

It doesn't look like it is context free either! Evidence comes from a Swiss German dialect and Bambara, a language spoken in Mali.

Crossing dependencies

omdat Wim₁ Jan₂ Henk₃ de kinderen₄ zag₁ helpen₂ leren₃ zwemmen₄
 because Wim₁ Jan₂ Henk₃ the children₄ saw₁ help learn swim₄
 because Wim saw Jan help Henk teach the children to learn to swim
 | zag | depends on | Wim |, and | helpen | depends on | Jan |, etc.

Idea of Proof

Languages $\{xx|x \in \{a, b\}^*\}$ are not context-free.
 Related $a^n b^m c^n d^m$ language also not context -free.
 Swiss German crossing dependencies equivalent to $a^n b^m c^n d^m$

11 / 24

Linear Indexed Grammars

A **linear indexed grammar** (LIG) is more powerful than a CFG, but much less powerful than an arbitrary CSG; it is a “**mildly CS**” grammar.

Definition

An indexed grammar has **three** disjoint sets of symbols: terminals, non-terminals and **indices**.

An index is a **stack** of symbols that can be passed from the LHS of a rule to its RHS, allowing counting and recording what rules were applied in what order.

12 / 24

Linear Indexed Grammars

$S \rightarrow D_f$ pushes an f onto the index on D
 $D \rightarrow D_g$ pushes a g onto the index on D
 $D \rightarrow ABC$ passes the index on D to A , B and C

$g = \langle A \rightarrow Aa \mid B \rightarrow Bb \mid C \rightarrow Cc \rangle$ pops g from an index
 $f = \langle A \rightarrow a \mid B \rightarrow b \mid C \rightarrow c \rangle$ pops f from an index

13 / 24

Derivation in an Indexed Grammar

$S \rightarrow D_f$ $g = \langle A \rightarrow Aa \mid B \rightarrow Bb \mid C \rightarrow Cc \rangle$
 $D \rightarrow D_g$ $f = \langle A \rightarrow a \mid B \rightarrow b \mid C \rightarrow c \rangle$
 $D \rightarrow ABC$

S

14 / 24

Derivation in an Indexed Grammar

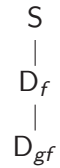
$S \rightarrow D_f$ $g = \langle A \rightarrow Aa \mid B \rightarrow Bb \mid C \rightarrow Cc \rangle$
 $D \rightarrow D_g$ $f = \langle A \rightarrow a \mid B \rightarrow b \mid C \rightarrow c \rangle$
 $D \rightarrow ABC$

S
|
D_f

15 / 24

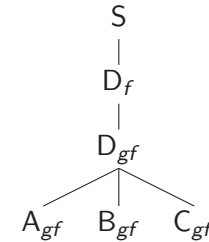
Derivation in an Indexed Grammar

$$\begin{aligned}
 S &\rightarrow D_f & g &= \langle A \rightarrow Aa \mid B \rightarrow Bb \mid C \rightarrow Cc \rangle \\
 D &\rightarrow D_g & f &= \langle A \rightarrow a \mid B \rightarrow b \mid C \rightarrow c \rangle \\
 D &\rightarrow ABC
 \end{aligned}$$



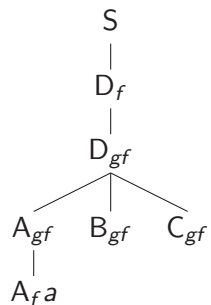
Derivation in an Indexed Grammar

$$\begin{aligned}
 S &\rightarrow D_f & g &= \langle A \rightarrow Aa \mid B \rightarrow Bb \mid C \rightarrow Cc \rangle \\
 D &\rightarrow D_g & f &= \langle A \rightarrow a \mid B \rightarrow b \mid C \rightarrow c \rangle \\
 D &\rightarrow ABC
 \end{aligned}$$



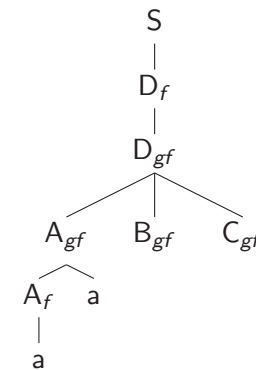
Derivation in an Indexed Grammar

$$\begin{aligned}
 S &\rightarrow D_f & g &= \langle A \rightarrow Aa \mid B \rightarrow Bb \mid C \rightarrow Cc \rangle \\
 D &\rightarrow D_g & f &= \langle A \rightarrow a \mid B \rightarrow b \mid C \rightarrow c \rangle \\
 D &\rightarrow ABC
 \end{aligned}$$



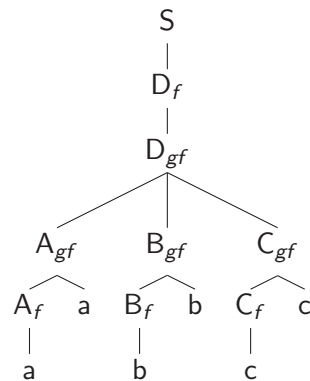
Derivation in an Indexed Grammar

$$\begin{aligned}
 S &\rightarrow D_f & g &= \langle A \rightarrow Aa \mid B \rightarrow Bb \mid C \rightarrow Cc \rangle \\
 D &\rightarrow D_g & f &= \langle A \rightarrow a \mid B \rightarrow b \mid C \rightarrow c \rangle \\
 D &\rightarrow ABC
 \end{aligned}$$



Derivation in an Indexed Grammar

$S \rightarrow D_f$ $g = \langle A \rightarrow Aa \mid B \rightarrow Bb \mid C \rightarrow Cc \rangle$
 $D \rightarrow D_g$ $f = \langle A \rightarrow a \mid B \rightarrow b \mid C \rightarrow c \rangle$
 $D \rightarrow ABC$



20 / 24

Linear Indexed Grammars

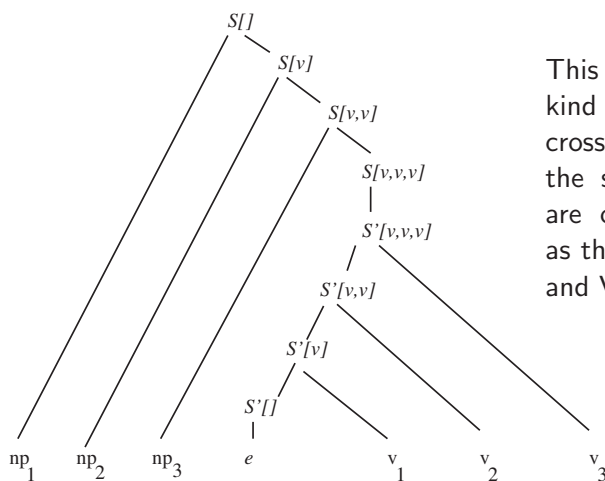
Linear Indexed Grammars (LIGs) allow an index to pass to only **one** non-terminal on the RHS (not three, as in previous example).

Here we'll push numbers onto an index.
An LIG for crossing dependencies in $np^k v^k$:

$S_{[...]} \rightarrow np_i S_{[i,...]}$ emit NP, push a number
 $S_{[...]} \rightarrow S'_{[...]}$ switch to verb sequence rule
 $S'_{[i,...]} \rightarrow S'_{[...]} v_i$ pop a number, emit a verb
 $S'_{[]} \rightarrow \epsilon$ stop if stack is empty

21 / 24

Example: LIG derivation for $np^3 v^3$



This grammar produces the kind of strings we want for crossing dependencies, but the structures it generates are only **weakly adequate**, as they don't associate NPs and Vs directly.

22 / 24

Linear Indexed Grammars

As a consequence of the weak adequacy of LIGs, other **"mildly CS"** grammar formalisms have been developed that are strongly adequate for NL:

- Tree Adjoining Grammar (TAG): a system of **tree** re-writing rules (ie, not string re-writing rules) in which elementary trees are combined by substitution and adjunction;
- Combinatory Categorical Grammar (CCG): a system that links words to complex categories that specify how adjacent words fit together, in terms of combinators like **apply** a functor to an argument, **compose** two functors, etc..

23 / 24

Summary

- The faculty of language contains a computational system that generates syntactic representations that can be mapped onto meanings.
- This raises the question of the complexity of this system (its position in the Chomsky hierarchy).
- A weakly adequate grammar generates the correct strings, while a strongly adequate one also generates the correct structures.
- There are structures in NLS which can be mapped on formal languages which are not context-free.
- NL probably belongs to the class of mildly context-sensitive languages, whose least powerful member (LIGs) is weakly adequate for NL.

Next Lecture: models of human parsing.