

Overview of Natural Language Processing

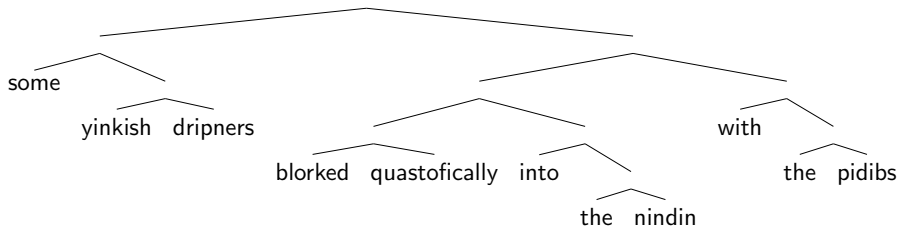
Part II & ACS L90

Lecture 4: Phrase Structure and Structured Prediction

Weiwei Sun

Department of Computer Science and Technology
University of Cambridge

Michaelmas 2023/24



Words are organized into nested blocks

Lecture 4: Phrase Structure and Structured Prediction

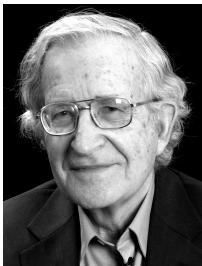
1. Phrase structure
2. Dependency structure
3. Structured prediction
4. Probabilistic Context-free grammars
5. Context-sensitive languages

Phrase Structure

Interview of Noam Chomsky by Lex Fridman

- (1) a. the guy who **fixed** the car very carefully **packed** his tools
b. very carefully, the guy who **fixed** the car **packed** his tools
c. *very carefully, the guy who **fixed** the car **is tall**

*I think the deepest property of language and puzzling property that's been discovered is what is sometimes called **structure dependence**. [...] Linear closeness is an easy computation, but here you're doing a much more, what looks like a more complex computation.*



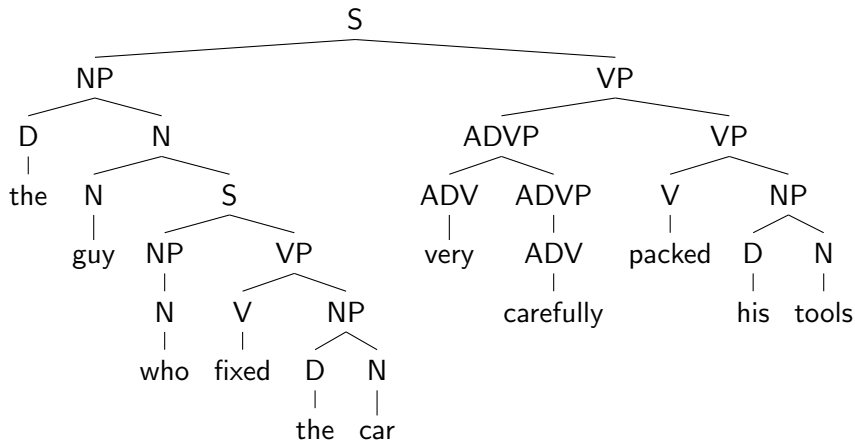
Noam Chomsky: Language, Cognition, and Deep Learning

www.youtube.com/watch?v=cMscNuSUy0I

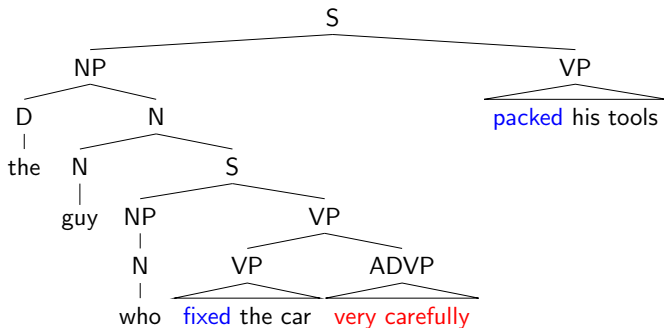
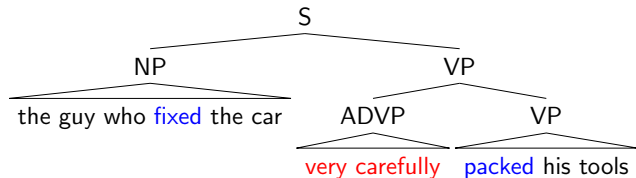
Constituency (phrase structure)

The basic idea

Phrase structure organizes words into *nested constituents*, which can be represented as **a tree**.

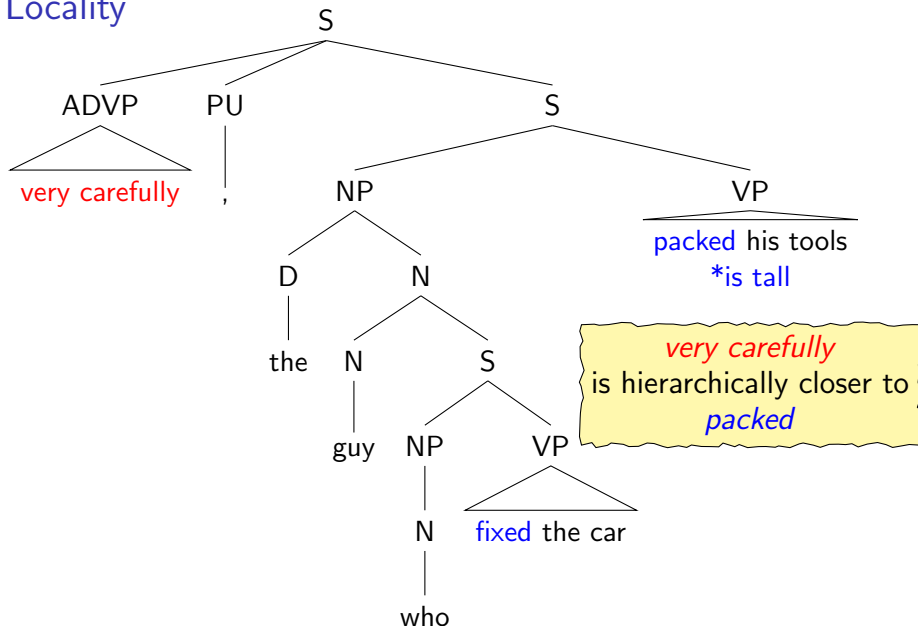


Different structures, different meaning

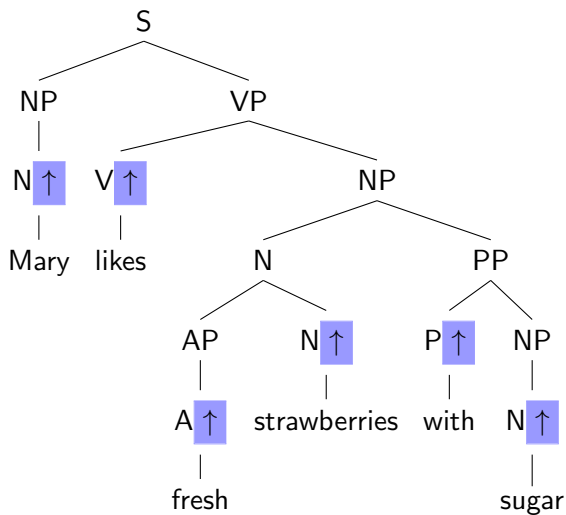


Results by a cool parser: <http://erg.delph-in.net/logon>

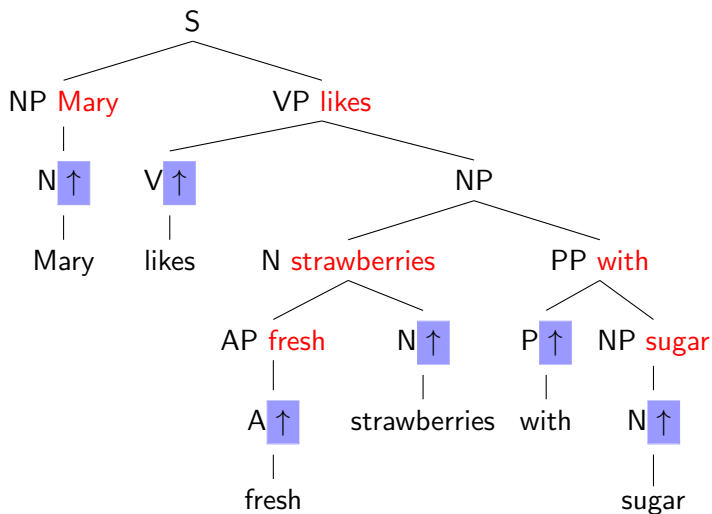
Locality



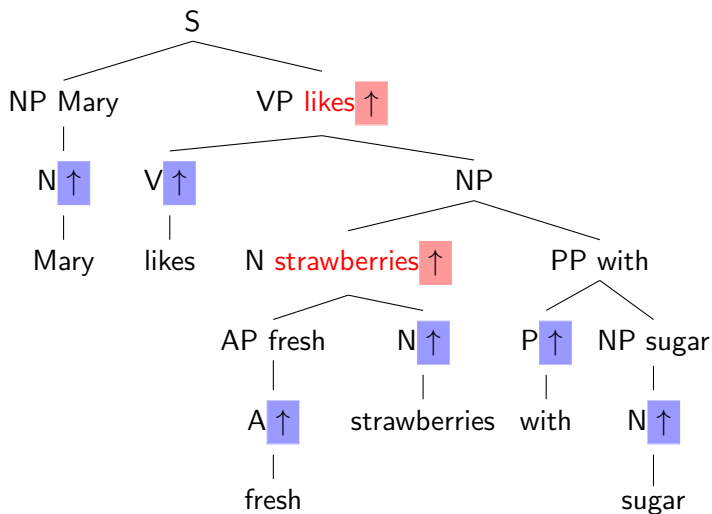
Projection and bi-lexical dependencies



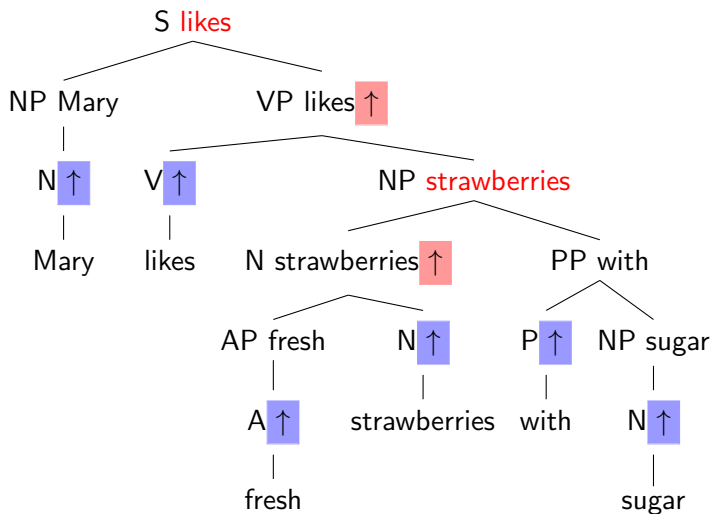
Projection and bi-lexical dependencies



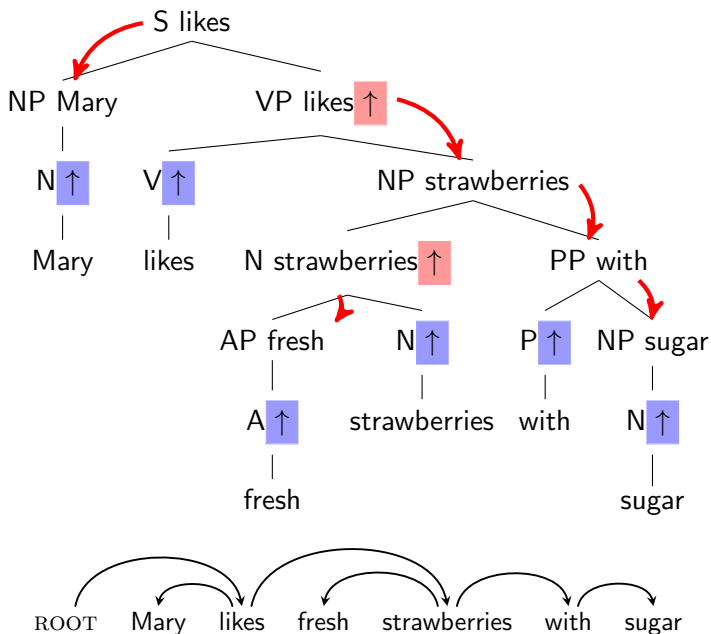
Projection and bi-lexical dependencies



Projection and bi-lexical dependencies



Projection and bi-lexical dependencies



Applications of parsing

Modern parsers are quite accurate

For some languages, automatic syntactic parsing is good enough to help in a range of NLP tasks

- Machine translation
- Information extraction
- Grammar checking
- etc.

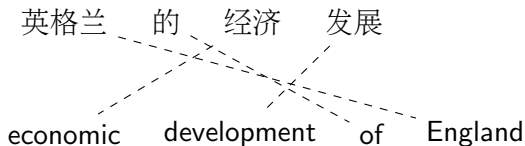
Applications of parsing

Modern parsers are quite accurate

For some languages, automatic syntactic parsing is good enough to help in a range of NLP tasks

- Machine translation
- Information extraction
- Grammar checking
- etc.

Translate “英格兰的经济发展” into English



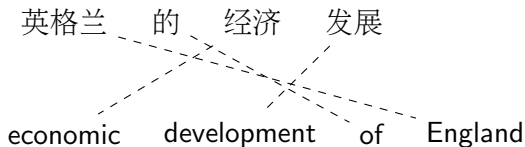
Applications of parsing

Modern parsers are quite accurate

For some languages, automatic syntactic parsing is good enough to help in a range of NLP tasks

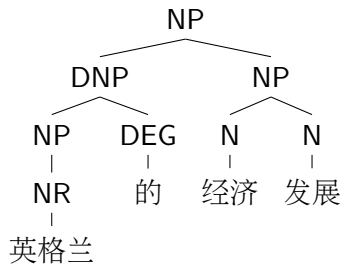
- Machine translation
- Information extraction
- Grammar checking
- etc.

Translate “英格兰的经济发展” into English



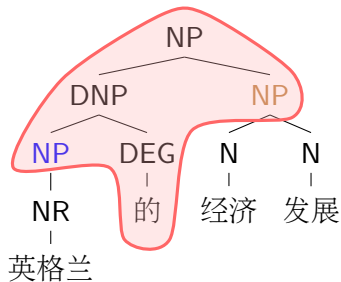
finite-state transducer?

An example: Tree-to-string transduction



R1	NP NR 英格兰	X England
R2	NP / \ DNP NP / \ NP DEG 的	X / \ X of X
R3	NP / \ N N	X / \ X X
R4	N 经济	X economic
R5	N 发展	X development

An example: Tree-to-string transduction

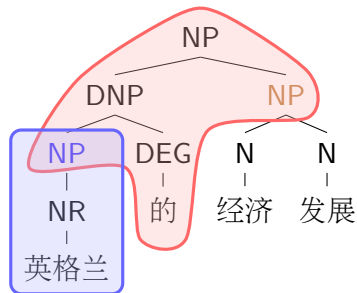


NP of NP

▷ R2

R1	NP └── NR 英格兰	X └── England
R2	NP ├── DNP │ ├── NP │ │ └── DEG │ │ └── 的 └── NP	X ├── X │ └── of └── X
R3	NP ├── N │ └── 经济 └── N └── 发展	X ├── X │ └── economic └── X
R4	N └── 经济	X └── economic
R5	N └── 发展	X └── development

An example: Tree-to-string transduction



NP of NP

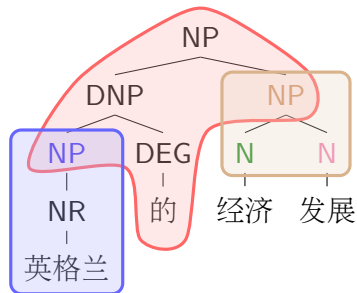
NP of England

▷R2

▷R1

R1	<p>NP</p> <p>NR</p> <p>英格兰</p>	<p>X</p> <p>England</p>
R2	<p>NP</p> <p>DNP NP</p> <p>NP DEG</p> <p>的</p>	<p>X</p> <p>X of X</p>
R3	<p>NP</p> <p>N N</p> <p>经济 发展</p>	<p>X</p> <p>X X</p>
R4	<p>N</p> <p>经济</p>	<p>X</p> <p>economic</p>
R5	<p>N</p> <p>发展</p>	<p>X</p> <p>development</p>

An example: Tree-to-string transduction



NP of NP

NP of England

N N of England

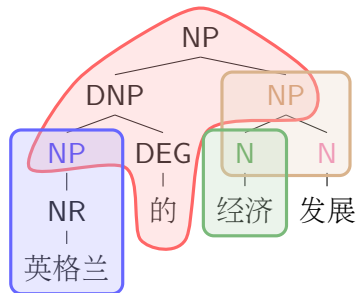
▷R2

▷R1

▷R3

R1	NP NR 英格兰	X England
R2	NP / \ DNP NP / \ NP DEG 的	X / \ X of X
R3	NP / \ N N 经济 发展	X / \ X X
R4	N 经济	X economic
R5	N 发展	X development

An example: Tree-to-string transduction



NP of NP

NP of England

N N of England

economic N of England

▷R2

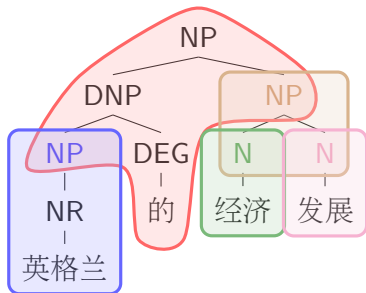
▷R1

▷R3

▷R4

R1	<pre> NP NR 英格兰 </pre>	<pre> X England </pre>
R2	<pre> NP / \ DNP NP / \ NP DEG / \ NR DEG 的 英格兰 </pre>	<pre> X / \ X of X </pre>
R3	<pre> NP / \ N N / \ 经济 发展 </pre>	<pre> X / \ X X </pre>
R4	<pre> N 经济 </pre>	<pre> X economic </pre>
R5	<pre> N 发展 </pre>	<pre> X development </pre>

An example: Tree-to-string transduction



NP of NP

NP of England

N N of England

economic N of England

economic development of England

▷R2

▷R1

▷R3

▷R4

▷R5

recursive
form transformation

R1	<pre> NP NR 英格兰 </pre>	<pre> X England </pre>
R2	<pre> NP / \ DNP NP / \ NP DEG / \ NR 的 英格兰 </pre>	<pre> X / \ X of X </pre>
R3	<pre> NP / \ N N / \ 经济 发展 </pre>	<pre> X / \ X X </pre>
R4	<pre> N 经济 </pre>	<pre> X economic </pre>
R5	<pre> N 发展 </pre>	<pre> X development </pre>

Structured Prediction

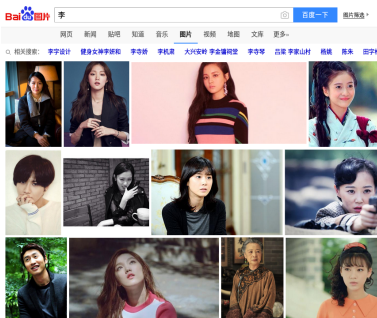
pre-lecture: watch this video

▶ www.youtube.com/watch?v=bjUwSHGsG9o

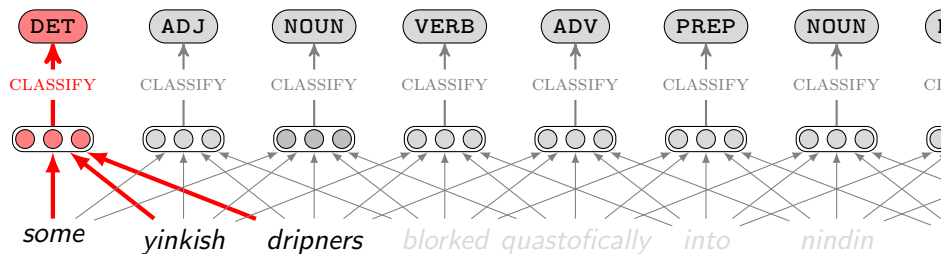
Muhammad Li

Howard Who's Muhammad Li?

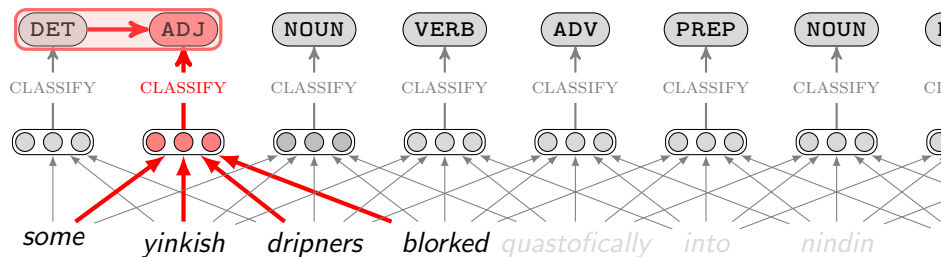
Sheldon Muhammad is the most common first name in the world, Li, the most common surname. As I didn't know the answer, I thought that gave me a mathematical edge.



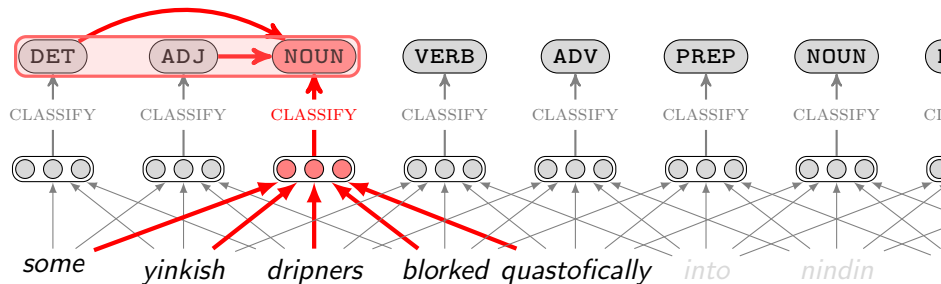
POS tagging and prediction



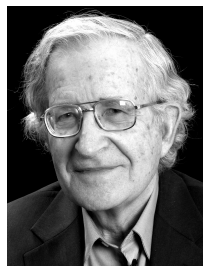
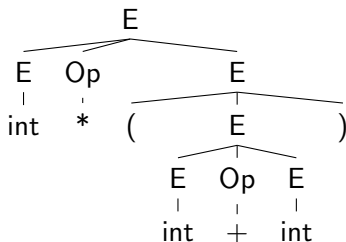
POS tagging and prediction



POS tagging and prediction



Two perspectives \approx Possible vs Probable



[...] Therefore the *true logic* for this world is the calculus of *probabilities*, which takes account of the magnitude of the probability which is, or ought to be, in a reasonable man's mind.



Linguistic structure prediction

As a structured prediction problem

- Search space: Is this analysis possible? \triangleright CFG (today)
- Measurement: Is this analysis *good*? \triangleright PCFG (today)


$$\mathbf{y}^*(\mathbf{x}; \theta) = \arg \max_{\mathbf{y} \in \mathcal{Y}(\mathbf{x})} \text{SCORE}(\mathbf{x}, \mathbf{y})$$

- Decoding: find the analysis that obtains the highest score
- Parameter estimation: find good parameters

Linguistic structure prediction

As a structured prediction problem

- Search space: Is this analysis possible? ▷ CFG (today)
- Measurement: Is this analysis *good*? ▷ PCFG (today)


$$\mathbf{y}^*(\mathbf{x}; \theta) = \underset{\mathbf{y} \in \mathcal{Y}(\mathbf{x})}{\operatorname{arg max}} \operatorname{SCORE}(\mathbf{x}, \mathbf{y})$$


- Decoding: find the analysis that obtains the highest score
- Parameter estimation: find good parameters

Linguistic structure prediction

As a structured prediction problem

- Search space: Is this analysis possible? \triangleright CFG (today)
- Measurement: Is this analysis *good*? \triangleright PCFG (today)



$$\mathbf{y}^*(\mathbf{x}; \theta) = \arg \max_{\mathbf{y} \in \mathcal{Y}(\mathbf{x})} \text{SCORE}(\mathbf{x}, \mathbf{y})$$

- Decoding: find the analysis that obtains the highest score
- Parameter estimation: find good parameters

Linguistic structure prediction

As a structured prediction problem

- Search space: Is this analysis possible? \triangleright CFG (today)
- Measurement: Is this analysis *good*? \triangleright PCFG (today)


$$\mathbf{y}^*(\mathbf{x}; \theta) = \underset{\mathbf{y} \in \mathcal{Y}(\mathbf{x})}{\text{arg max}} \text{ SCORE}(\mathbf{x}, \mathbf{y})$$


- Decoding: find the analysis that obtains the highest score
- Parameter estimation: find good parameters

Linguistic structure prediction

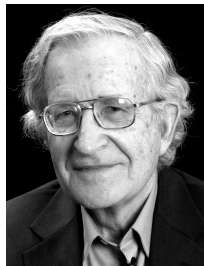
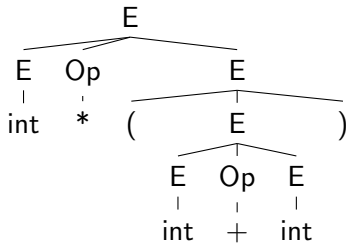
As a structured prediction problem

- Search space: Is this analysis possible? \triangleright CFG (today)
- Measurement: Is this analysis *good*? \triangleright PCFG (today)

$$\mathbf{y}^*(\mathbf{x}; \theta) = \arg \max_{\mathbf{y} \in \mathcal{Y}(\mathbf{x})} \text{SCORE}(\mathbf{x}, \mathbf{y})$$


- Decoding: find the analysis that obtains the highest score
- Parameter estimation: find good parameters

Context-Free Grammar



Formal grammars

Formally specify a grammar that can generate all and only the acceptable sentences of a natural language.

Formal grammars

Formally specify a grammar that can generate all and only the acceptable sentences of a natural language.

A grammar G consists of the following components:

1. A finite set Σ of terminal symbols.
2. A finite set N of nonterminal symbols that is disjoint from Σ .
3. A distinguished nonterminal symbol that is the `START` symbol.
4. A finite set R of production rules, each rule of the form

$$(\Sigma \cup N)^+ \rightarrow (\Sigma \cup N)^*$$

Each production rule maps from one string of symbols to another.

Context-Free Grammars

- ① N : variables
- ② Σ : terminals
- ③ R : productions

$$A \rightarrow (N \cup \Sigma)^*$$

$$A \in N$$

- ④ S : START

A linguistic example (1)

$N = \{S, NP, VP, AdjP, AdvP\} \cup \{N, Adj, Adv\}$

$\Sigma = \{colorless, green, ideas, sleep, furiously\}$

R

$S \rightarrow NP \ VP$ $VP \rightarrow VP \ AdvP$ $VP \rightarrow V$ $AdvP \rightarrow Adv$	$NP \rightarrow AdjP \ NP$ $NP \rightarrow N$ $AdjP \rightarrow Adj$
$Adj \rightarrow colorless$ $N \rightarrow ideas$ $Adv \rightarrow furiously$	$Adj \rightarrow green$ $V \rightarrow sleep$

$S = S$

A linguistic example (1)

$N = \{S, NP, VP, AdjP, AdvP\} \cup \{N, Adj, Adv\}$

$\Sigma = \{colorless, green, ideas, sleep, furiously\}$

R

$S \rightarrow NP \ VP$ $VP \rightarrow VP \ AdvP$ $VP \rightarrow V$ $AdvP \rightarrow Adv$	$NP \rightarrow AdjP \ NP$ $NP \rightarrow N$ $AdjP \rightarrow Adj$
$Adj \rightarrow colorless$ $N \rightarrow ideas$ $Adv \rightarrow furiously$	$Adj \rightarrow green$ $V \rightarrow sleep$

$S = S$

We can **derive** the structure of a string.

A linguistic example (1)

$N = \{S, NP, VP, AdjP, AdvP\} \cup \{N, Adj, Adv\}$

$\Sigma = \{colorless, green, ideas, sleep, furiously\}$

R

$S \rightarrow NP VP$ $VP \rightarrow VP AdvP$ $VP \rightarrow V$ $AdvP \rightarrow Adv$	$NP \rightarrow AdjP NP$ $NP \rightarrow N$ $AdjP \rightarrow Adj$
$Adj \rightarrow colorless$ $N \rightarrow ideas$ $Adv \rightarrow furiously$	$Adj \rightarrow green$ $V \rightarrow sleep$

$S = S$

We can **derive** the structure of a string.

$S \Rightarrow NP VP$

A linguistic example (1)

$N = \{S, NP, VP, AdjP, AdvP\} \cup \{N, Adj, Adv\}$

$\Sigma = \{colorless, green, ideas, sleep, furiously\}$

R

$S \rightarrow NP VP$ $VP \rightarrow VP AdvP$ $VP \rightarrow V$ $AdvP \rightarrow Adv$	$NP \rightarrow AdjP NP$ $NP \rightarrow N$ $AdjP \rightarrow Adj$
$Adj \rightarrow colorless$ $N \rightarrow ideas$ $Adv \rightarrow furiously$	$Adj \rightarrow green$ $V \rightarrow sleep$

$S = S$

We can **derive** the structure of a string.

$S \Rightarrow NP VP$

$\Rightarrow N VP$

A linguistic example (1)

$N = \{S, NP, VP, AdjP, AdvP\} \cup \{N, Adj, Adv\}$

$\Sigma = \{colorless, green, ideas, sleep, furiously\}$

R

$S \rightarrow NP VP$ $VP \rightarrow VP AdvP$ $VP \rightarrow V$ $AdvP \rightarrow Adv$	$NP \rightarrow AdjP NP$ $NP \rightarrow N$ $AdjP \rightarrow Adj$
$Adj \rightarrow colorless$ $N \rightarrow ideas$ $Adv \rightarrow furiously$	$Adj \rightarrow green$ $V \rightarrow sleep$

$S = S$

We can **derive** the structure of a string.

$S \Rightarrow NP VP$
 $\Rightarrow N VP$
 $\Rightarrow ideas VP$
 $\Rightarrow ideas VP AdvP$
 $\Rightarrow ideas V AdvP$
 $\Rightarrow ideas sleep AdvP$
 $\Rightarrow ideas sleep Adv$
 $\Rightarrow ideas sleep furiously$

A linguistic example (1)

$N = \{S, NP, VP, AdjP, AdvP\} \cup \{N, Adj, Adv\}$

$\Sigma = \{colorless, green, ideas, sleep, furiously\}$

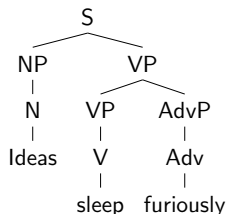
R

$S \rightarrow NP VP$	$NP \rightarrow AdjP NP$
$VP \rightarrow VP AdvP$	
$VP \rightarrow V$	$NP \rightarrow N$
$AdvP \rightarrow Adv$	$AdjP \rightarrow Adj$
$Adj \rightarrow colorless$	$Adj \rightarrow green$
$N \rightarrow ideas$	$V \rightarrow sleep$
$Adv \rightarrow furiously$	

$S = S$

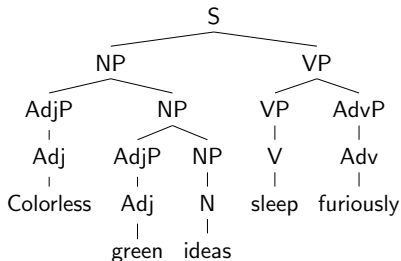
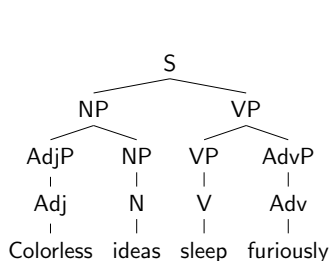
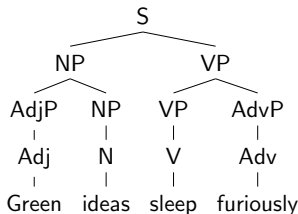
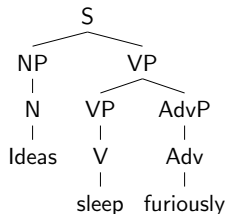
We can **derive** the structure of a string.

$S \Rightarrow NP VP$
 $\Rightarrow N VP$
 $\Rightarrow ideas VP$
 $\Rightarrow ideas VP AdvP$
 $\Rightarrow ideas V AdvP$
 $\Rightarrow ideas sleep AdvP$
 $\Rightarrow ideas sleep Adv$
 $\Rightarrow ideas sleep furiously$



A linguistic example (2)

We can define the language of a grammar by applying the productions.



Recursion (1)



from *Inception* (<https://www.imdb.com/title/tt1375666/>)

recursion

place one component inside another component of the same type

Recursion (2)

Natural numbers

- $0 \leftarrow \emptyset$
- If n is a natural number, let $n + 1 \leftarrow n \cup \{n\}$

$$0 = \emptyset$$

$$1 = \{0\} = \{\emptyset\}$$

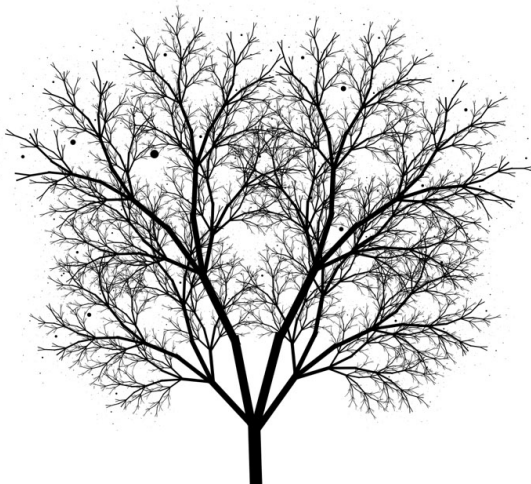
$$2 = \{0, 1\} = \{\emptyset, \{\emptyset\}\}$$

$$3 = \{0, 1, 2\} = \{\emptyset, \{\emptyset\}, \{\emptyset, \{\emptyset\}\}\}$$

recursion

place one component inside another component of the same type

Recursion (3)



<https://matthewjamestaylor.com/recursive-drawing>

Recursion (4)

We hypothesize that FLN (faculty of language in the narrow sense) only includes recursion and is the only uniquely human component of the faculty of language.

M Hauser, N Chomsky and W Fitch (2002)

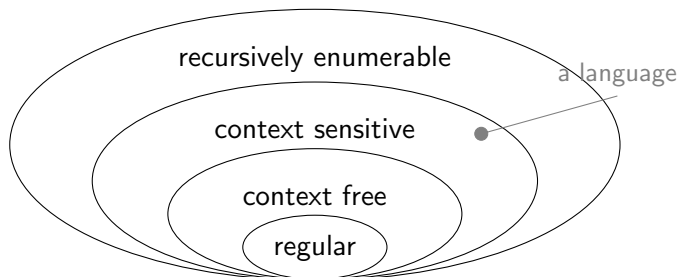
science.sciencemag.org/content/298/5598/1569

- (2) a. The dog bit the cat [which chased the mouse [which died]]. (right)
b. [[the dog] 's owner] 's friend (left)
c. The mouse [the cat [the dog bit] chased] died. (center)

Reminder: Chomsky Hierarchy

Grammar	Languages	Production rules
Type-0	Recursively enumerable	$\alpha \rightarrow \gamma$
Type-1	Context-sensitive	$\alpha A \beta \rightarrow \alpha \gamma \beta$
Type-2	Context-free	$A \rightarrow \gamma$
Type-3	Regular	$A \rightarrow a$ $A \rightarrow aB$

$$a \in N; \alpha, \beta \in (N \cup \Sigma)^*, \gamma \in (N \cup \Sigma)^+$$



Where can I get a grammar?

English Treebank

- Penn Treebank = ca. 50,000 sentences with associated trees
- Usual set-up: ca. 40,000 training sentences, ca. 2,400 test sentences
- Cut all trees into 2-level subtrees.

Probabilistic Context-Free Grammars

*[...] Therefore the **true logic** for this world is the calculus of **probabilities**, which takes account of the magnitude of the probability which is, or ought to be, in a reasonable man's mind.*



Probabilistic CFGs

Probability of a tree t with rules $A_1 \rightarrow \beta_1, A_2 \rightarrow \beta_2, \dots$ is

$$p(t) = \prod_{i=1}^n q(A_i \rightarrow \beta_i)$$

where $q(A_i \rightarrow \beta_i)$ is the probability for rule $A_i \rightarrow \beta_i$.

- When we expand A_i , how likely is it that we choose $A_i \rightarrow \beta_i$?
- For each nonterminal A_i ,

$$\sum_{\beta} q(A_i \rightarrow \beta | A_i) = 1$$

- PCFG generates random derivations of CFG.
- Each event (expanding nonterminal by production rules) is statistically independent of all the others.

An example (1)

S	→	NP VP	0.8
S	→	Aux NP VP	0.15
S	→	VP	0.05
NP	→	AdjP NP	0.2
NP	→	D N	0.7
NP	→	N	0.1
VP	→	VP AdvP	0.3
VP	→	V	0.2
VP	→	V NP	0.3
VP	→	V NP NP	0.2
AdvP	→	Adv	1.0
AdjP	→	Adj	1.0

Adj	→	<i>colorless</i>	0.4
Adj	→	<i>green</i>	0.6
N	→	<i>ideas</i>	1.0
V	→	<i>sleep</i>	1.0
Adv	→	<i>furiously</i>	1.0

An example (2)

S

$S \rightarrow NP VP$

0.8

An example (2)

\Rightarrow S
NP VP

S \rightarrow NP VP 0.8
NP \rightarrow N 0.1

An example (2)

	S	$S \rightarrow NP\ VP$	0.8
\Rightarrow	NP VP	$NP \rightarrow N$	0.1
\Rightarrow	N VP	$N \rightarrow ideas$	1.0
\Rightarrow	ideas VP	$VP \rightarrow VP\ AdvP$	0.3
\Rightarrow	ideas VP AdvP	$VP \rightarrow V$	0.2
\Rightarrow	ideas V AdvP	$V \rightarrow sleep$	1.0
\Rightarrow	ideas sleep AdvP	$AdvP \rightarrow Adv$	1.0
\Rightarrow	ideas sleep Adv	$Adv \rightarrow furiously$	1.0

An example (2)

	S	$S \rightarrow NP VP$	0.8
\Rightarrow	NP VP	$NP \rightarrow N$	0.1
\Rightarrow	N VP	$N \rightarrow ideas$	1.0
\Rightarrow	ideas VP	$VP \rightarrow VP AdvP$	0.3
\Rightarrow	ideas VP AdvP	$VP \rightarrow V$	0.2
\Rightarrow	ideas V AdvP	$V \rightarrow sleep$	1.0
\Rightarrow	ideas sleep AdvP	$AdvP \rightarrow Adv$	1.0
\Rightarrow	ideas sleep Adv	$Adv \rightarrow furiously$	1.0

$$0.8 \times 0.1 \times 1.0 \times 0.3 \times 0.2 \times 1.0 \times 1.0 \times 1.0$$

Properties of PCFGs

- Assigns a probability to each parse-tree, allowed by the underlying CFG
- Say we have a sentence s , set of derivations for that sentence is $\mathcal{T}(s)$, as defined by a CFG. Then a PCFG assigns a probability $p(t)$ to each member of $\mathcal{T}(s)$.
- We now have a SCORE function (probability) that can ranks trees.
- The most likely parse tree for a sentence s is

$$\arg \max_{t \in \mathcal{T}(s)} p(t)$$

“correct” means more probable parse tree

“language” means set of grammatical sentences

Deriving a PCFG from a Treebank

Given a set of example trees (a treebank), the underlying CFG can simply be all rules seen in the corpus

Maximum Likelihood Estimates

$$q_{ML}(\alpha \rightarrow \beta) = \frac{\text{COUNT}(\alpha \rightarrow \beta)}{\text{COUNT}(\alpha)}$$

The counts are taken from a training set of example trees.

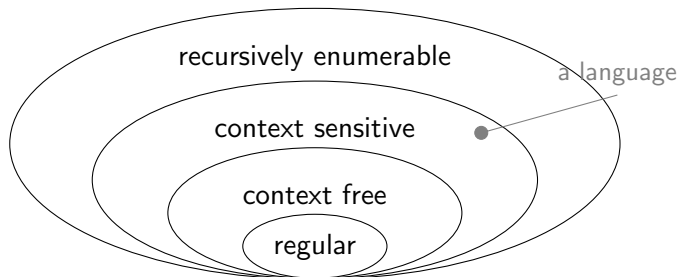
If the training data is generated by a PCFG, then as the training data size goes to infinity, the maximum-likelihood PCFG will converge to the same distribution as the “true” PCFG.

Rethink Part-of-Speech Tagging

Chomsky Hierarchy

Grammar	Languages	Production rules
Type-0	Recursively enumerable	$\alpha \rightarrow \gamma$
Type-1	Context-sensitive	$\alpha A \beta \rightarrow \alpha \gamma \beta$
Type-2	Context-free	$A \rightarrow \gamma$
Type-3	Regular	$A \rightarrow a$ $A \rightarrow aB$

$a \in N; \alpha, \beta, \gamma \in (N \cup \Sigma)^*$



An example

Max bit the cat [which chased the mouse [which died]].

A toy grammar

- $VP \rightarrow \text{bit} | \text{chased} | \dots DP$
- $VP \rightarrow \text{died}$
- $DP \rightarrow \text{the} | \text{a} | \text{this} | \dots NP$
- $NP \rightarrow \text{dog} | \text{cat} | \text{mouse} | \dots RC$
- $RC \rightarrow \text{which} | \text{that} | \dots VP$

An example

Max bit the cat [which chased the mouse [which died]].

A toy grammar

- $VP \rightarrow \text{bit} | \text{chased} | \dots DP$
- $VP \rightarrow \text{died}$
- $DP \rightarrow \text{the} | \text{a} | \text{this} | \dots NP$
- $NP \rightarrow \text{dog} | \text{cat} | \text{mouse} | \dots RC$
- $RC \rightarrow \text{which} | \text{that} | \dots VP$

VP

An example

Max **bit** the cat [which chased the mouse [which died]].

A toy grammar

- **VP** \rightarrow **bit** | chased | ... **DP**
- **VP** \rightarrow died
- **DP** \rightarrow the | a | this | ... **NP**
- **NP** \rightarrow dog | cat | mouse | ... **RC**
- **RC** \rightarrow which | that | ... **VP**

VP

\Rightarrow *bit* DP

An example

Max bit the cat [which chased the mouse [which died]].

A toy grammar

- $VP \rightarrow \text{bit} | \text{chased} | \dots DP$
- $VP \rightarrow \text{died}$
- $DP \rightarrow \text{the} | \text{a} | \text{this} | \dots NP$
- $NP \rightarrow \text{dog} | \text{cat} | \text{mouse} | \dots RC$
- $RC \rightarrow \text{which} | \text{that} | \dots VP$

VP

$\Rightarrow \text{bit } \underline{DP} \Rightarrow \text{bit the } \underline{NP}$

An example

Max **bit the cat [which chased the mouse [which died]]**.

A toy grammar

- $VP \rightarrow \text{bit} | \text{chased} | \dots DP$
- $VP \rightarrow \text{died}$
- $DP \rightarrow \text{the} | \text{a} | \text{this} | \dots NP$
- $NP \rightarrow \text{dog} | \text{cat} | \text{mouse} | \dots RC$
- $RC \rightarrow \text{which} | \text{that} | \dots VP$

VP

$\Rightarrow \text{bit } \underline{DP} \Rightarrow \text{bit } \text{the } \underline{NP} \Rightarrow \text{bit the cat } \underline{RC}$

An example

Max bit the cat [which chased the mouse [which died]].

A toy grammar

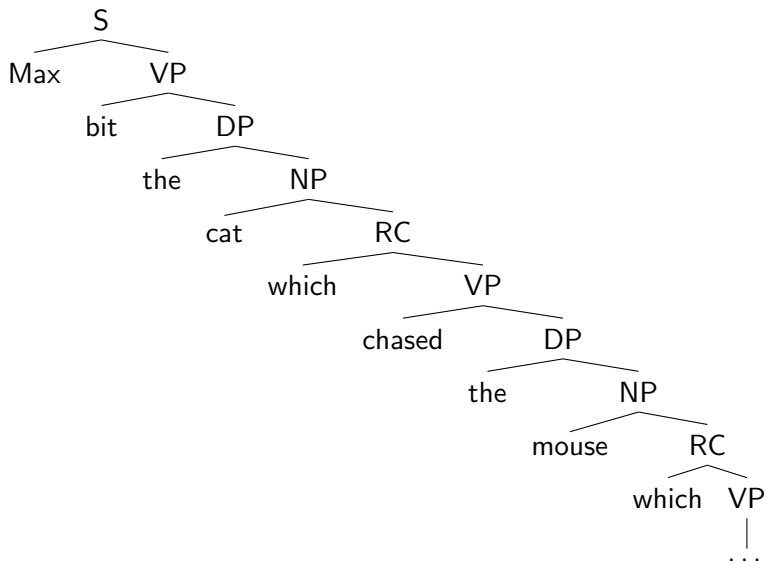
- $VP \rightarrow \text{bit} | \text{chased} | \dots DP$
- $VP \rightarrow \text{died}$
- $DP \rightarrow \text{the} | \text{a} | \text{this} | \dots NP$
- $NP \rightarrow \text{dog} | \text{cat} | \text{mouse} | \dots RC$
- $RC \rightarrow \text{which} | \text{that} | \dots VP$

VP

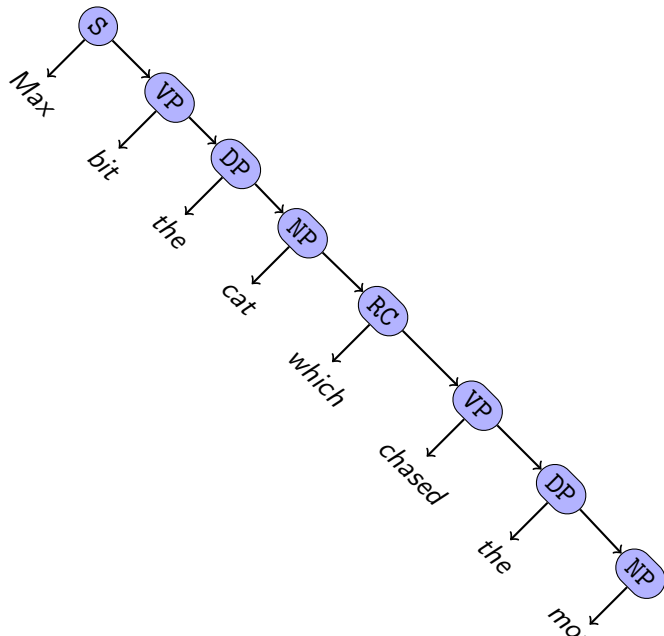
$\Rightarrow \text{bit } \underline{DP} \Rightarrow \text{bit } \text{the } \underline{NP} \Rightarrow \text{bit the cat } \underline{RC}$

$\Rightarrow \text{bit the cat } \text{which } \underline{VP}$

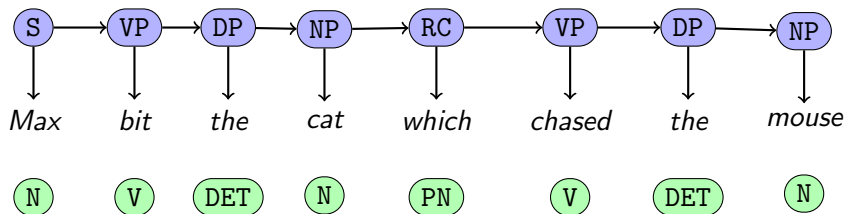
Finite state machines?



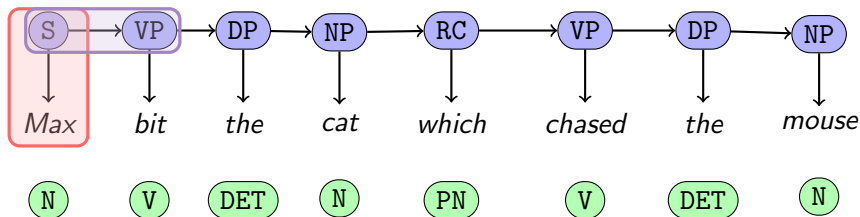
Finite state machines?



Word tagging is very powerful



Word tagging is very powerful



Generative models: Hidden Markov Models and PCFG

- $p_e(\text{Max}|\text{S}) \times p_t(\text{VP}|\text{S})$
- $p(\text{S} \rightarrow \text{Max VP})$

Probabilistic models for sequence pairs

- We have two sequences of random variables: X_1, X_2, \dots, X_n and S_1, S_2, \dots, S_n
- Intuitively, each X_i corresponds to an **observation** and each S_i corresponds to an underlying **state** that generated the observation. Assume that each S_i is in $\{1, 2, \dots, k\}$, and each X_i is in $\{1, 2, \dots, o\}$.
- How do we model the joint distribution

$$P(X_1 = x_1, \dots, X_n = x_n, S_1 = s_1, \dots, S_n = s_n)$$

Hidden Markov Models

An HMM takes the following form

$$p(x_1 \dots x_n, s_1 \dots s_n; \theta) = p_t(s_1) \prod_{j=2}^n p_t(s_j | s_{j-1}) \prod_{j=1}^n p_e(x_j | s_j)$$

Parameters in the model

- 1 Initial state parameters ϕ_s for $s \in \{1, 2, \dots, k\}$
- 2 Transition parameters $\phi_{s'|s}$ for $s, s' \in \{1, 2, \dots, k\}$
- 3 Emission parameters $\phi_{e|s}$ for $s \in \{1, 2, \dots, k\}$ and $e \in \{1, 2, \dots, o\}$

If we use a specific symbol to denote *stop of a sequence*: $s_0 = *$

- Initial state parameters $\phi_{s|*}$
- Just look like transition parameters

Harmonic word order

Morphology

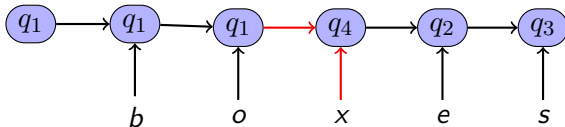
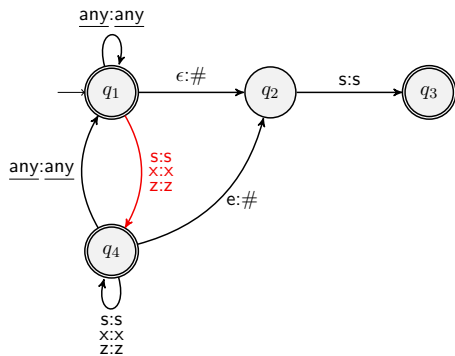
- Postpositional and head-final languages use suffixes and no prefixes.
- Prepositional and head-initial languages use not only prefixes but also suffixes.

Greenberg's word order universals

- Universal 3: Languages with dominant VSO order are always prepositional.
- Universal 4: With overwhelmingly greater than chance frequency, languages with normal SOV order are postpositional.
- Universal 5: If a language has dominant SOV order and the genitive follows the governing noun, then the adjective likewise follows the noun.
- Universal 17: With overwhelmingly more than chance frequency, languages with dominant order VSO have the adjective after the noun.

Empirical data can be found at <https://wals.info>.

Connection to Finite State Machines



Mildly Context-Sensitive Languages

Challenge

Cross-serial dependencies in Swiss German

... das mer em Hans es huus hääfed aastrüiche

... that we Hans_{Dat} house_{Acc} help paint

... that we helped Hans paint the house

... das mer d'chind em Hans es huus lönd hääfe aastrüiche

... that we the children_{Acc} Hans_{Dat} house_{Acc} let help paint

... that we let the children help Hans paint the house

Cross-serial dependencies in Dutch

... dat Wim Jan Marie de kinderen zag helpen leren zwemmen

... that Wim Jan Marie the children saw help teach swim

... that Wim saw Jan help Marie teach the children to swim

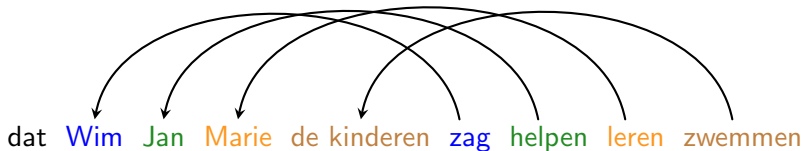
Cross-serial dependencies

Cross-serial dependencies in Dutch

...dat Wim Jan Marie de kinderen zag helpen leren zwemmen

...that Wim Jan Marie the children saw help teach swim

...that Wim saw Jan help Marie teach the children to swim

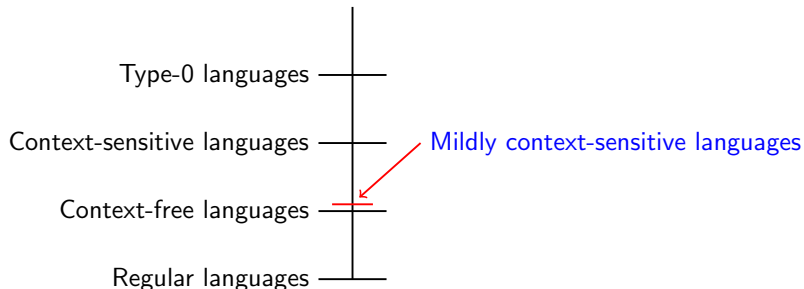


Mildly Context-Sensitive Languages

With a *possibility* perspective

Natural languages are provably **non-context-free**.

Natural languages = mildly context-sensitive languages?

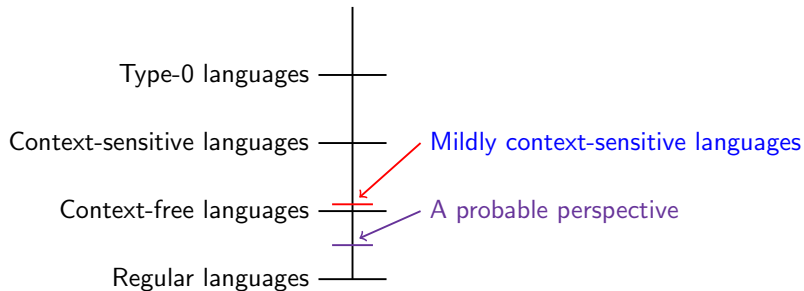


Mildly Context-Sensitive Languages

With a *possibility* perspective

Natural languages are provably **non-context-free**.

Natural languages = mildly context-sensitive languages?



Reading

D Jurafsky and J Martin. *Speech and Language Processing*.

- §17.1–§17.5, and §17.8. Context-free Grammars and Constituency Parsing. *Speech and Language Processing*. D Jurafsky and J Martin.
<https://web.stanford.edu/~jurafsky/slp3/17.pdf>
- §18.1 and §18.4. Dependency Parsing. *Speech and Language Processing*. D Jurafsky and J Martin.
<https://web.stanford.edu/~jurafsky/slp3/18.pdf>