

# Overview of Natural Language Processing

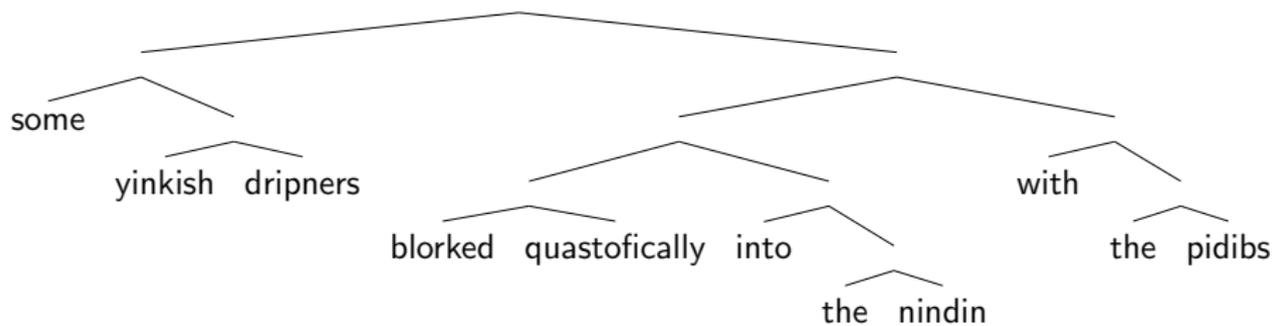
## Part II & ACS L390

### Lecture 5: Phrase Structure and Structured Prediction

Weiwei Sun

Department of Computer Science and Technology  
University of Cambridge

Michaelmas 2024/25



Words are organized into nested blocks

## Lecture 5: Phrase Structure and Structured Prediction

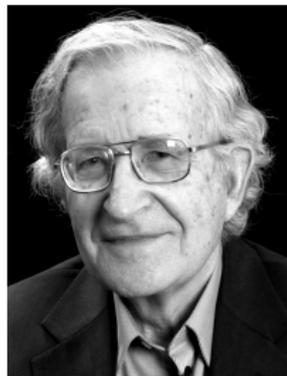
1. Phrase structure
3. Structured prediction
4. Probabilistic Context-free grammars
5. Neural parameterisation

# 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.*



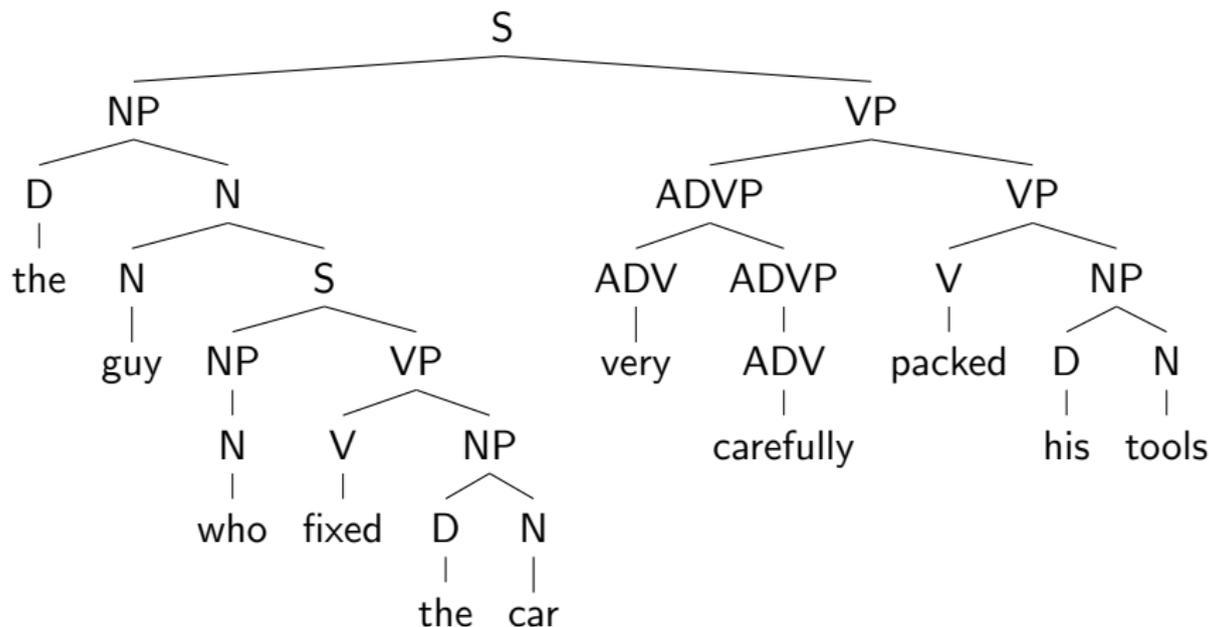
Noam Chomsky: Language, Cognition, and Deep Learning

[www.youtube.com/watch?v=cMscNuSUy0I](https://www.youtube.com/watch?v=cMscNuSUy0I)

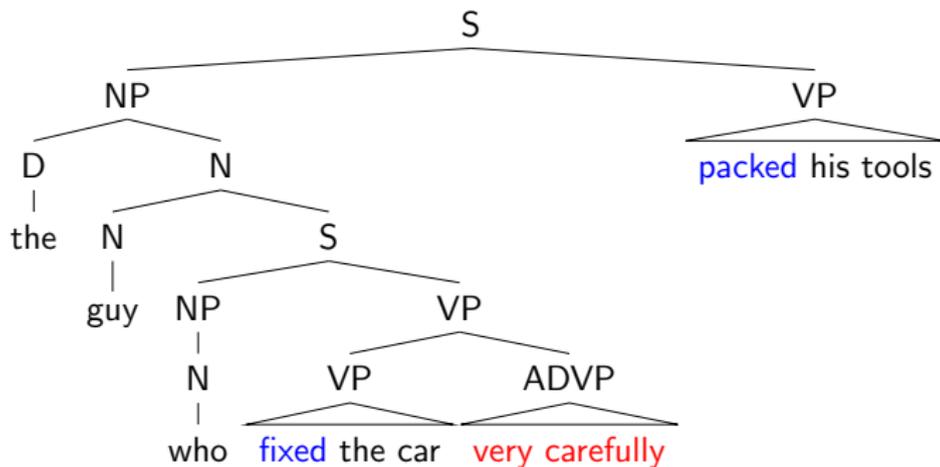
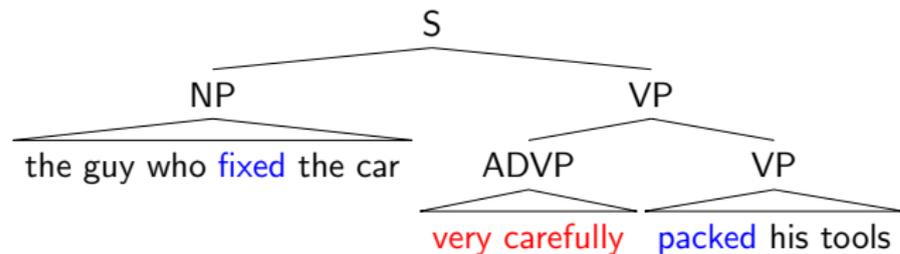
# Constituency (phrase structure)

## The basic idea

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

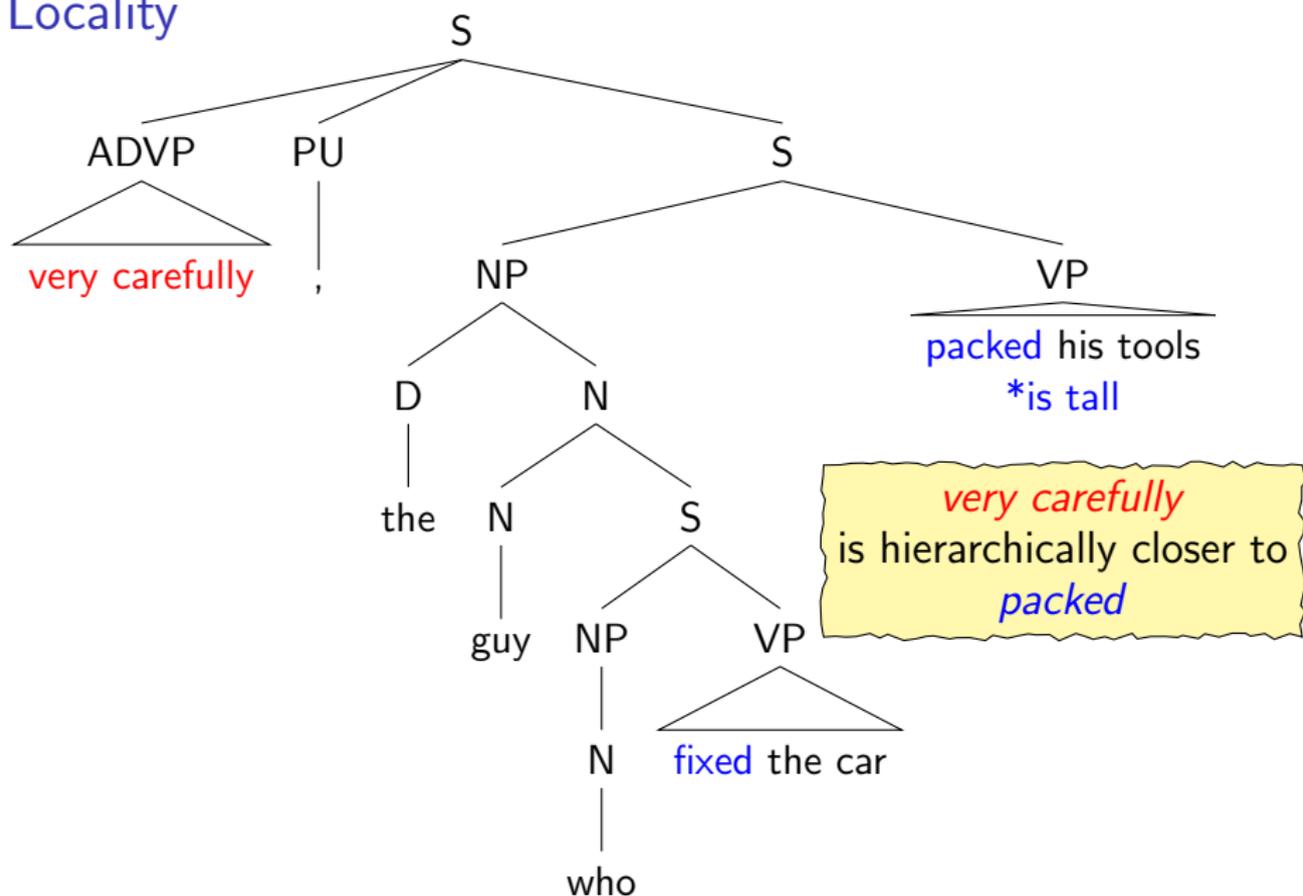


## Different structures, different meaning



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

# Locality



# 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.

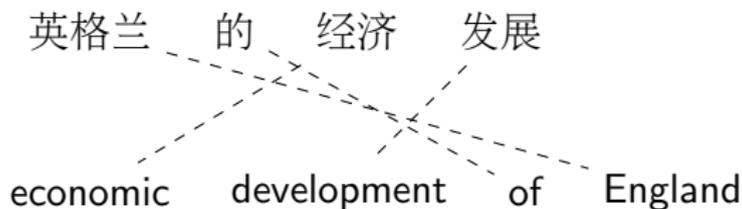
# Applications of parsing

## 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



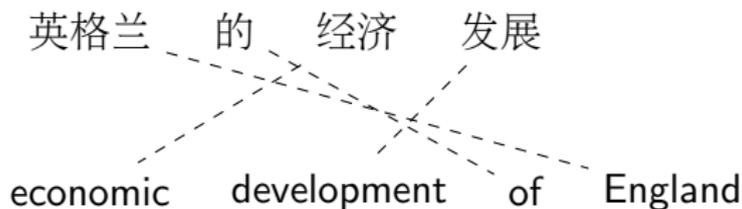
# Applications of parsing

## Modern parsers are quite accurate

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

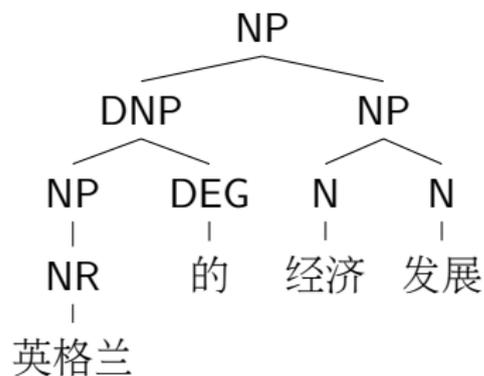
- Machine translation
- Information extraction
- Grammar checking
- etc.

Translate “英格兰的经济发展” into English



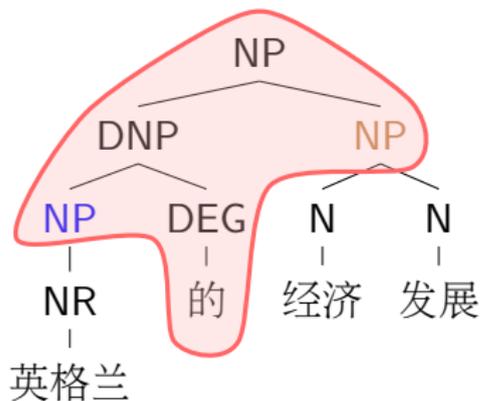
finite-state transducer?

# An example: Tree-to-string transduction



R1	<pre>           graph TD             NP1[NP] --- NR[NR]             NR --- EN[英格兰]           </pre>	<pre>           graph TD             X1[X] --- EN[England]           </pre>
R2	<pre>           graph TD             NP1[NP] --- DNP[DNP]             NP1 --- NP2[NP]             DNP --- NP3[NP]             DNP --- DEG1[的]             NP3 --- NP4[NP]             NP3 --- DEG2[的]             NP4 --- EN[英格兰]           </pre>	<pre>           graph TD             X1[X] --- X2[X]             X1 --- OF[of]             X1 --- X3[X]           </pre>
R3	<pre>           graph TD             NP1[NP] --- N1[N]             NP1 --- N2[N]             N1 --- JI[经济]             N2 --- FA[发展]           </pre>	<pre>           graph TD             X1[X] --- X2[X]             X1 --- X3[X]           </pre>
R4	<pre>           graph TD             N1[N] --- JI[经济]           </pre>	<pre>           graph TD             X1[X] --- ECON[economic]           </pre>
R5	<pre>           graph TD             N1[N] --- FA[发展]           </pre>	<pre>           graph TD             X1[X] --- DEV[development]           </pre>

# An example: Tree-to-string transduction



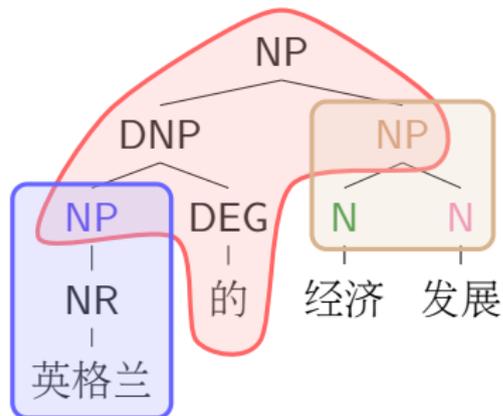
NP of NP

▷R2

R1	NP   NR   英格兰	X   England
R2	NP / \ DNP NP / \ / \ NP DEG     的 经济 发展	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

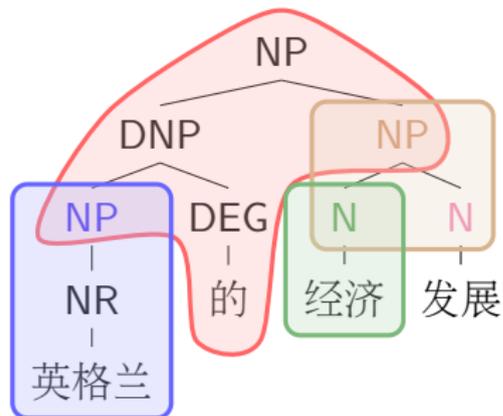


NP of NP  
 NP of England  
 N N of England

▷R2  
 ▷R1  
 ▷R3

R1	<pre>           NP                       NR                     英格兰           </pre>	<pre>           X                     England           </pre>
R2	<pre>           NP          /  \         DNP  NP        /  \       NP  DEG                   的           </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

▷R2

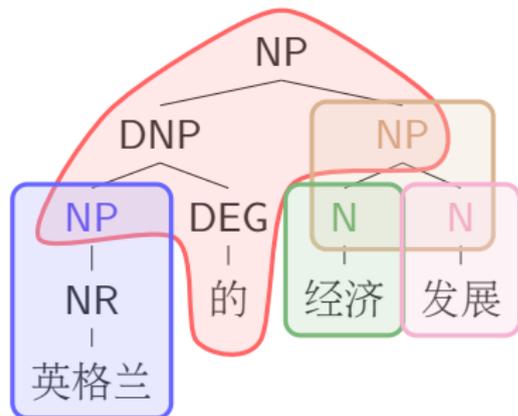
▷R1

▷R3

▷R4

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>

# 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

R1	<pre>           NP                       NR                     英格兰                     </pre>	<pre>           X                     England                     </pre>
R2	<pre>           NP          /  \         DNP  NP        /  \       NP  DEG                       的                     </pre>	<pre>           X          /   \         X of X                     </pre>
R3	<pre>           NP          /  \         N   N   </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

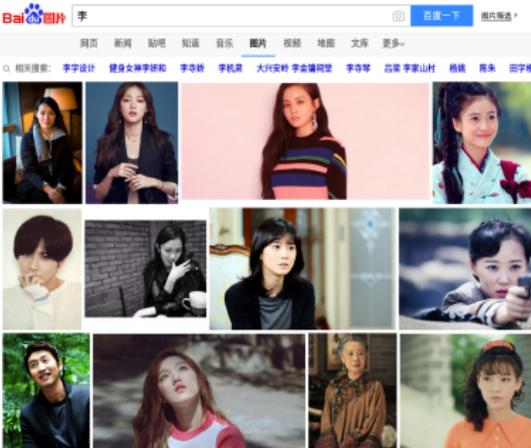
pre-lecture: watch this video

▶ [www.youtube.com/watch?v=bjUwSHGsG9o](https://www.youtube.com/watch?v=bjUwSHGsG9o)

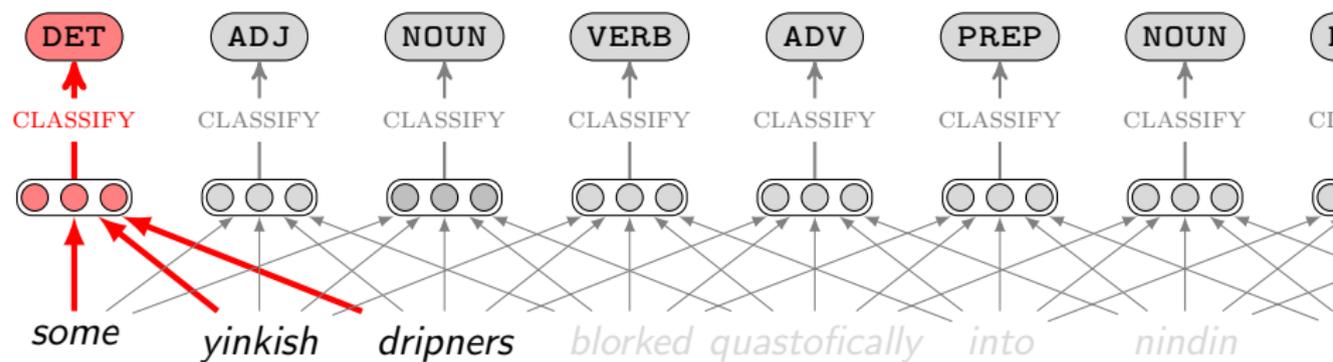
# Muhammad Li

Howard Who's Muhammad Li?

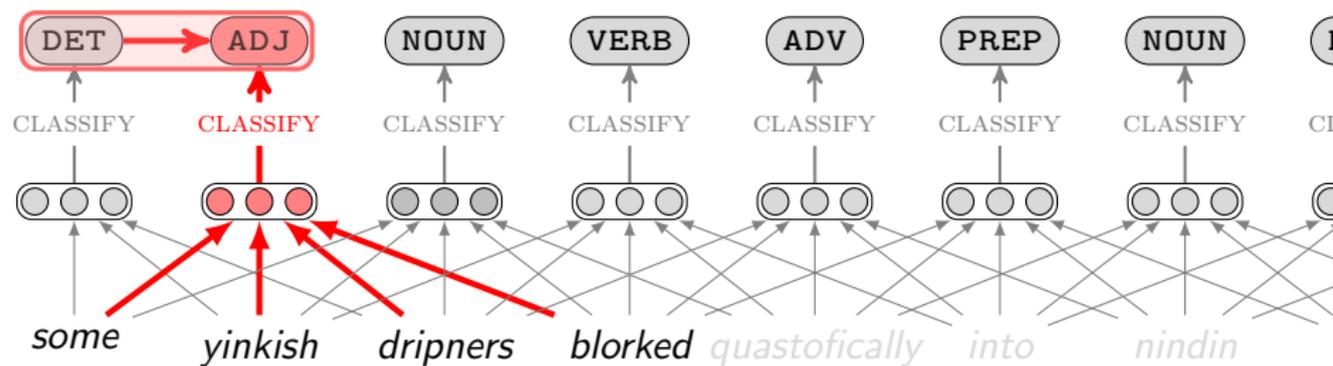
Sheldon Muhammad is the most common first name in the world, Li, the most common surname. As I didn't know the answer, I thought that gave me a mathematical edge.



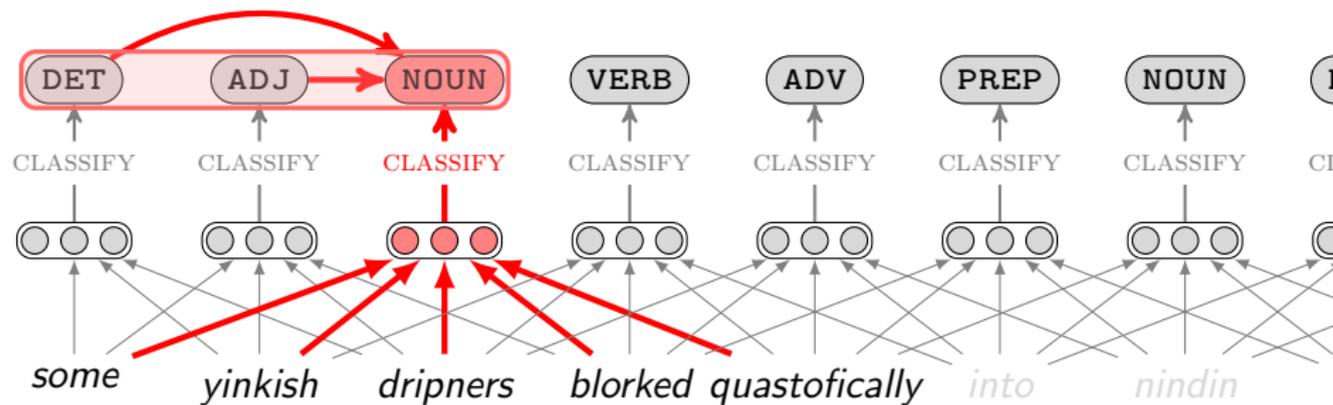
# POS tagging and prediction



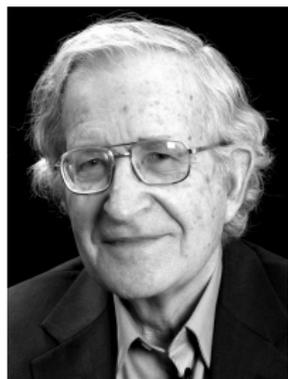
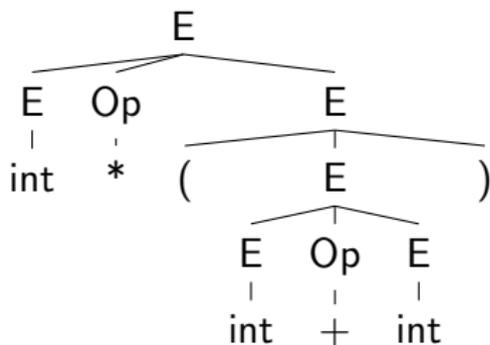
# POS tagging and prediction



# 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? ▷CFG (today)
- Measurement: Is this analysis *good*? ▷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)

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


- 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? ▷ CFG (today)
- Measurement: Is this analysis *good*? ▷ 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? ▷CFG (today)
- Measurement: Is this analysis *good*? ▷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

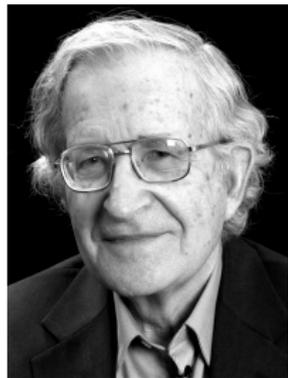
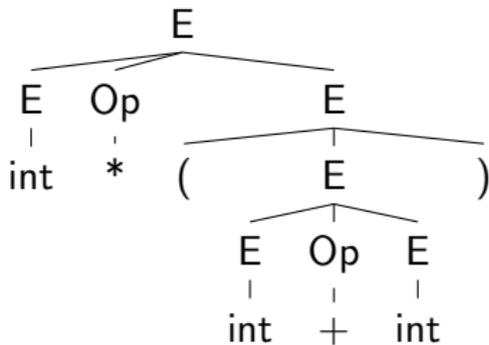
## As a structured prediction problem

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

# 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  $\Sigma$  of terminal symbols.
2. A finite set  $N$  of nonterminal symbols that is disjoint from  $\Sigma$ .
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  $\Sigma$ : terminals
- 3  $R$ : productions

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

$$A \in N$$

- 4  $S$ : START

# A linguistic example (1)

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

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

$R$

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

$S = S$

# 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.

# A linguistic example (1)

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

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

$R$

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

$S = S$

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

$S \Rightarrow NP VP$

## A linguistic example (1)

$$N = \{S, NP, VP, AdjP, AdvP\} \cup \{N, Adj, Adv\}$$
$$\Sigma = \{\textit{colorless}, \textit{green}, \textit{ideas}, \textit{sleep}, \textit{furiously}\}$$

$R$

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

$$S = S$$

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

$$S \Rightarrow NP VP$$
$$\Rightarrow N VP$$

# A linguistic example (1)

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

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

$R$

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

$S = S$

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

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

# A linguistic example (1)

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

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

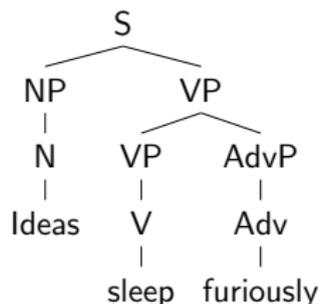
$R$

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

$S = S$

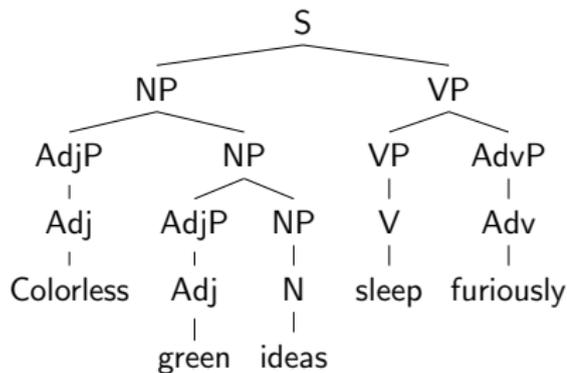
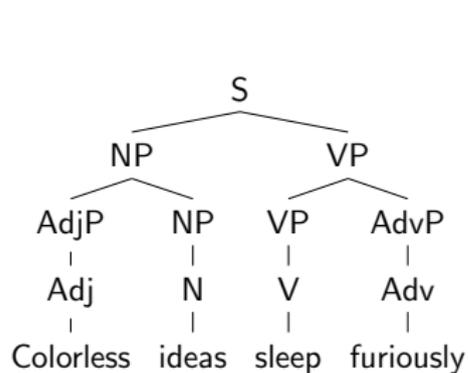
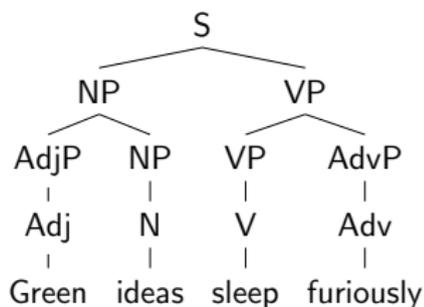
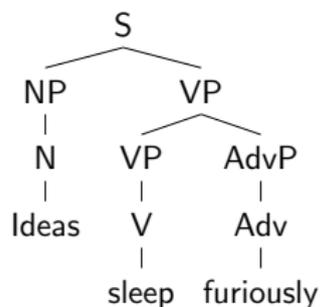
We can **derive** the structure of a string.

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



## 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  $n + 1 \leftarrow n \cup \{n\}$

$$0 = \emptyset$$

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

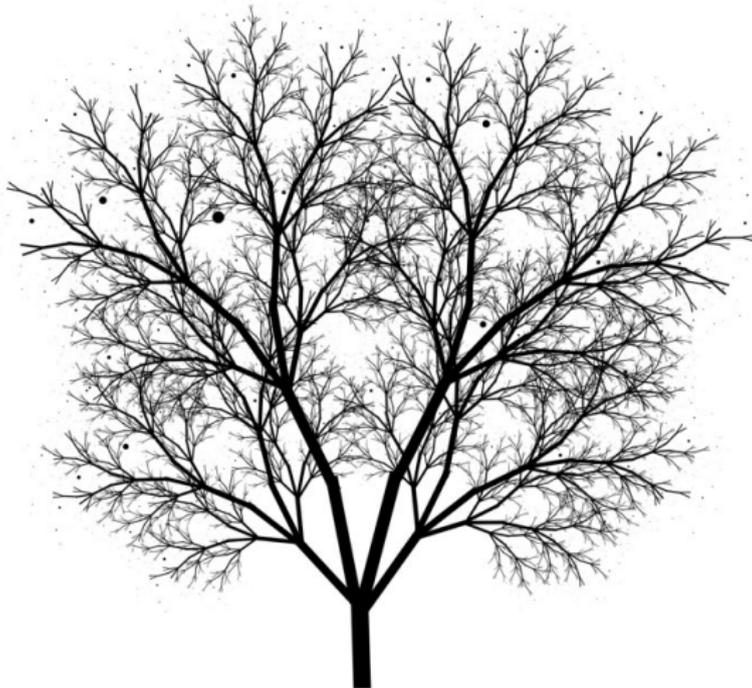
$$2 = \{0, 1\} = \{\emptyset, \{\emptyset\}\}$$

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

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)*

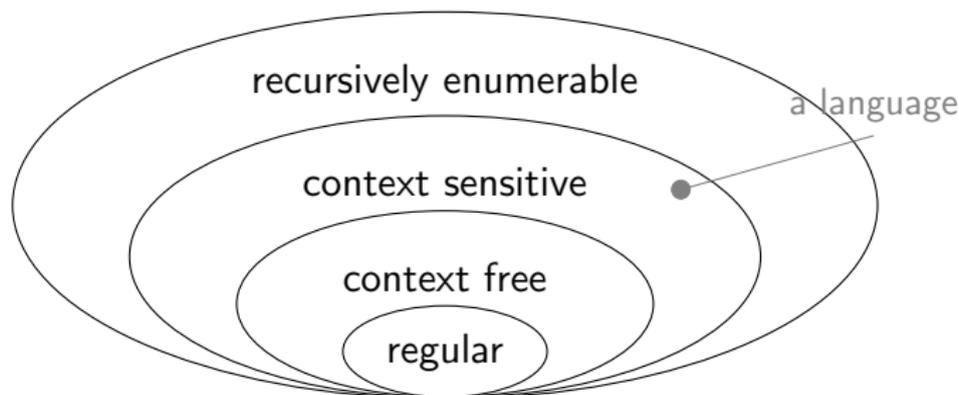
*[science.sciencemag.org/content/298/5598/1569](http://science.sciencemag.org/content/298/5598/1569)*

- (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)

	S	S → NP VP	0.8
⇒	NP VP	NP → N	0.1

## An example (2)

	S	S→NP VP	0.8
⇒	NP VP	NP→N	0.1
⇒	N VP	N→ <i>ideas</i>	1.0
⇒	ideas VP	VP→VP AdvP	0.3
⇒	ideas VP AdvP	VP→V	0.2
⇒	ideas V AdvP	V→ <i>sleep</i>	1.0
⇒	ideas sleep AdvP	AdvP→Adv	1.0
⇒	ideas sleep Adv	Adv→ <i>furiously</i>	1.0

## An example (2)

	S	S→NP VP	0.8
⇒	NP VP	NP→N	0.1
⇒	N VP	N→ <i>ideas</i>	1.0
⇒	ideas VP	VP→VP AdvP	0.3
⇒	ideas VP AdvP	VP→V	0.2
⇒	ideas V AdvP	V→ <i>sleep</i>	1.0
⇒	ideas sleep AdvP	AdvP→Adv	1.0
⇒	ideas sleep Adv	Adv→ <i>furiously</i>	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

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

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.

# Discriminative vs Generative

We have learned two probabilistic models:

- Log-linear:  $P(Y|X)$
- PCFG:  $P(X, Y)$
  
- Generative models can generate new data instances. It includes the distribution of the data itself, and tells you how likely a given example is.
- Discriminative models discriminate between different kinds of data instances. A discriminative model just tells you how likely a label is for a given instance.

# Discriminative vs Generative

We have learned two probabilistic models:

- Log-linear:  $P(Y|X)$  an example of discriminative model
- PCFG:  $P(X, Y)$  an example of generative model
  
- Generative models can generate new data instances. It includes the distribution of the data itself, and tells you how likely a given example is.
- Discriminative models discriminate between different kinds of data instances. A discriminative model just tells you how likely a label is for a given instance.

# Neural Parameterisation

## Neural parameterisation

- Simple parameterisation of PCFG:  $q(A_i \rightarrow \beta_i)$  is a real number.
- Alternative neural parameterisation:  $q(A_i \rightarrow \beta_i)$  is a trainable NN.

Kim et al. (2019)

<https://aclanthology.org/P19-1228.pdf>

$$q(A \rightarrow BC) = \frac{\exp(\mathbf{u}_{BC}^\top \mathbf{w}_A)}{\sum_{B'C'} \exp(\mathbf{u}_{B'C'}^\top \mathbf{w}_A)}$$

- $A \in \Sigma, B, C \in \Sigma \cup N$
- $\{\mathbf{w}_x | x \in \Sigma \cup N\}$  and  $\{\mathbf{u}_{xy} | x, y \in \Sigma \cup N\}$ : the set of input symbol embeddings for a grammar.

## Compound PCFG

$$q(A \rightarrow BC; \mathbf{z}) = \frac{\exp(\mathbf{u}_{BC}^\top [\mathbf{w}_A; \mathbf{z}])}{\sum_{B'C'} \exp(\mathbf{u}_{B'C'}^\top [\mathbf{w}_A; \mathbf{z}])}$$

# Grammar induction

- Children acquire their grammars in a more or less **unsupervised** manner.
- Grammar induction is a research task in computational linguistics with the following research question: to what extent can the structure of human language be distributionally identified?
- Recent progress: Applying neural parameterisation to PCFG and estimate parameters in an unsupervised way.

## 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>