

# *Syntax - Introduction*

Pawan Goyal

CSE, IIT Kharagpur

Week 5: Lecture 1

# *What is Syntax?*

# *What is Syntax?*

- Refers to the way words are arranged together, and the relationship between them.

# What is Syntax?

- Refers to the way words are arranged together, and the relationship between them.
- **Language Models:** Importance of modeling word order

# What is Syntax?

- Refers to the way words are arranged together, and the relationship between them.
- **Language Models:** Importance of modeling word order
- **POS categories:** An equivalence class for words

# What is Syntax?

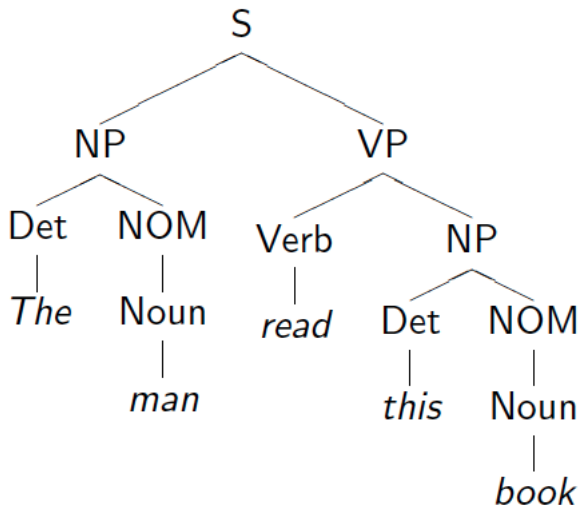
- Refers to the way words are arranged together, and the relationship between them.
- **Language Models:** Importance of modeling word order
- **POS categories:** An equivalence class for words
- More complex notions: constituency, grammatical relations, subcategorization etc.

# What is Syntax?

- Refers to the way words are arranged together, and the relationship between them.
- **Language Models:** Importance of modeling word order
- **POS categories:** An equivalence class for words
- More complex notions: constituency, grammatical relations, subcategorization etc.



# Syntax Tree: Example





# Defining the notions: Constituency

## *Constituent*

A group of words acts as a single unit - phrases, clauses etc.

# Defining the notions: Constituency

## *Constituent*

A group of words acts as a single unit - phrases, clauses etc.

## *Part of Speech - "Substitution Test"*

The {sad, intelligent, green, fat, ...} one is in the corner.

# Defining the notions: Constituency

## *Constituent*

A group of words acts as a single unit - phrases, clauses etc.

## *Part of Speech - "Substitution Test"*

The {sad, intelligent, green, fat, ...} one is in the corner.

## *Constituency: Noun Phrase*

- *Kermit the frog*
- *they*
- *December twenty-sixth*
- *the reason he is running for president*

# Constituent Phrases

Usually named based on the word that heads the constituent:

<i>the man from Amherst</i>	is a Noun Phrase (NP) because the head <i>man</i> is a noun
<i>extremely clever</i>	is an Adjective Phrase (AP) because the head <i>clever</i> is an adjective
<i>down the river</i>	is a Prepositional Phrase (PP) because the head <i>down</i> is a preposition
<i>killed the rabbit</i>	is a Verb Phrase (VP) because the head <i>killed</i> is a verb

# Constituent Phrases

Usually named based on the word that heads the constituent:

<i>the man from Amherst</i>	is a Noun Phrase (NP) because the head <i>man</i> is a noun
<i>extremely clever</i>	is an Adjective Phrase (AP) because the head <i>clever</i> is an adjective
<i>down the river</i>	is a Prepositional Phrase (PP) because the head <i>down</i> is a preposition
<i>killed the rabbit</i>	is a Verb Phrase (VP) because the head <i>killed</i> is a verb

*Words can also act as phrases*

*Joe grew potatoes*

# Constituent Phrases

Usually named based on the word that heads the constituent:

<i>the man from Amherst</i>	is a Noun Phrase (NP) because the head <i>man</i> is a noun
<i>extremely clever</i>	is an Adjective Phrase (AP) because the head <i>clever</i> is an adjective
<i>down the river</i>	is a Prepositional Phrase (PP) because the head <i>down</i> is a preposition
<i>killed the rabbit</i>	is a Verb Phrase (VP) because the head <i>killed</i> is a verb

## *Words can also act as phrases*

*Joe grew potatoes*

*Joe* and *potatoes* are both nouns and noun phrases

# Constituent Phrases

Usually named based on the word that heads the constituent:

<i>the man from Amherst</i>	is a Noun Phrase (NP) because the head <i>man</i> is a noun
<i>extremely clever</i>	is an Adjective Phrase (AP) because the head <i>clever</i> is an adjective
<i>down the river</i>	is a Prepositional Phrase (PP) because the head <i>down</i> is a preposition
<i>killed the rabbit</i>	is a Verb Phrase (VP) because the head <i>killed</i> is a verb

## *Words can also act as phrases*

*Joe grew potatoes*

*Joe* and *potatoes* are both nouns and noun phrases

Compare with: *The man from Amherst grew beautiful russet potatoes.*

# Constituent Phrases

Usually named based on the word that heads the constituent:

<i>the man from Amherst</i>	is a Noun Phrase (NP) because the head man is a noun
<i>extremely clever</i>	is an Adjective Phrase (AP) because the head clever is an adjective
<i>down the river</i>	is a Prepositional Phrase (PP) because the head down is a preposition
<i>killed the rabbit</i>	is a Verb Phrase (VP) because the head killed is a verb

## *Words can also act as phrases*

*Joe grew potatoes*

*Joe* and *potatoes* are both nouns and noun phrases

Compare with: *The man from Amherst grew beautiful russet potatoes.*

*Joe* appears in a place that a larger noun phrase could have been.



# *Evidence that constituency exists*

*They appear in similar environments*

# Evidence that constituency exists

## They appear in similar environments

Kermit the frog comes on stage

They come to Massachusetts every summer

December twenty-sixth comes after Christmas

The reason he is running for president comes out only now.

But not each individual word in the constituent

\*The comes out... \*is comes out... \*for comes out...

# Evidence that constituency exists

## They appear in similar environments

Kermit the frog comes on stage

They come to Massachusetts every summer

December twenty-sixth comes after Christmas

The reason he is running for president comes out only now.

But not each individual word in the constituent

\*The comes out... \*is comes out... \*for comes out...

## Can be placed in a number of different locations

# Evidence that constituency exists

## They appear in similar environments

Kermit the frog comes on stage

They come to Massachusetts every summer

December twenty-sixth comes after Christmas

The reason he is running for president comes out only now.

But not each individual word in the constituent

\*The comes out... \*is comes out... \*for comes out...

## Can be placed in a number of different locations

Constituent = Prepositional phrase: On December twenty-sixth

On December twenty-sixth I'd like to fly to Florida.

I'd like to fly on December twenty-sixth to Florida.

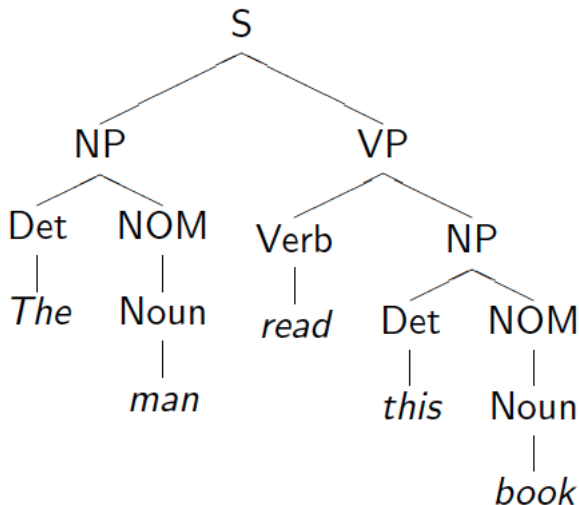
I'd like to fly to Florida on December twenty-sixth.

But not split apart

\*On December I'd like to fly twenty-sixth to Florida.

\*On I'd like to fly December twenty-sixth to Florida.

# Modeling Constituency: what tool do we need?



# *Modeling Constituency*

# Modeling Constituency

## *Context-free grammar*

The most common way of modeling constituency

# Modeling Constituency

## *Context-free grammar*

The most common way of modeling constituency

## *Consists of production Rules*

These rules express the ways in which the symbols of the language can be grouped and ordered together



# Modeling Constituency

## *Context-free grammar*

The most common way of modeling constituency

## *Consists of production Rules*

These rules express the ways in which the symbols of the language can be grouped and ordered together

## *Example*

*Noun phrase can be composed of either a ProperNoun or a determiner (Det) followed by a Nominal; a Nominal can be more than one nouns*

# Modeling Constituency

## *Context-free grammar*

The most common way of modeling constituency

## *Consists of production Rules*

These rules express the ways in which the symbols of the language can be grouped and ordered together

## *Example*

*Noun phrase can be composed of either a ProperNoun or a determiner (Det) followed by a Nominal; a Nominal can be more than one nouns*

$NP \rightarrow \text{Det Nominal}$

$NP \rightarrow \text{ProperNoun}$

$\text{Nominal} \rightarrow \text{Noun} \mid \text{Noun Nominal}$

# CFG for Languages

CFG:  $G = (T, N, S, R)$

- $T$ : set of terminals
- $N$ : set of non-terminals
  - ▶ For NLP, we distinguish out a set  $P \subset N$  of pre-terminals, which always rewrite as terminals
- $S$ : start symbol
- $R$ : Rules/productions of the form  $X \rightarrow \gamma$ ,  $X \in N$  and  $\gamma \in (T \cup N)^*$

## CFG: $G = (T, N, S, R)$

- $T$ : set of terminals
- $N$ : set of non-terminals
  - ▶ For NLP, we distinguish out a set  $P \subset N$  of pre-terminals, which always rewrite as terminals
- $S$ : start symbol
- $R$ : Rules/productions of the form  $X \rightarrow \gamma$ ,  $X \in N$  and  $\gamma \in (T \cup N)^*$

## Terminals and pre-terminals

Terminals mainly correspond to words in the language while pre-terminals mainly correspond to POS categories

## CFG: $G = (T, N, S, R)$

- $T$ : set of terminals
- $N$ : set of non-terminals
  - ▶ For NLP, we distinguish out a set  $P \subset N$  of pre-terminals, which always rewrite as terminals
- $S$ : start symbol
- $R$ : Rules/productions of the form  $X \rightarrow \gamma$ ,  $X \in N$  and  $\gamma \in (T \cup N)^*$

## Terminals and pre-terminals

Terminals mainly correspond to words in the language while pre-terminals mainly correspond to POS categories

## Example

$NP \rightarrow \text{Det Nominal}$

$NP \rightarrow \text{ProperNoun}$

$\text{Nominal} \rightarrow \text{Noun} \mid \text{Noun Nominal}$

## Example

$NP \rightarrow \text{Det Nominal}$

$NP \rightarrow \text{ProperNoun}$

$\text{Nominal} \rightarrow \text{Noun} \mid \text{Noun Nominal}$

Now, these can be combined with other rules, that express facts about a lexicon.



## Example

$NP \rightarrow \text{Det Nominal}$

$NP \rightarrow \text{ProperNoun}$

$\text{Nominal} \rightarrow \text{Noun} \mid \text{Noun Nominal}$

Now, these can be combined with other rules, that express facts about a lexicon.

$\text{Det} \rightarrow \text{a}$

$\text{Det} \rightarrow \text{the}$

$\text{Noun} \rightarrow \text{flight}$

## Example

$NP \rightarrow \text{Det Nominal}$

$NP \rightarrow \text{ProperNoun}$

$\text{Nominal} \rightarrow \text{Noun} \mid \text{Noun Nominal}$

Now, these can be combined with other rules, that express facts about a lexicon.

$\text{Det} \rightarrow \text{a}$

$\text{Det} \rightarrow \text{the}$

$\text{Noun} \rightarrow \text{flight}$

Can you identify the terminal, non-terminals and preterminals?

# *CFG as a generator*

*NP* → Det Nominal

*NP* → ProperNoun

*Nominal* → Noun | Noun Nominal

*Det* → a

*Det* → the

*Noun* → flight

# CFG as a generator

$NP \rightarrow \text{Det Nominal}$

$NP \rightarrow \text{ProperNoun}$

$\text{Nominal} \rightarrow \text{Noun} \mid \text{Noun Nominal}$

$\text{Det} \rightarrow \text{a}$

$\text{Det} \rightarrow \text{the}$

$\text{Noun} \rightarrow \text{flight}$

Generating 'a flight':

# *CFG as a generator*

*NP* → Det Nominal

*NP* → ProperNoun

*Nominal* → Noun | Noun Nominal

*Det* → a

*Det* → the

*Noun* → flight

Generating 'a flight':

*NP*

# CFG as a generator

$NP \rightarrow \text{Det Nominal}$

$NP \rightarrow \text{ProperNoun}$

$\text{Nominal} \rightarrow \text{Noun} \mid \text{Noun Nominal}$

$\text{Det} \rightarrow \text{a}$

$\text{Det} \rightarrow \text{the}$

$\text{Noun} \rightarrow \text{flight}$

Generating 'a flight':

$NP \rightarrow \text{Det Nominal}$

# CFG as a generator

$NP \rightarrow \text{Det Nominal}$

$NP \rightarrow \text{ProperNoun}$

$\text{Nominal} \rightarrow \text{Noun} \mid \text{Noun Nominal}$

$\text{Det} \rightarrow \text{a}$

$\text{Det} \rightarrow \text{the}$

$\text{Noun} \rightarrow \text{flight}$

Generating 'a flight':

$NP \rightarrow \text{Det Nominal}$

$\rightarrow \text{Det Noun}$

# CFG as a generator

$NP \rightarrow \text{Det Nominal}$

$NP \rightarrow \text{ProperNoun}$

$\text{Nominal} \rightarrow \text{Noun} \mid \text{Noun Nominal}$

$\text{Det} \rightarrow \text{a}$

$\text{Det} \rightarrow \text{the}$

$\text{Noun} \rightarrow \text{flight}$

Generating 'a flight':

$NP \rightarrow \text{Det Nominal}$

$\rightarrow \text{Det Noun} \rightarrow \text{a Noun}$



# CFG as a generator

$NP \rightarrow \text{Det Nominal}$

$NP \rightarrow \text{ProperNoun}$

$\text{Nominal} \rightarrow \text{Noun} \mid \text{Noun Nominal}$

$\text{Det} \rightarrow \text{a}$

$\text{Det} \rightarrow \text{the}$

$\text{Noun} \rightarrow \text{flight}$

Generating 'a flight':

$NP \rightarrow \text{Det Nominal}$

$\rightarrow \text{Det Noun} \rightarrow \text{a Noun} \rightarrow \text{a flight}$

# CFG as a generator

$NP \rightarrow \text{Det Nominal}$

$NP \rightarrow \text{ProperNoun}$

$\text{Nominal} \rightarrow \text{Noun} \mid \text{Noun Nominal}$

$\text{Det} \rightarrow \text{a}$

$\text{Det} \rightarrow \text{the}$

$\text{Noun} \rightarrow \text{flight}$

Generating 'a flight':

$NP \rightarrow \text{Det Nominal}$

$\rightarrow \text{Det Noun} \rightarrow \text{a Noun} \rightarrow \text{a flight}$

- Thus a CFG can be used to randomly generate a series of strings

# CFG as a generator

$NP \rightarrow \text{Det Nominal}$

$NP \rightarrow \text{ProperNoun}$

$\text{Nominal} \rightarrow \text{Noun} \mid \text{Noun Nominal}$

$\text{Det} \rightarrow \text{a}$

$\text{Det} \rightarrow \text{the}$

$\text{Noun} \rightarrow \text{flight}$

Generating 'a flight':

$NP \rightarrow \text{Det Nominal}$

$\rightarrow \text{Det Noun} \rightarrow \text{a Noun} \rightarrow \text{a flight}$

- Thus a CFG can be used to randomly generate a series of strings
- This sequence of rule expansions is called a derivation of the string of words, usually represented as a tree

A CFG defines a formal language = set of all sentences (string of words) that can be derived by the grammar

A CFG defines a formal language = set of all sentences (string of words) that can be derived by the grammar

- Sentences in this set are said to be **grammatical**
- Sentences outside this set are said to be **ungrammatical**

## *Recursive Definition*

- $PP \rightarrow \text{Prep NP}$
- $NP \rightarrow \text{Noun PP}$

## *Recursive Definition*

- $PP \rightarrow \text{Prep NP}$
- $NP \rightarrow \text{Noun PP}$

## *Example Sentence*

$[_S \text{The mailman ate his } [_{NP} \text{lunch } [_{PP} \text{with his friend } [_{PP} \text{from the cleaning staff } [_{PP} \text{of the building } [_{PP} \text{at the intersection } [_{PP} \text{on the north end } [_{PP} \text{of town}}]]]]]]]$ .

# *What does Context stand for in CFG?*



# What does *Context* stand for in CFG?

- The notion of *context* has nothing to do with the ordinary meaning of word context in language

# What does *Context* stand for in CFG?

- The notion of *context* has nothing to do with the ordinary meaning of word context in language
- All it really means is that the non-terminal on the left-hand side of a rule is out there all by itself (free of context)

# What does Context stand for in CFG?

- The notion of *context* has nothing to do with the ordinary meaning of word context in language
- All it really means is that the non-terminal on the left-hand side of a rule is out there all by itself (free of context)

$A \rightarrow BC$

- I can rewrite  $A$  as  $B$  followed by  $C$  regardless of the context in which  $A$  is found
- Or when I see a  $B$  followed by a  $C$ , I can infer an  $A$  regardless of the surrounding context