

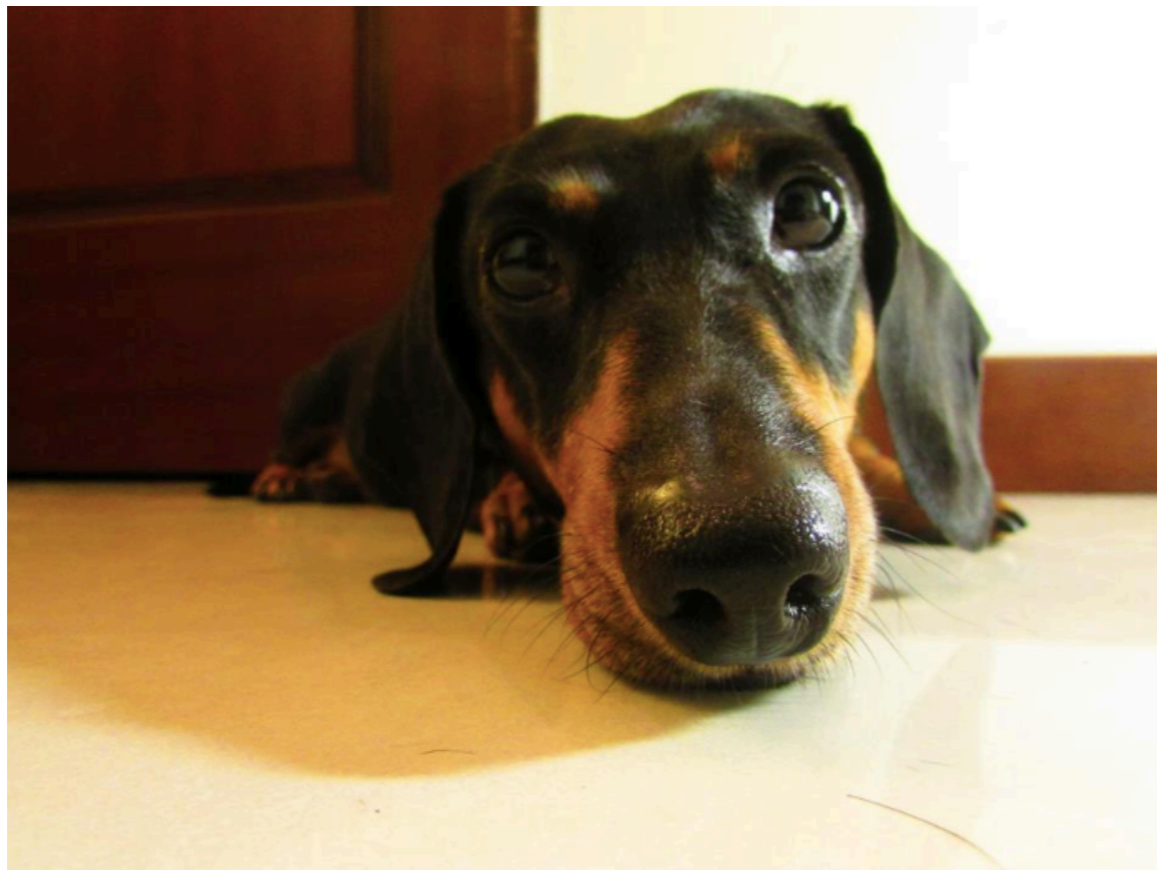
Modeling Grammaticality

Kevin Duh
Fall 2019

Today's Outline

1. What does it mean to be grammatical?
2. Modeling grammaticality with Context-Free Grammars

Is this grammatical?



a cute dog
a very cute dog
super cute puppy
adorable puppy looking at me

....

dog cute a
dog cut a very
puppy cute super
me at puppy looking adorable at

....

Grammar and Syntax

- **Grammar**: Formal rules, principles, or processes that determine valid and invalid structure in language
- **Syntax**: Grammar of sentences
 - (We'll focus on this today)

Prescriptive vs Descriptive

- Prescriptive Grammar
 - How you “ought” to speak. Otherwise, you’re ungrammatical!
 - e.g. Don’t split infinitives! (e.g. “to go”)
- Descriptive Grammar
 - Focus on describing the language as it’s used
 - e.g. “To boldy go where no man has gone before”
- In NLP, we do a bit of both, with probabilities

Grammaticality in perspective

- Dialect differences:
 - I didn't eat dinner
 - I didn't eat no dinner
- Changes in usage:
 - She said, "I want to go!"
 - She was, like, "I want to go!"

“Chinese has no grammar” — false!

- 他喝茶 (Literal: He drinks tea)
 - Grammar rule: Subject, Verb, Object
- 我有一件黑色襯衫 (I have a black shirt)
 - Grammar rule: modifier “black” comes before modified
- All languages have grammar!

Today's Outline

1. What does it mean to be grammatical?
2. Modeling grammaticality with Context-Free Grammars

Goal of modeling grammaticality

- We need a way to mathematically or formally capture what is grammatical and what is not.
- There are many “formalisms” for doing so. We’ll cover:
 - Constituency grammar (today)
 - Dependency grammar (later)

Constituency grammar

(aka phrase structure grammar)

- Focuses on groups of words (**constituent**)
- A **sentence (S)** is made of:
 - subject, typically a **noun phrase (NP)**
 - predicate, typically **verb phrase (VP)**
- NP and VP are in turn made of groups of words

Sentence

Noun Phrase

Verb Phrase

the man walked to the park

Bracketing notation:

((the man) (walked to the park))

Sentence

Noun Phrase

Verb Phrase

the man walked to the park

Bracketing notation:

((the man) (walked (to (the park))))

Sentence

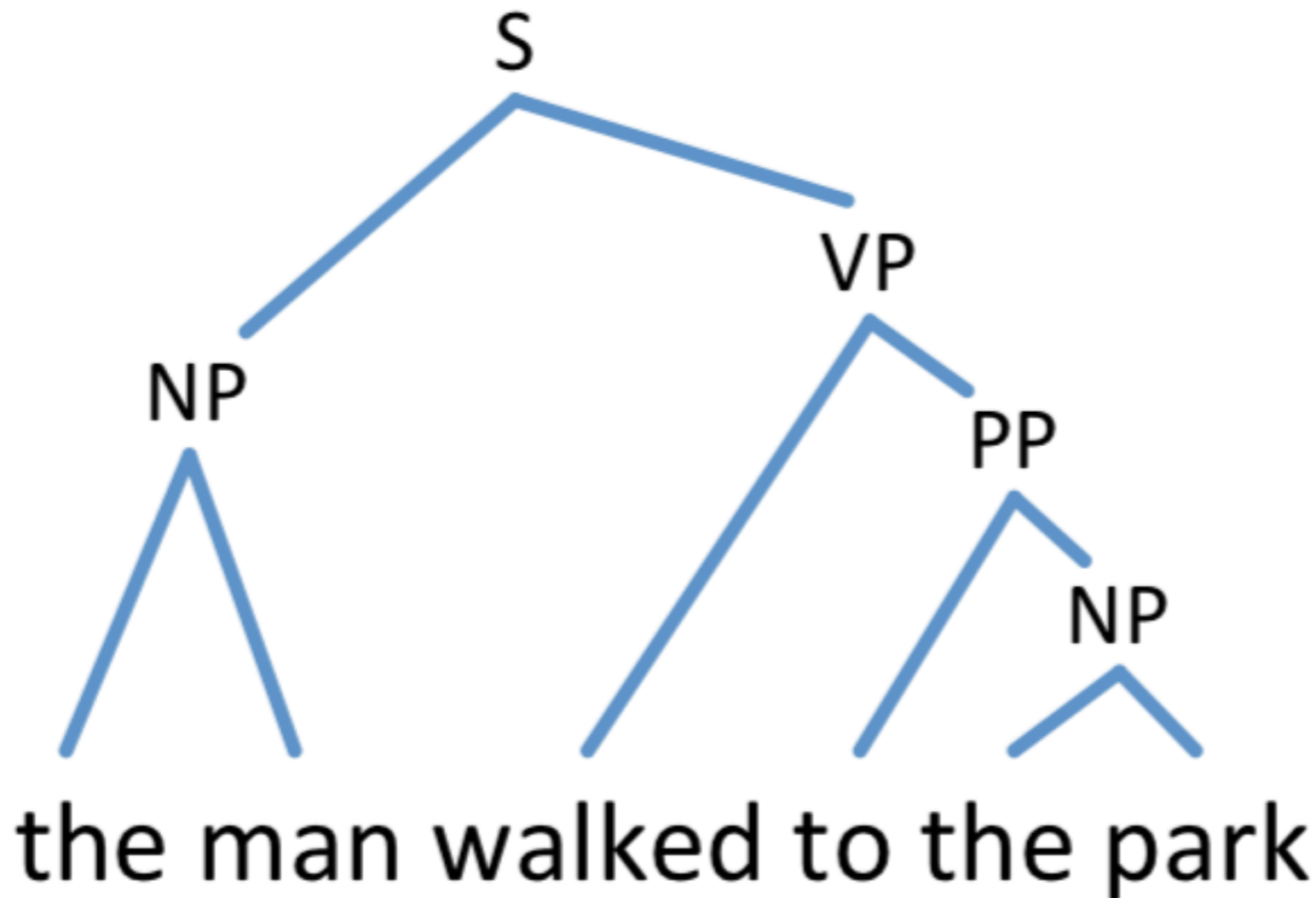
Noun Phrase

Verb Phrase

the man walked to the park

Add labels to each constituent

(S (NP the man) (VP walked (PP to (NP the park))))



Key:

S = sentence

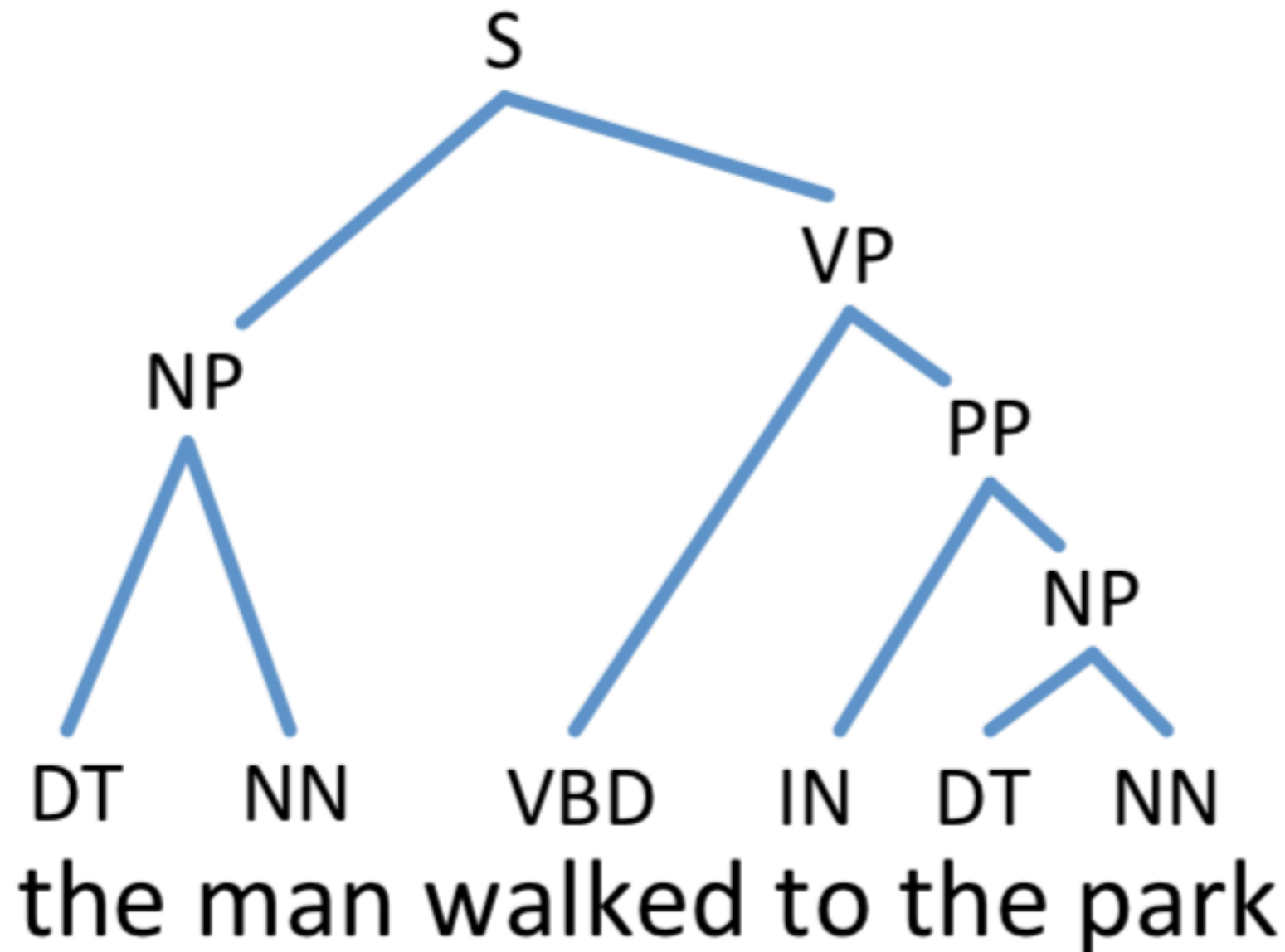
NP = noun phrase

VP = verb phrase

PP = prepositional phrase

Include pre-terminals (part-of-speech labels)

(S (NP the man) (VP walked (PP to (NP the park))))



Key:

S = sentence

NP = noun phrase

VP = verb phrase

PP = prepositional phrase

DT = determiner

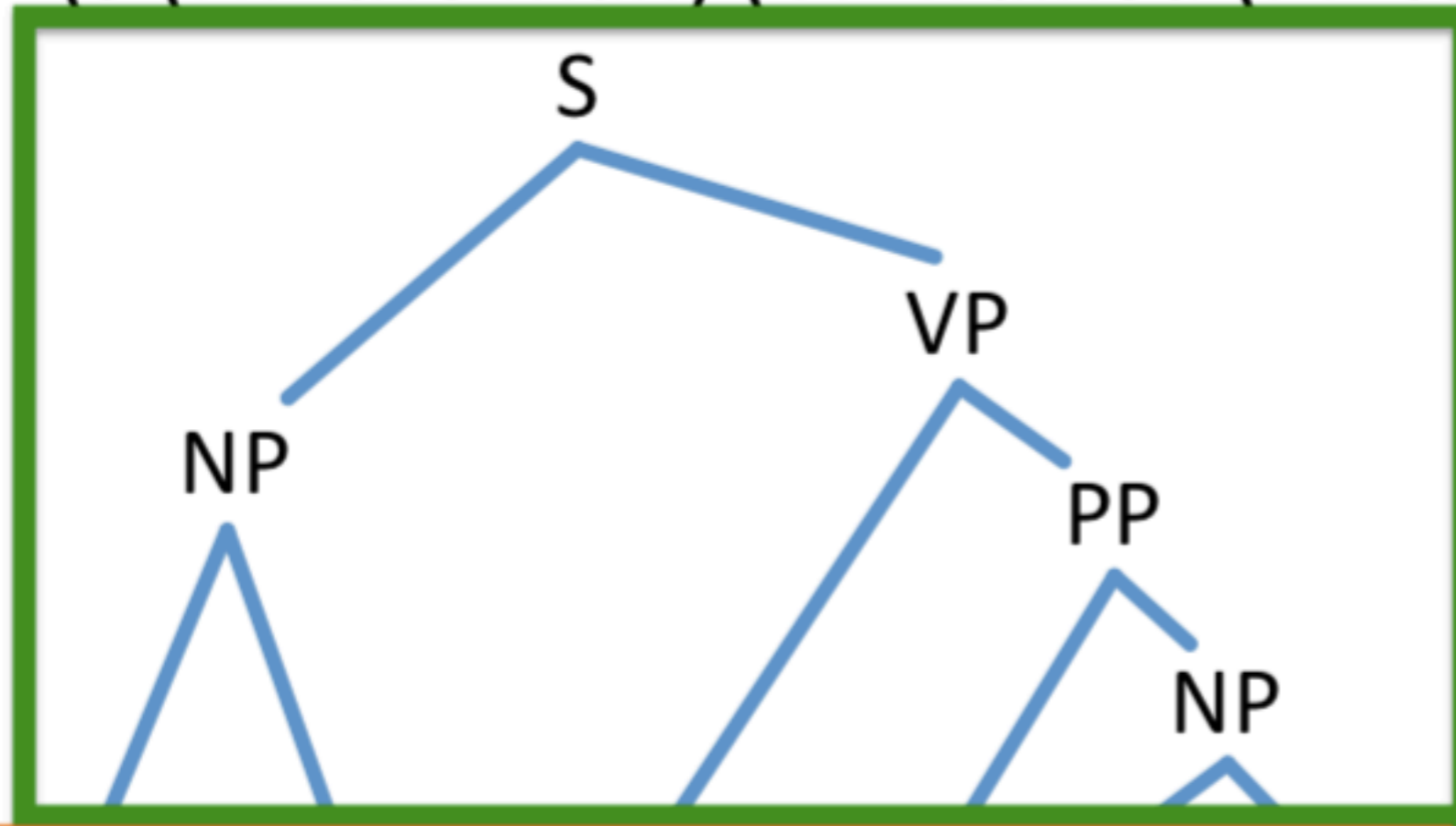
NN = noun

VBD = verb (past tense)

IN = preposition

Now we have a constituency tree!

(S (NP the man) (VP walked (PP to (NP the park))))



nonterminals

DT NN VBD IN DT NN

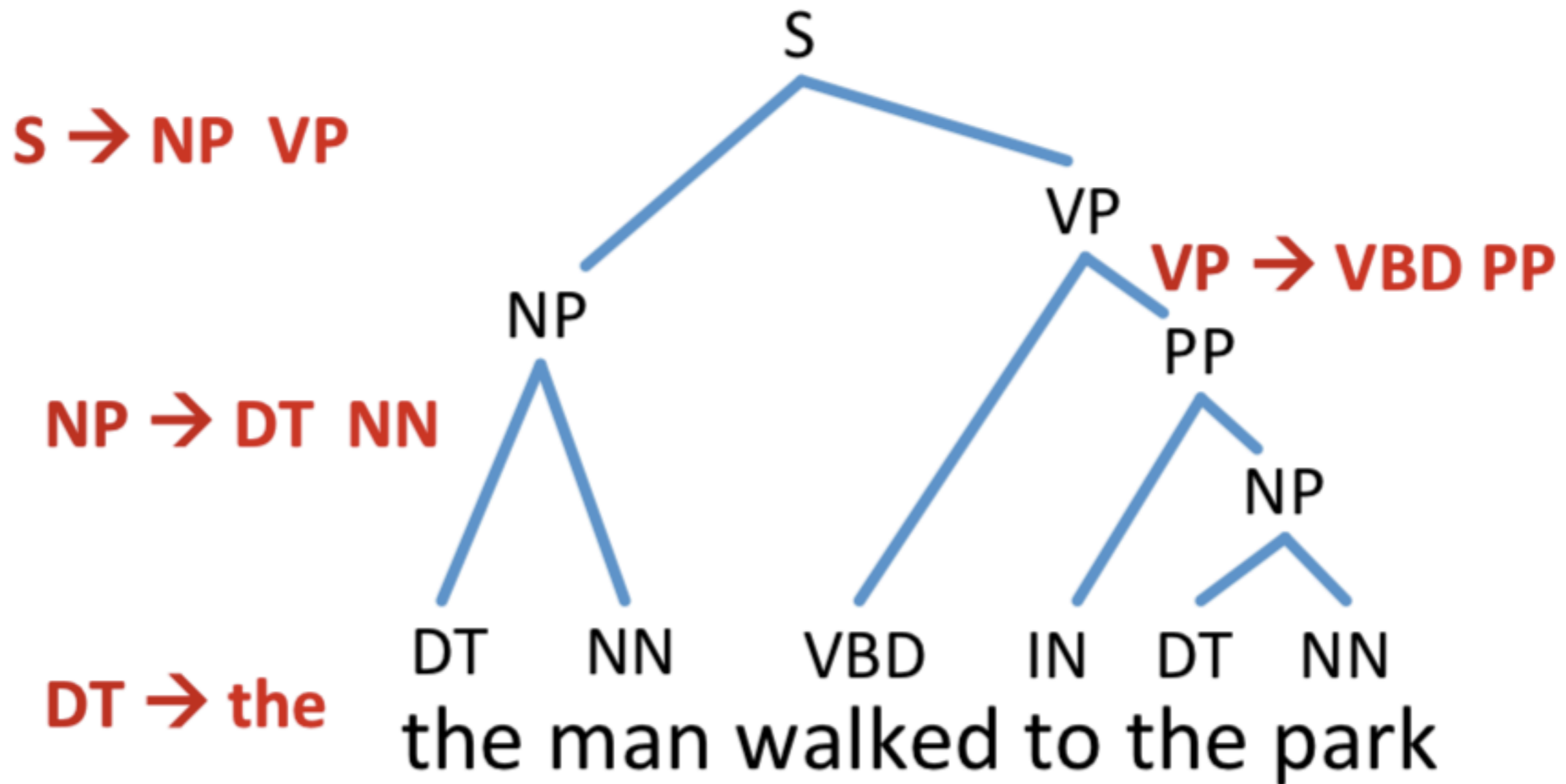
preterminals

the man walked to the park

terminals

Context-Free Grammar

- Syntactic Re-write Rules
 - $S \rightarrow NP VP$
 - $NP \rightarrow DT NN$
 - $VP \rightarrow VBD PP$
 - $PP \rightarrow IN NP$
 - etc
- Lexical Re-write Rules
 - $NN \rightarrow \text{man}$
 - $DT \rightarrow \text{the}$
 - $VBD \rightarrow \text{walked}$
 - $IN \rightarrow \text{to}$
 - $NN \rightarrow \text{park}$



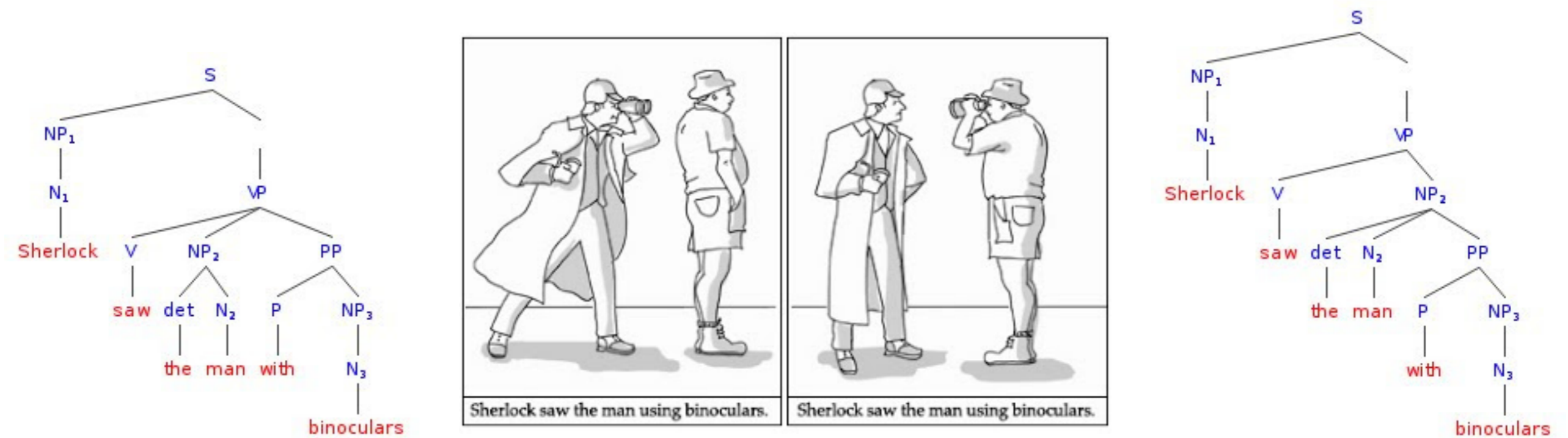
Probabilistic CFG

- Syntactic Re-write Rules
 - $S \rightarrow NP VP$ **Probability=1.0**
 - $NP \rightarrow DT NN$ **Probability=0.7**
 - $VP \rightarrow VBD PP$ **Probability=1.0**
 - $PP \rightarrow IN NP$ **Probability=1.0**
 - $NP \rightarrow NNP$ **Probability=0.3**
 - etc
- Lexical Re-write Rules
 - $NN \rightarrow \text{man}$ **Probability=0.4**
 - $DT \rightarrow \text{the}$ **Probability=1.0**
 - $VBD \rightarrow \text{walked}$ **Probability=1.0**
 - $IN \rightarrow \text{to}$ **Probability=1.0**
 - $NN \rightarrow \text{park}$ **Probability=0.4**
 - $NN \rightarrow \text{John}$ **Probability=0.2**

Top-down generation

Ambiguities - Prepositional Phrase (PP) Attachment

Sherlock saw the man using binoculars



Ambiguities - more examples

- Coordination:
 - ((laptop and monitor) with the Apple logo)
 - (laptop and (monitor with the Apple logo))
- Noun compound
 - ((Natural Language) Processing)
 - (Natural (Language Processing))

CFG Formalism

- **$G=(\Sigma,N,S,R)$**
 - Σ is finite set of terminal, e.g. a, b
 - N is finite set of nonterminal, e.g. A, B ($V = \Sigma \cup N$)
 - S is start symbol
 - R is production rule $A \rightarrow a$ where a is V^*
 - For PCFG, probability is attached to each R
- **Chomsky Normal Form (CNF) – only these rules are allowed**
 - unary terminal rule $A \rightarrow w$
 - binary nonterminal rule $A \rightarrow B C$

Why is it called context-free?

- A rule like $NP \rightarrow DT\ NN$ applies regardless of the neighboring context of NP
- i.e. left-hand-side of each rule is a single non-terminal symbol

Today's Outline

1. What does it mean to be grammatical?
2. Modeling grammaticality with Context-Free Grammars