



CS 4650/7650: Natural Language Processing

Constituency Parsing

Diyi Yang

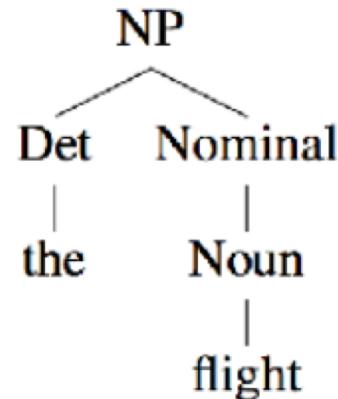
Many slides from Yulia Tsvetkov, Greg Durrett, David Bamman, and others

Logistics

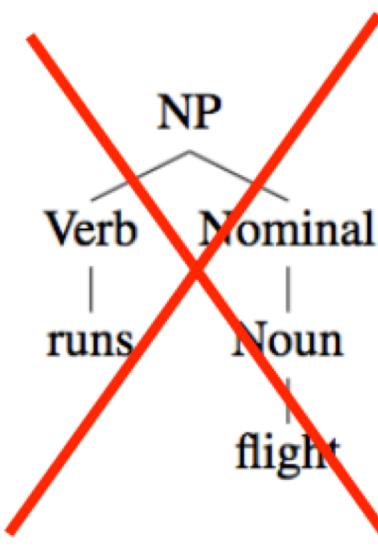
- Midterm review: Feb 24th
- Midterm: Feb 26th, 3:00-4:15pm
 - One-page cheat sheet

Context-free Grammar

- A CFG gives a formal way to define what meaningful constituents are and exactly how a constituent is formed out of other constituents (or words). It defines **valid structure** in a language.



$NP \rightarrow Det\ Nominal$



$NP \rightarrow Verb\ Nominal$

Context-free Grammar

- A context-free grammar defines how symbols in a language combine to form valid structures

NP	→	Det Nominal
NP	→	ProperNoun
Nominal	→	Noun Nominal Noun
Det	→	a the
Noun	→	flight

non-terminals

lexicon/
terminals

Context-free Grammar

N Finite set of non-terminal symbols NP, VP, S

Σ Finite alphabet of terminal symbols the, dog, a

R Set of production rules, each
 $A \rightarrow \beta$
 $\beta \in (\Sigma, N)$ $S \rightarrow NP\ VP$
Noun \rightarrow dog

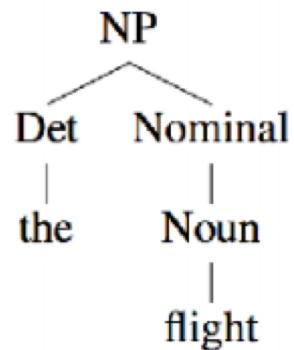
S Start symbol

Infinite Strings with Finite Productions

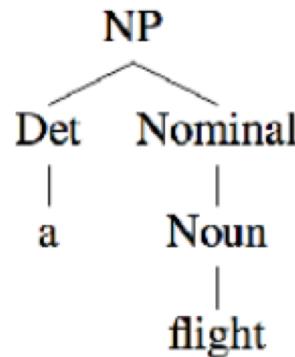
- This is the house
- This is the house that Jack built
- This is the cat that lives in the house that Jack built
- This is the dog that chased the cat that lives in the house that Jack built
- This is the flea that bit the dog that chased the cat that lives in the house the Jack built
- This is the virus that infected the flea that bit the dog that chased the cat that lives in the house that Jack built

Derivation

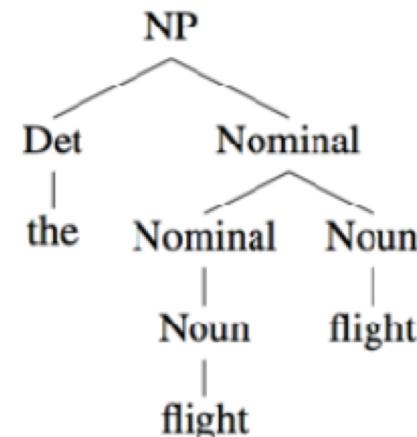
Given a CFG, a derivation is the sequence of productions used to generate a string of words (e.g., a sentence), often visualized as a parse tree



the flight

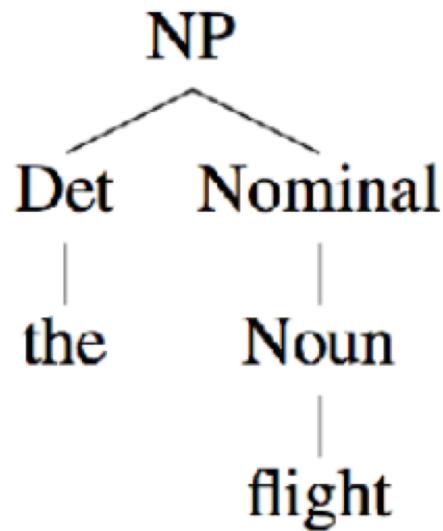


a flight



the flight flight

Bracketed Notation



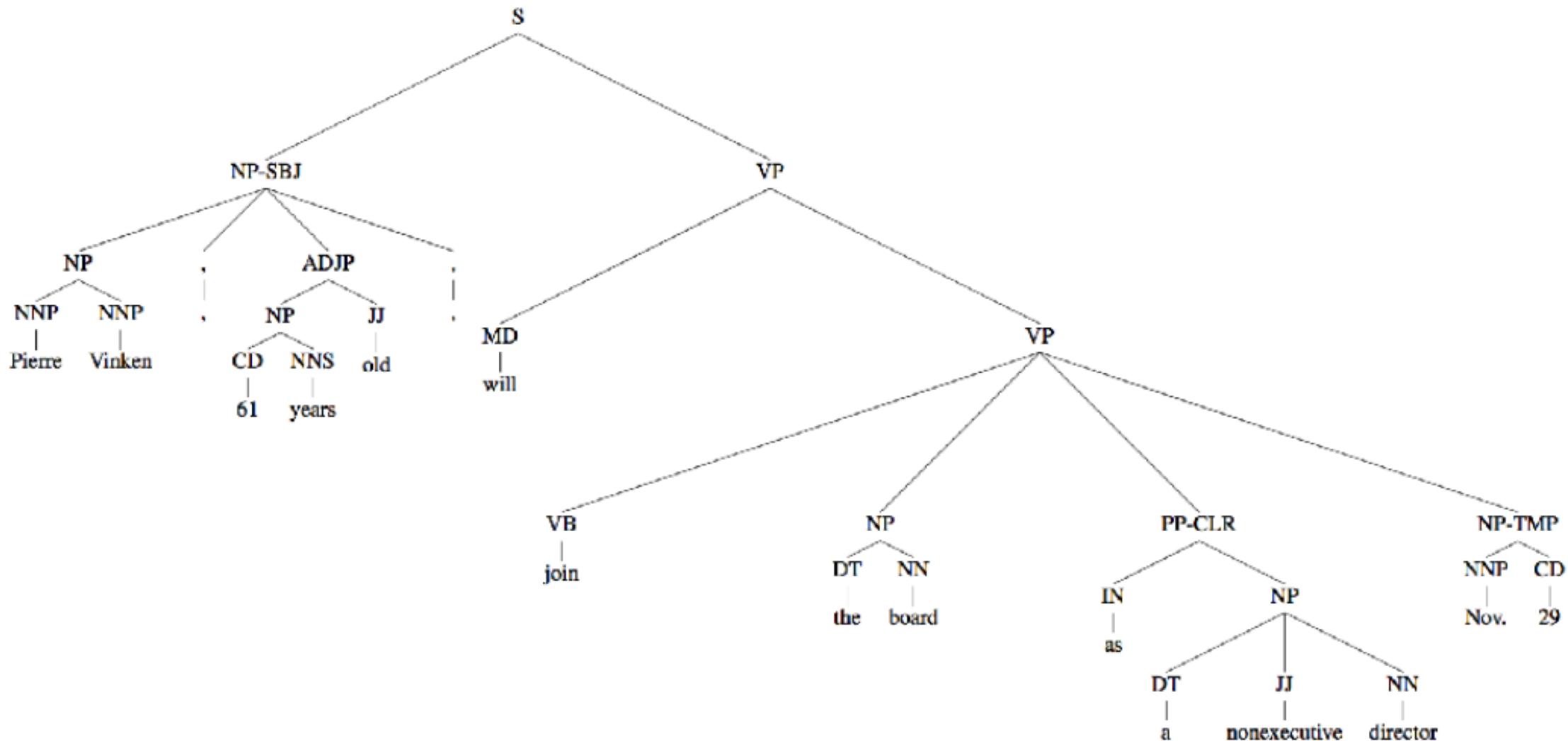
[NP [Det **the**] [Nominal [Noun **flight**]]]

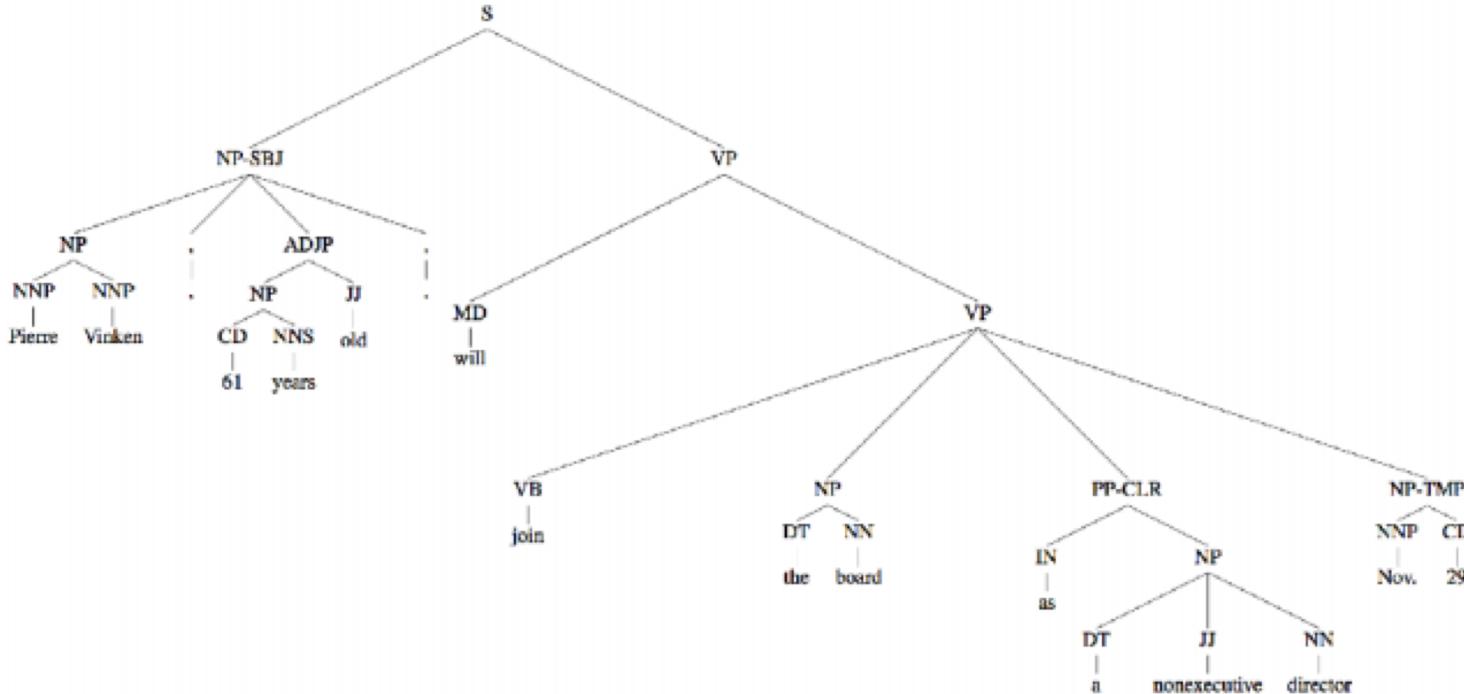
Treebanks

- Rather than create the rules by hand, we can annotate sentences with their syntactic structure and then extract the rules from the annotations
- Treebanks:

Collections of sentences annotated with **syntactic structure**

Penn Treebank





NP	\rightarrow	NNP NNP
NP-SBJ	\rightarrow	NP , ADJP ,
S	\rightarrow	NP-SBJ VP
VP	\rightarrow	VB NP PP-CLR NP-TMP

Summary: Syntax Through The Use Of Context-free Grammar

- Groups of consecutive words act as a group or a constituent
- A context-free grammar consists of a set of rules or productions, expressed over a set of non-terminal symbols and a set of terminal symbols.
- A particular context-free language is the set of strings that can be derived from a particular context-free grammar
- Verbs can be subcategorized by the types of complements they expect. Simple subcategories are transitive and intransitive
- Treebanks of parsed sentences exist for many genres of English and for many languages.

Summary: Syntax Through The Use Of Context-free Grammar

- Groups of consecutive words act as a group or a **constituent**
- A **context-free grammar** consists of a set of **rules** or **productions**, expressed over a set of **non-terminal** symbols and a set of **terminal** symbols.
- A particular **context-free language** is the set of strings that can be derived from a particular **context-free grammar**
- Verbs can be **subcategorized** by the types of **complements** they expect. Simple subcategories are **transitive** and **intransitive**
- **Treebanks** of parsed sentences exist for many genres of English and for many languages.

CFG

- A basic CFG allows us to check whether a sentence is grammatical in the language it defines
- **Binary decision:** a sentence is either in the language (a series of productions yields the words we see) or it is not
- Where would this be useful?

PCFG

- Probabilistic context-free grammar: each production is also associated with a probability
- This lets us calculate the probability of a parse for a given sentence; for a given parse tree T for sentence S comprised of n rules from R (each $A \rightarrow \beta$):

$$P(T, S) = \prod_{i=1}^n P(\beta | A)$$

PCFG

N Finite set of non-terminal symbols NP, VP, S

Σ Finite alphabet of terminal symbols the, dog, a

R Set of production rules, each
 $A \rightarrow \beta$ [p]
 $p = P(\beta | A)$ $S \rightarrow NP\ VP$
 Noun \rightarrow dog

S Start symbol

Estimating PCFGs

How do we calculate $P(A \rightarrow \beta)$?

Maximum likelihood estimates

Estimating PCFGs

$$\sum_{\beta} P(\beta \mid A) = \frac{C(A \rightarrow \beta)}{\sum_{\gamma} C(A \rightarrow \gamma)}$$

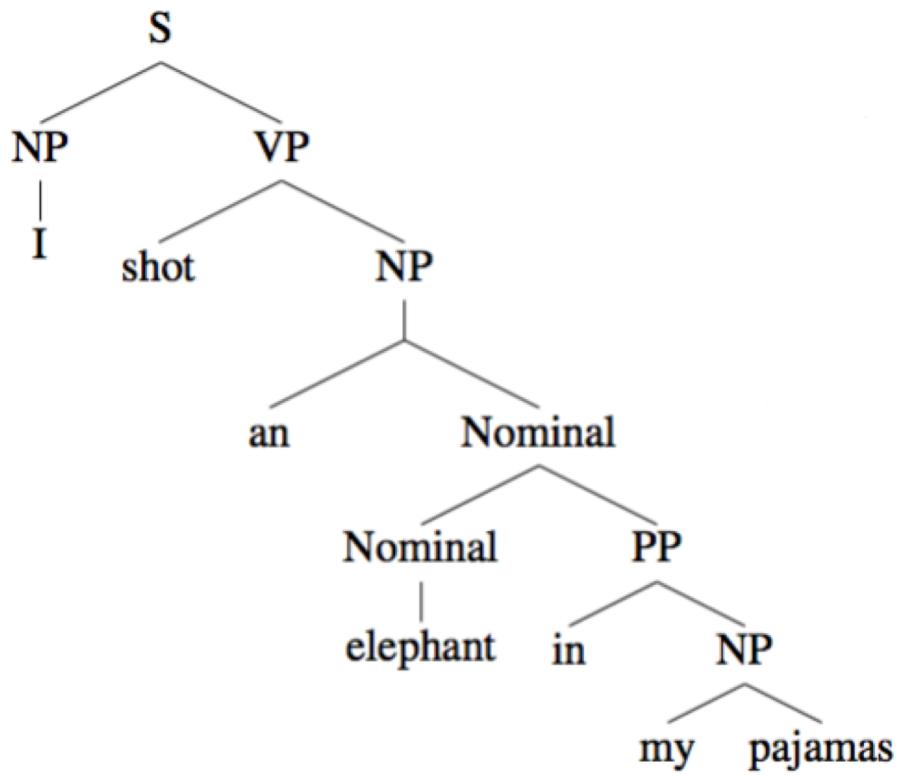
(equivalently)

$$\sum_{\beta} P(\beta \mid A) = \frac{C(A \rightarrow \beta)}{C(A)}$$

A		β	$P(\beta NP)$
NP	\rightarrow	NP PP	0.092
NP	\rightarrow	DT NN	0.087
NP	\rightarrow	NN	0.047
NP	\rightarrow	NNS	0.042
NP	\rightarrow	DT JJ NN	0.035
NP	\rightarrow	NNP	0.034
NP	\rightarrow	NNP NNP	0.029
NP	\rightarrow	JJ NNS	0.027
NP	\rightarrow	QP -NONE-	0.018
NP	\rightarrow	NP SBAR	0.017
NP	\rightarrow	NP PP-LOC	0.017
NP	\rightarrow	JJ NN	0.015
NP	\rightarrow	DT NNS	0.014
NP	\rightarrow	CD	0.014
NP	\rightarrow	NN NNS	0.013
NP	\rightarrow	DT NN NN	0.013
NP	\rightarrow	NP CC NP	0.013

- CNF and CKY Parsing

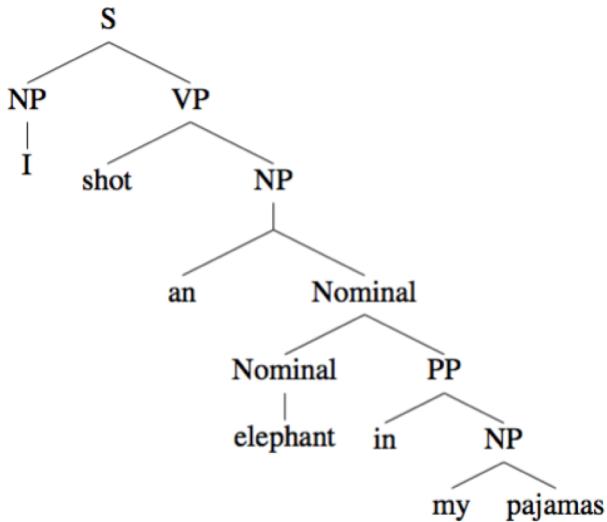
Constituents



Every internal node is a phrase

- My pajamas
- In my pajamas
- Elephant in my pajamas
- an elephant in my pajamas
- Shot an elephant in my pajamas
- I shot an elephant in my pajamas

Constituents



Each phrase could be replaced by another of the same type of constituent

Every internal node is a phrase

- My pajamas
- In my pajamas
- Elephant in my pajamas
- an elephant in my pajamas
- Shot an elephant in my pajamas
- I shot an elephant in my pajamas

Context-Free Grammar

N	Finite set of non-terminal symbols	NP, VP, S
Σ	Finite alphabet of terminal symbols	the, dog, a
R	Set of production rules, each $A \rightarrow \beta$ $\beta \in (\Sigma, N)$	$NP \rightarrow DT\ JJ\ NN$ Noun \rightarrow dog
S	Start symbol	

Chomsky Normal Form (CNF)

N Finite set of non-terminal symbols NP, VP, S

Σ Finite alphabet of terminal symbols the, dog, a

R Set of production rules, each $A \rightarrow \beta$
 β = single terminal (from Σ) or two
non-terminals (from N) S → NP VP
 Noun → dog

S Start symbol

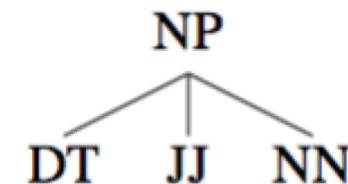
Context-free Grammar vs. Chomsky Normal Form

- A context-free grammar is in Chomsky normal form (CNF) if
 - each production is either $A \rightarrow BC$, $A \rightarrow a$, or $S \rightarrow \epsilon$

Chomsky Normal Form (CNF)

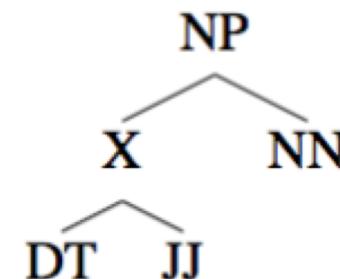
- Any CFG can be converted into weakly equivalent CNF (recognizing the same set of sentences as existing in the grammar but differing in their derivation)

$NP \rightarrow DT\ JJ\ NN$



$NP \rightarrow X\ NN$

$X \rightarrow DT\ JJ$



S	→	NP VP
VP	→	VBD NP
VP	→	VP PP
Nominal	→	Nominal PP
Nominal	→	NN
Nominal	→	NNS
Nominal	→	PRP
PP	→	IN NP
NP	→	DT NN
NP	→	Nominal
NP	→	PRP\$ Nominal

VBD	→	shot
DT	→	an my
NN	→	elephant
NNS	→	pajamas
PRP	→	I
PRP\$	→	my
IN	→	in

I shot an elephant in my pajamas

S → NP VP
VP → VBD NP
VP → VP PP
Nominal → Nominal PP
Nominal → pajamas elephant I
PP → IN NP
NP → DT NN
NP → pajamas elephant I
NP → PRP\$ Nominal

VBD → shot
DT → an my
PRP → I
PRP\$ → my
IN → in

I shot an elephant in my pajamas

CKY

- Cocke-Kasami-Younger algorithm (also CYK) for parsing for a grammar expressed in CNF.
- Bottom-up dynamic programming

CKY

- The subproblems for our dynamic program are all possible spans of contiguous words in the sentence.

S → NP VP
VP → VBD NP
VP → VP PP
Nominal → Nominal PP
Nominal → pajamas elephant I
PP → IN NP
NP → DT NN
NP → pajamas elephant I
NP → PRP\$ Nominal

VBD → shot
DT → an my
PRP → I
PRP\$ → my
IN → in

I shot an elephant in my pajamas

0 I 1 shot 2 an 3 elephant 4 in 5 my 6 pajamas 7

I	shot	an	elephant	in	my	pajamas
---	------	----	----------	----	----	---------

NP, PRP [0,1]						
	VBD [1,2]					
		DT [2,3]				
			NP, NN [3,4]			
				IN [4,5]		
					PRP\$ [5,6]	
						NNS [6,7]

Each cell i,j keeps track of all phrase types that can be formed from *all* words from position i through position j

I	shot	an	elephant	in	my	pajamas
---	------	----	----------	----	----	---------

NP, PRP [0,1]						
	VBD [1,2]					
		DT [2,3]				
			NP, NN [3,4]			
				IN [4,5]		
					PRP\$ [5,6]	
						NNS [6,7]

What phrases can be formed from “shot an elephant in”

I	shot	an	elephant	in	my	pajamas
---	------	----	----------	----	----	---------

NP, PRP [0,1]						
	VBD [1,2]					
		DT [2,3]				
			NP, NN [3,4]			
				IN [4,5]		
					PRP\$ [5,6]	
						NNS [6,7]
What phrases can be formed from “I shot an elephant in my pajamas”						

CKY

- In CNF, each non-terminal generates two non-terminals

$$S \rightarrow NP\ VP$$

[S [NP I] [VP shot an elephant in my pajamas]]

- The recursive step of our dynamic program needs to select a split point and rule
- If the parent spans tokens i-j, then there is a k where the children span i-k and k-j.

Completed Parse Table

<i>Book</i>	<i>the</i>	<i>flight</i>	<i>through</i>	<i>Houston</i>
S, VP, Verb Nominal, Noun [0,1]	[0,2]	S,VP,X2 [0,3]	[0,4]	S,VP,X2 [0,5]
Det [1,2]	NP [1,3]	[1,4]	NP [1,5]	
Nominal, Noun [2,3]	[2,4]	[2,5]	Nominal [2,6]	
Prep [3,4]	PP [3,5]			
NP, Proper- Noun [4,5]				

Figure 13.4 Completed parse table for *Book the flight through Houston.*

CKY Algorithm

```
function CKY-PARSE(words, grammar) returns table
    for  $j \leftarrow$  from 1 to LENGTH(words) do
        for all  $\{A \mid A \rightarrow words[j] \in grammar\}$ 
            table[ $j - 1, j$ ]  $\leftarrow$  table[ $j - 1, j$ ]  $\cup$  A
        for  $i \leftarrow$  from  $j - 2$  downto 0 do
            for  $k \leftarrow i + 1$  to  $j - 1$  do
                for all  $\{A \mid A \rightarrow BC \in grammar \text{ and } B \in table[i, k] \text{ and } C \in table[k, j]\}$ 
                    table[ $i, j$ ]  $\leftarrow$  table[ $i, j$ ]  $\cup$  A
```

Figure 13.5 The CKY algorithm.

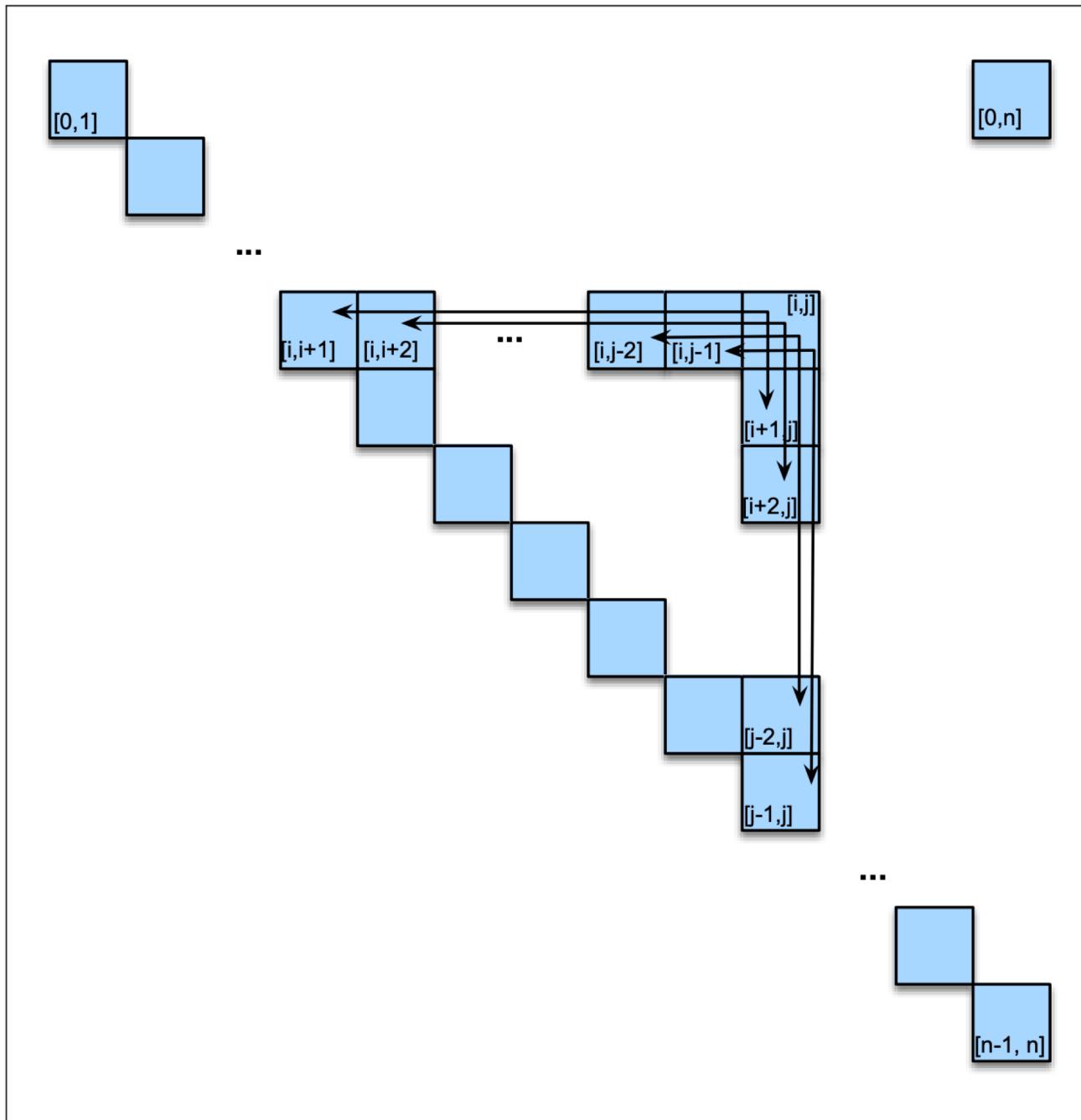


Figure 13.6 All the ways to fill the $[i, j]$ th cell in the CKY table.

I	shot	an	elephant	in	my	pajamas
---	------	----	----------	----	----	---------

NP, PRP [0,1]						
	VBD [1,2]					
		DT [2,3]				
			NP, NN [3,4]			
				IN [4,5]		
					PRP\$ [5,6]	
						NNS [6,7]

Does any rule generate PRP
VBD?

I	shot	an	elephant	in	my	pajamas
---	------	----	----------	----	----	---------

NP, PRP [0,1]						
	VBD [1,2]					
		DT [2,3]				
			NP, NN [3,4]			
				IN [4,5]		
					PRP\$ [5,6]	
						NNS [6,7]

Does any rule generate PRP VBD?

S → NP VP
VP → VBD NP
VP → VP PP
Nominal → Nominal PP
Nominal → pajamas | elephant | I
PP → IN NP
NP → DT NN
NP → pajamas | elephant | I
NP → PRP\$ Nominal

VBD → shot
DT → an | my
PRP → I
PRP\$ → my
IN → in

I shot an elephant in my pajamas

I	shot	an	elephant	in	my	pajamas
---	------	----	----------	----	----	---------

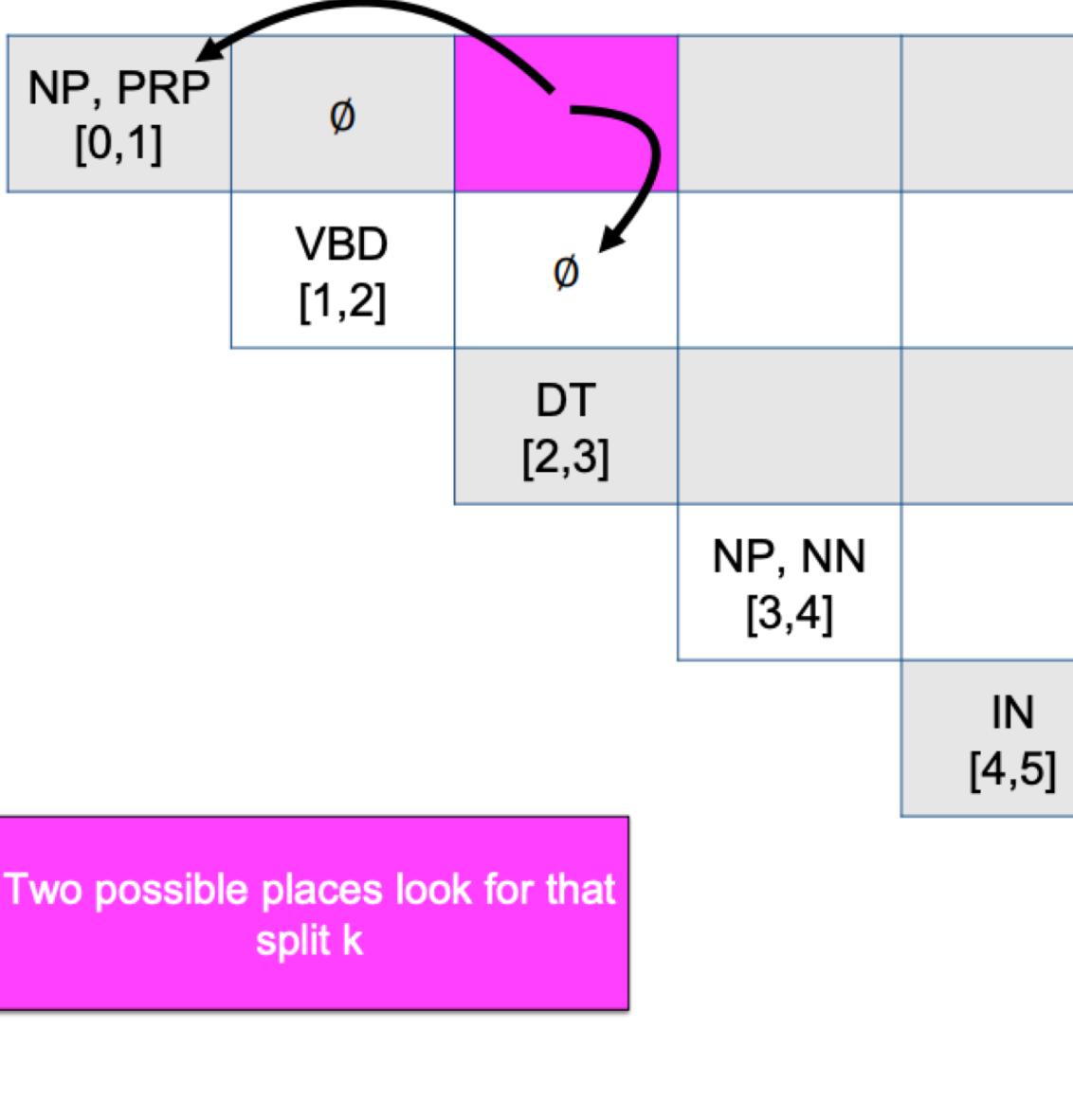
NP, PRP [0,1]	\emptyset					
	VBD [1,2]					
		DT [2,3]				
			NP, NN [3,4]			

S → NP VP		
VP → VBD NP		
VP → VP PP		
Nominal → Nominal PP		
Nominal → pajamas elephant I		
PP → IN NP		
NP → DT NN		
NP → pajamas elephant I		
NP → PRP\$ Nominal		

I shot an elephant in my pajamas

Does any rule generate
VBD DT?

I	shot	an	elephant	in	my	pajamas
---	------	----	----------	----	----	---------



Two possible places look for that split k

S → NP VP
VP → VBD NP
VP → VP PP
Nominal → Nominal PP
Nominal → pajamas elephant I
PP → IN NP
NP → DT NN
NP → pajamas elephant I
NP → PRP\$ Nominal

VBD → shot
DT → an my
PRP → I
PRP\$ → my
IN → in

I shot an elephant in my pajamas

I	shot	an	elephant	in	my	pajamas
---	------	----	----------	----	----	---------

NP, PRP [0,1]	\emptyset					
VBD [1,2]		\emptyset				
		DT [2,3]				
			NP, NN [3,4]			

S → NP VP			
VP → VBD NP			
VP → VP PP			
Nominal → Nominal PP			
Nominal → pajamas elephant I			
PP → IN NP			
NP → DT NN			
NP → pajamas elephant I			
NP → PRP\$ Nominal			

I shot an elephant in my pajamas

Two possible places look for that split k

IN [4,5]		
PRP\$ [5,6]		
NNS [6,7]		

I	shot	an	elephant	in	my	pajamas
---	------	----	----------	----	----	---------

NP, PRP [0,1]	\emptyset	\emptyset				
VBD [1,2]		\emptyset				
	DT [2,3]					

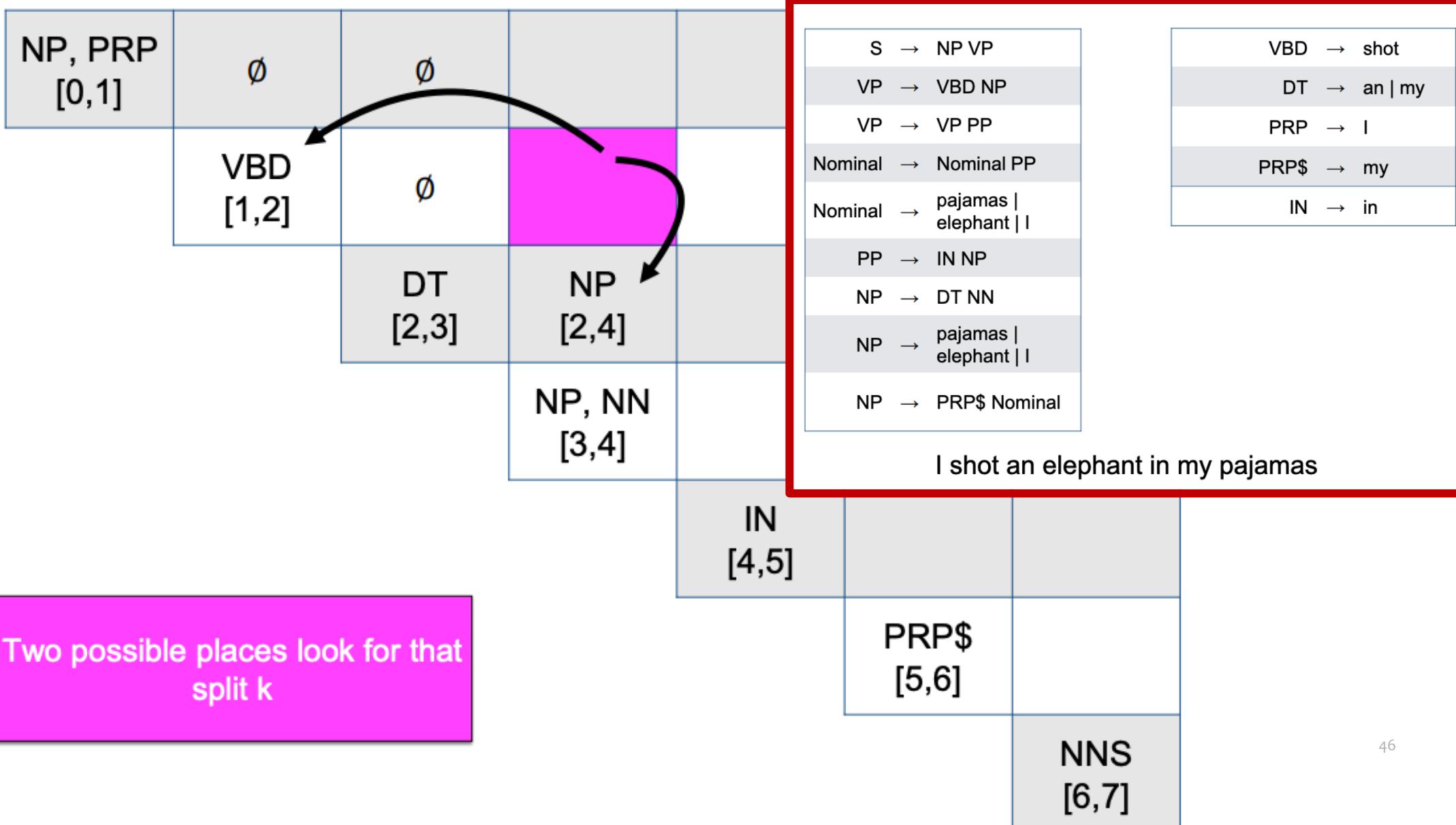
S → NP VP		
VP → VBD NP		
VP → VP PP		
Nominal → Nominal PP		
Nominal → pajamas elephant I		
PP → IN NP		
NP → DT NN		
NP → pajamas elephant I		
NP → PRP\$ Nominal		

I shot an elephant in my pajamas

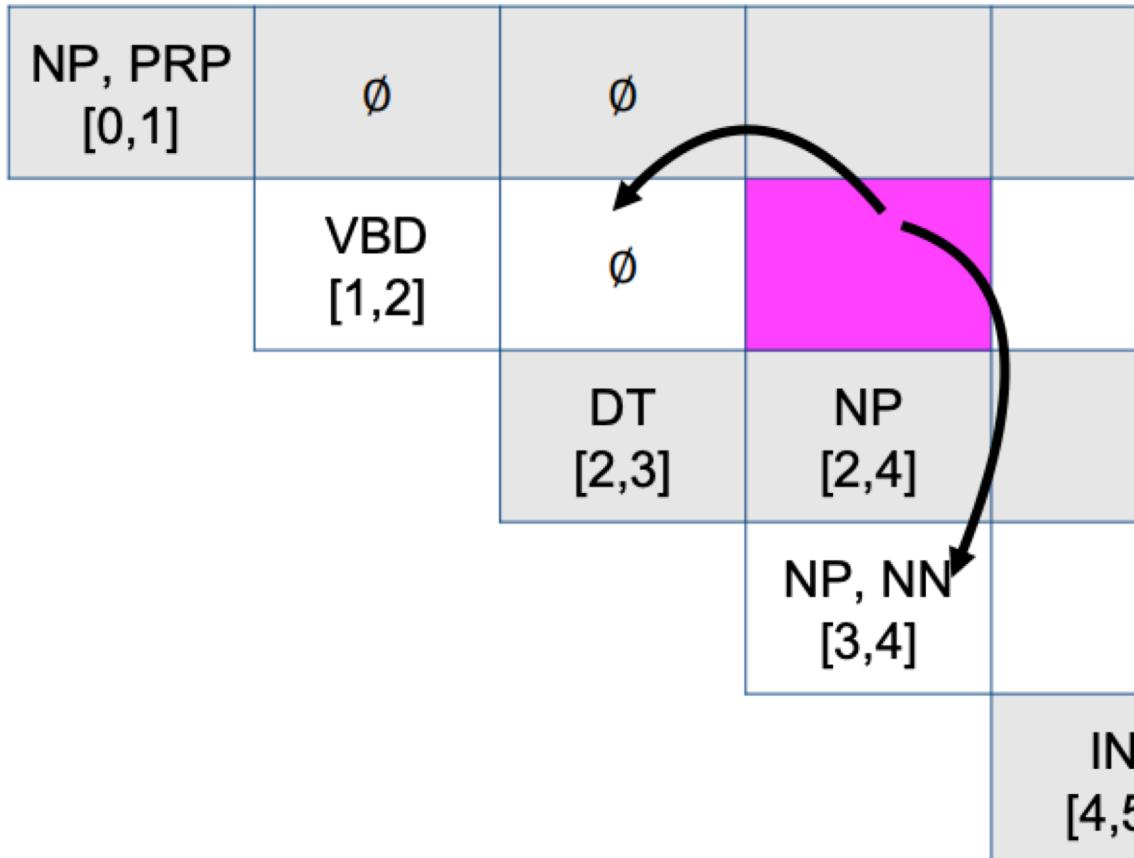
Does any rule generate
DT NN?

IN [4,5]		
PRP\$ [5,6]		
NNS [6,7]		

I	shot	an	elephant	in	my	pajamas
---	------	----	----------	----	----	---------



I	shot	an	elephant	in	my	pajamas
---	------	----	----------	----	----	---------



Two possible places look for that split k

S → NP VP		
VP → VBD NP		
VP → VP PP		
Nominal → Nominal PP		
Nominal → pajamas elephant I		
PP → IN NP		
NP → DT NN		
NP → pajamas elephant I		
NP → PRP\$ Nominal		

I shot an elephant in my pajamas

I	shot	an	elephant	in	my	pajamas
---	------	----	----------	----	----	---------

NP, PRP [0,1]	Ø	Ø				VBD → shot DT → an my PRP → I PRP\$ → my IN → in
	VBD [1,2]	Ø	VP [1,4]			
		DT [2,3]	NP [2,4]			
			NP, NN [3,4]			
				IN [4,5]		
					PRP\$ [5,6]	
						NNS [6,7]

Three possible places look for that split k

S → NP VP
 VP → VBD NP
 VP → VP PP
 Nominal → Nominal PP
 Nominal → pajamas | elephant | I
 PP → IN NP
 NP → DT NN
 NP → pajamas | elephant | I
 NP → PRP\$ Nominal

I shot an elephant in my pajamas

I	shot	an	elephant	in	my	pajamas
---	------	----	----------	----	----	---------

NP, PRP [0,1]	\emptyset	\emptyset				
VBD [1,2]	\emptyset		VP [1,4]			
	DT [2,3]		NP [2,4]			
			NP, NN [3,4]			

S → NP VP
VP → VBD NP
VP → VP PP
Nominal → Nominal PP
Nominal → pajamas elephant I
PP → IN NP
NP → DT NN
NP → pajamas elephant I
NP → PRP\$ Nominal

VBD → shot
DT → an my
PRP → I
PRP\$ → my
IN → in

I shot an elephant in my pajamas

IN [4,5]		
PRP\$ [5,6]		
NNS [6,7]		

Three possible places look for that split k

I	shot	an	elephant	in	my	pajamas
---	------	----	----------	----	----	---------

NP, PRP [0,1]	\emptyset	\emptyset				
VBD [1,2]	\emptyset		VP [1,4]			
	DT [2,3]		NP [2,4]			
		NP, NN [3,4]				

S → NP VP
VP → VBD NP
VP → VP PP
Nominal → Nominal PP
Nominal → pajamas elephant I
PP → IN NP
NP → DT NN
NP → pajamas elephant I
NP → PRP\$ Nominal

VBD → shot
DT → an my
PRP → I
PRP\$ → my
IN → in

I shot an elephant in my pajamas

Three possible places look for that split k

IN [4,5]		
PRP\$ [5,6]		
NNS [6,7]		

I	shot	an	elephant	in	my	pajamas
---	------	----	----------	----	----	---------

NP, PRP [0,1]	\emptyset	\emptyset				
VBD [1,2]	\emptyset	\emptyset	VP [1,4]			
	DT [2,3]		NP [2,4]			
			NP, NN [3,4]			
				IN [4,5]		

S → NP VP			
VP → VBD NP			
VP → VP PP			
Nominal → Nominal PP			
Nominal → pajamas elephant I			
PP → IN NP			
NP → DT NN			
NP → pajamas elephant I			
NP → PRP\$ Nominal			

I shot an elephant in my pajamas

Three possible places look for that split k

	PRP\$ [5,6]	
		NN [6,7]

I	shot	an	elephant	in	my	pajamas
---	------	----	----------	----	----	---------

NP, PRP [0,1]	\emptyset	\emptyset	S [0,4]			
	VBD [1,2]	\emptyset	VP [1,4]			
		DT [2,3]	NP [2,4]			
			NP, NN [3,4]			
				IN [4,5]		
					PRP\$ [5,6]	
						NNS [6,7]
<p>S → NP VP VP → VBD NP VP → VP PP Nominal → Nominal PP Nominal → pajamas elephant I PP → IN NP NP → DT NN NP → pajamas elephant I NP → PRP\$ Nominal</p> <p>VBD → shot DT → an my PRP → I PRP\$ → my IN → in</p> <p>I shot an elephant in my pajamas</p>						

I	shot	an	elephant	in	my	pajamas
---	------	----	----------	----	----	---------

NP, PRP [0,1]	Ø	Ø	S [0,4]	Ø	Ø	
	VBD [1,2]	Ø	VP [1,4]	Ø	Ø	
		DT [2,3]	NP [2,4]	Ø	Ø	
			NP, NN [3,4]	Ø	Ø	
				IN [4,5]	Ø	
*elephant in	*in my				PRP\$ [5,6]	
*an elephant in	*elephant in my					
*shot an elephant in	*an elephant in my					
*I shot an elephant in	*shot an elephant in my					
	*I shot an elephant in my					NNS [6,7]

I	shot	an	elephant	in	my	pajamas
---	------	----	----------	----	----	---------

NP, PRP [0,1]	\emptyset	\emptyset	S [0,4]	\emptyset	\emptyset	
	VBD [1,2]	\emptyset	VP [1,4]	\emptyset	\emptyset	
		DT [2,3]	NP [2,4]	\emptyset	\emptyset	
			NP, NN [3,4]	\emptyset	\emptyset	
				IN [4,5]	\emptyset	
					PRP\$ [5,6]	NP [5,7]
						NNS [6,7]

Grammar Rules:

- S → NP VP
- VP → VBD NP
- VP → VP PP
- Nominal → Nominal PP
- Nominal → pajamas | elephant | I
- PP → IN NP
- NP → DT NN
- NP → pajamas | elephant | I
- NP → PRP\$ Nominal

Lexical Entries:

- VBD → shot
- DT → an | my
- PRP → I
- PRP\$ → my
- IN → in

I shot an elephant in my pajamas

I	shot	an	elephant	in	my	pajamas
---	------	----	----------	----	----	---------

NP, PRP [0,1]	\emptyset	\emptyset	S [0,4]	\emptyset	\emptyset	
	VBD [1,2]	\emptyset	VP [1,4]	\emptyset	\emptyset	
		DT [2,3]	NP [2,4]	\emptyset	\emptyset	
			NP, NN [3,4]	\emptyset	\emptyset	
				IN [4,5]	\emptyset	PP [4,7]
					PRP\$ [5,6]	NP [5,7]
						NNS [6,7]

$S \rightarrow NP\ VP$
 $VP \rightarrow VBD\ NP$
 $VP \rightarrow VP\ PP$
 $Nominal \rightarrow Nominal\ PP$
 $Nominal \rightarrow pajamas\ |\ elephant\ |\ I$
 $PP \rightarrow IN\ NP$
 $NP \rightarrow DT\ NN$
 $NP \rightarrow pajamas\ |\ elephant\ |\ I$
 $NP \rightarrow PRP\$ Nominal$

$VBD \rightarrow shot$
 $DT \rightarrow an\ |\ my$
 $PRP \rightarrow I$
 $PRP\$ \rightarrow my$
 $IN \rightarrow in$

I shot an elephant in my pajamas

I	shot	an	elephant	in	my	pajamas
---	------	----	----------	----	----	---------

NP, PRP [0,1]	\emptyset	\emptyset	S [0,4]	\emptyset	\emptyset	
	VBD [1,2]	\emptyset	VP [1,4]	\emptyset	\emptyset	
		DT [2,3]	NP [2,4]	\emptyset	\emptyset	
			NP, NN [3,4]	\emptyset	\emptyset	NP [3,7]
				IN [4,5]	\emptyset	PP [4,7]
					PRP\$ [5,6]	NP [5,7]
						NNS [6,7]

$S \rightarrow NP\ VP$
 $VP \rightarrow VBD\ NP$
 $VP \rightarrow VP\ PP$
 $Nominal \rightarrow Nominal\ PP$
 $Nominal \rightarrow pajamas\ |\ elephant\ |\ I$
 $PP \rightarrow IN\ NP$
 $NP \rightarrow DT\ NN$
 $NP \rightarrow pajamas\ |\ elephant\ |\ I$
 $NP \rightarrow PRP\$ Nominal$

$VBD \rightarrow shot$
 $DT \rightarrow an\ |\ my$
 $PRP \rightarrow I$
 $PRP\$ \rightarrow my$
 $IN \rightarrow in$

I shot an elephant in my pajamas

I	shot	an	elephant	in	my	pajamas
---	------	----	----------	----	----	---------

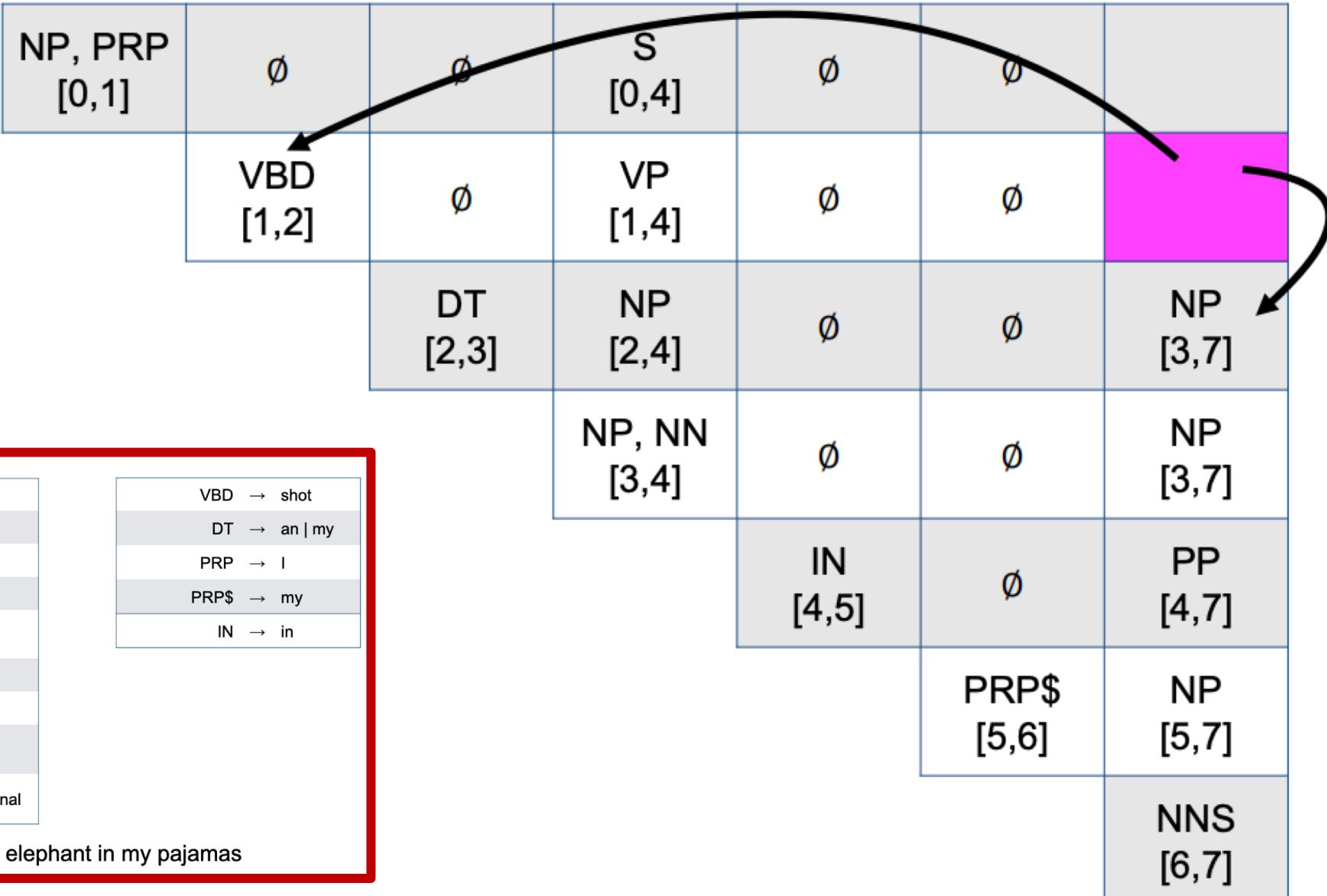
NP, PRP [0,1]	\emptyset	\emptyset	S [0,4]	\emptyset	\emptyset	
	VBD [1,2]	\emptyset	VP [1,4]	\emptyset	\emptyset	
		DT [2,3]	NP [2,4]	\emptyset	\emptyset	NP [3,7]
			NP, NN [3,4]	\emptyset	\emptyset	NP [3,7]
				IN [4,5]	\emptyset	PP [4,7]
					PRP\$ [5,6]	NP [5,7]
						NNS [6,7]

$S \rightarrow NP\ VP$
 $VP \rightarrow VBD\ NP$
 $VP \rightarrow VP\ PP$
 $Nominal \rightarrow Nominal\ PP$
 $Nominal \rightarrow pajamas\ |\ elephant\ |\ I$
 $PP \rightarrow IN\ NP$
 $NP \rightarrow DT\ NN$
 $NP \rightarrow pajamas\ |\ elephant\ |\ I$
 $NP \rightarrow PRP\$ Nominal$

$VBD \rightarrow shot$
 $DT \rightarrow an\ |\ my$
 $PRP \rightarrow I$
 $PRP\$ \rightarrow my$
 $IN \rightarrow in$

I shot an elephant in my pajamas

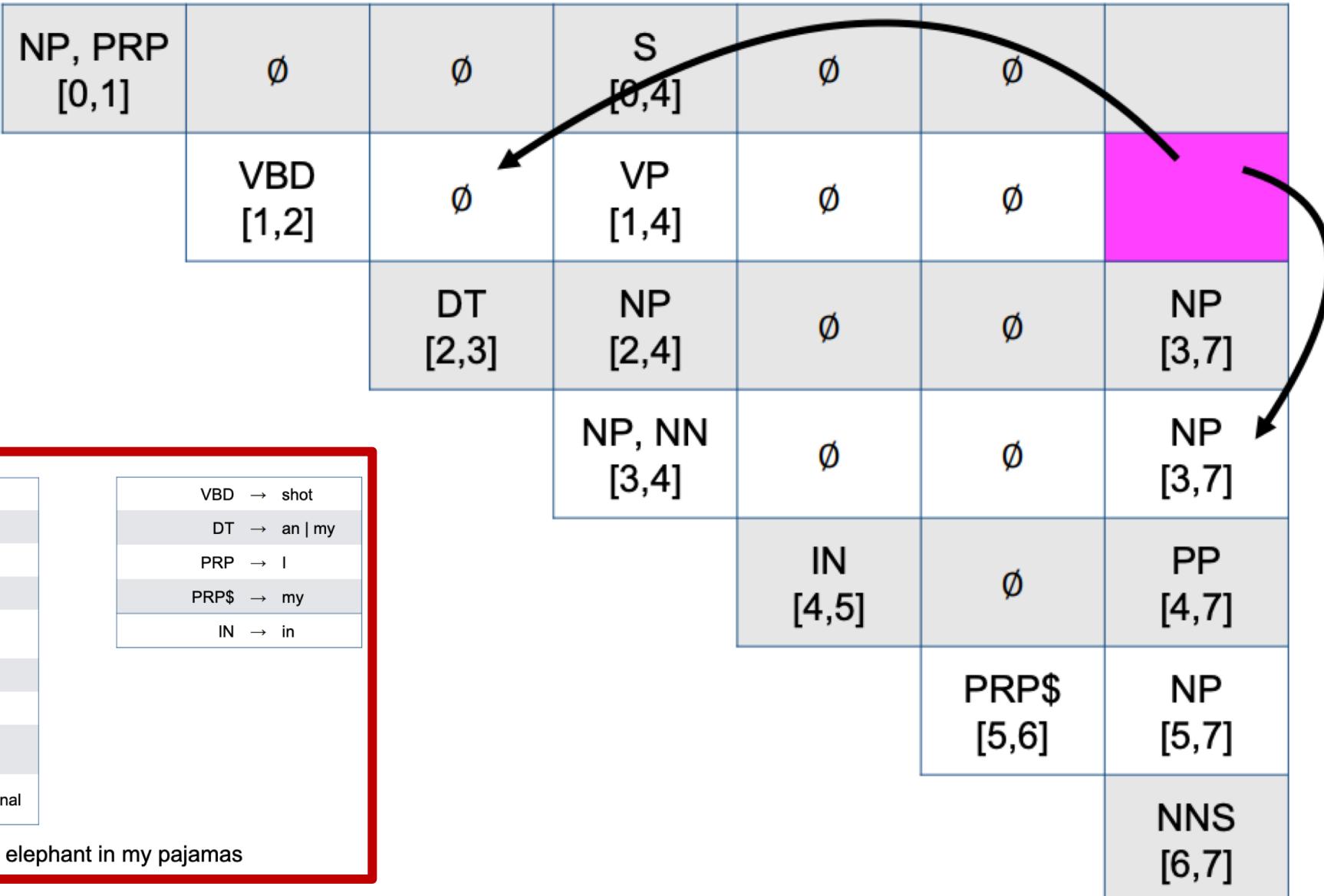
I	shot	an	elephant	in	my	pajamas
---	------	----	----------	----	----	---------



S → NP VP
VP → VBD NP
VP → VP PP
Nominal → Nominal PP
Nominal → pajamas elephant I
PP → IN NP
NP → DT NN
NP → pajamas elephant I
NP → PRP\$ Nominal

VBD → shot
DT → an my
PRP → I
PRP\$ → my
IN → in

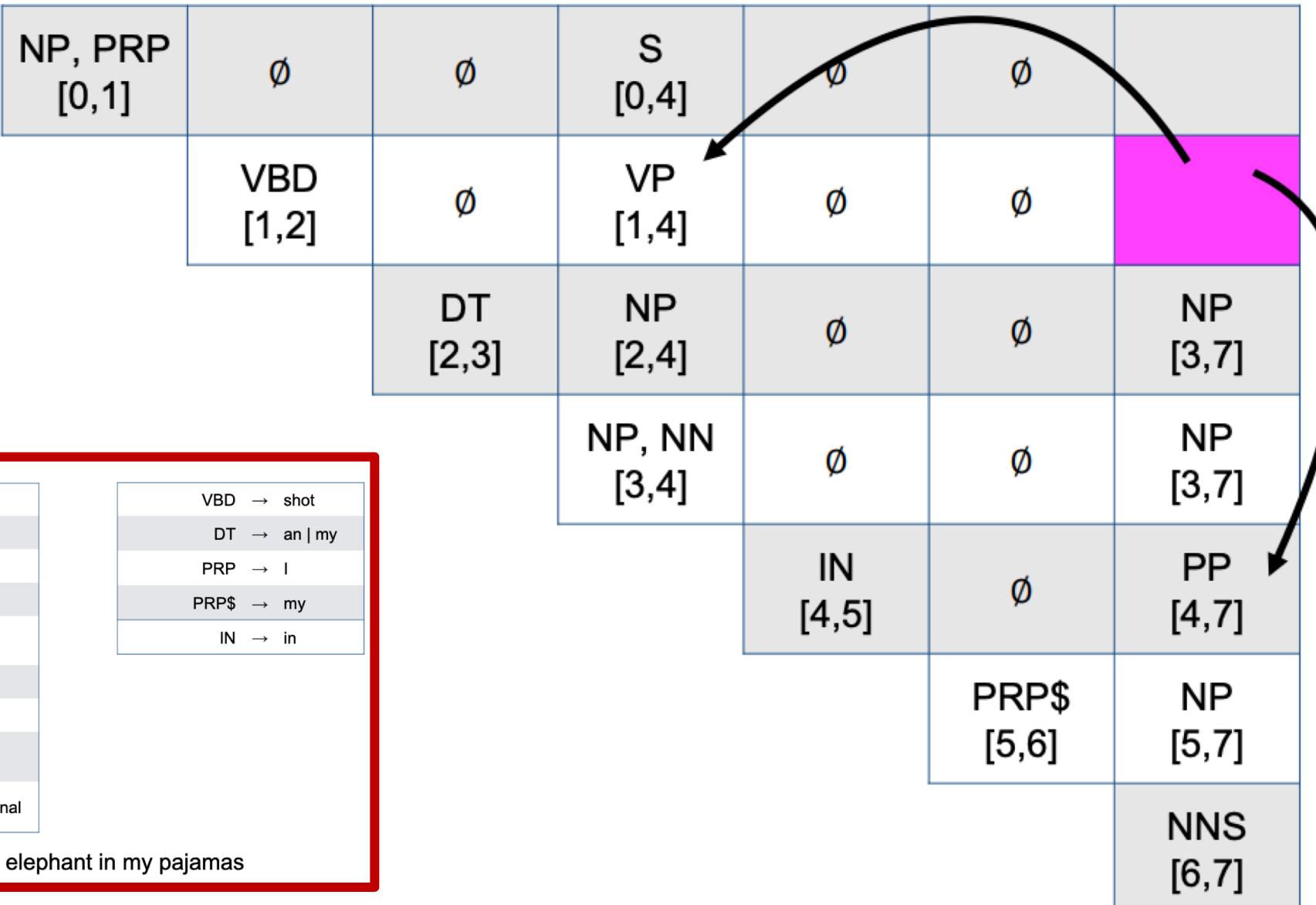
I	shot	an	elephant	in	my	pajamas
---	------	----	----------	----	----	---------



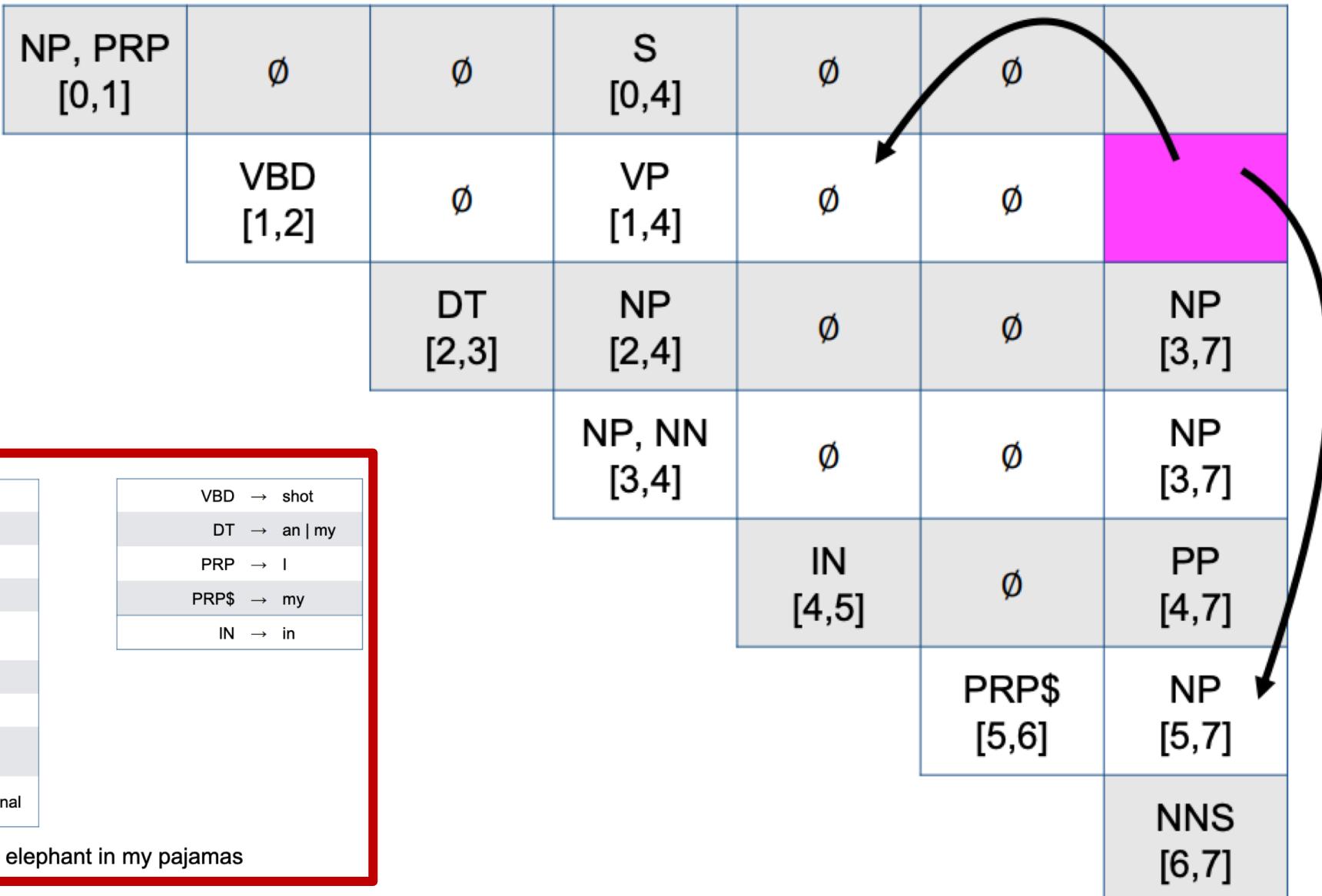
S → NP VP
VP → VBD NP
VP → VP PP
Nominal → Nominal PP
Nominal → pajamas elephant I
PP → IN NP
NP → DT NN
NP → pajamas elephant I
NP → PRP\$ Nominal

VBD → shot
DT → an my
PRP → I
PRP\$ → my
IN → in

I	shot	an	elephant	in	my	pajamas
---	------	----	----------	----	----	---------



I	shot	an	elephant	in	my	pajamas
---	------	----	----------	----	----	---------



I	shot	an	elephant	in	my	pajamas
---	------	----	----------	----	----	---------

NP, PRP [0,1]	\emptyset	\emptyset	S [0,4]	\emptyset	\emptyset	
	VBD [1,2]	\emptyset	VP [1,4]	\emptyset	\emptyset	
		DT [2,3]	NP [2,4]	\emptyset	\emptyset	NP [3,7]
			NP, NN [3,4]	\emptyset	\emptyset	NP [3,7]
				IN [4,5]	\emptyset	PP [4,7]
					PRP\$ [5,6]	NP [5,7]
						NNS [6,7]

Diagram illustrating the parse tree for the sentence "I shot an elephant in my pajamas". The tokens are shown in the first row, and the parse tree structure is shown in the subsequent rows. The tokens are: I, shot, an, elephant, in, my, pajamas.

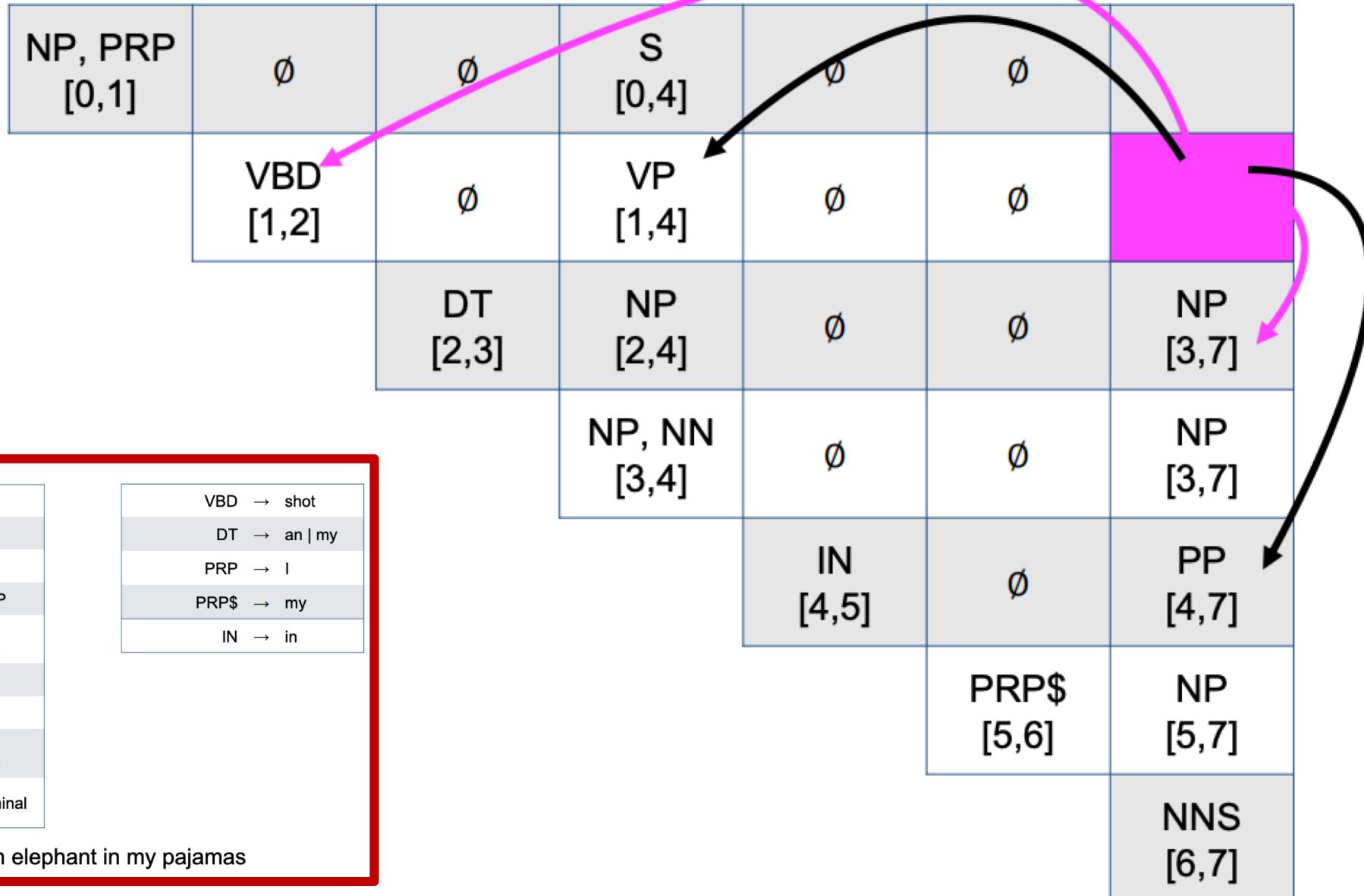
The parse tree structure is as follows:

- S → NP VP
- NP → VBD NP
- VP → VBD NP
- NP → DT NN
- NP → Nominal PP
- NP → Nominal | elephant | I
- NP → Nominal | pajamas | elephant | I
- NP → PRP\$ Nominal
- VBD → shot
- DT → an | my
- PRP → I
- PRP\$ → my
- IN → in

A red box highlights the first four rules of the grammar, which are used to derive the tokens. A pink box highlights the NP node for "pajamas", which is further expanded by a curved arrow to show its internal structure: NNS [6,7].

I shot an elephant in my pajamas

I	shot	an	elephant	in	my	pajamas
---	------	----	----------	----	----	---------



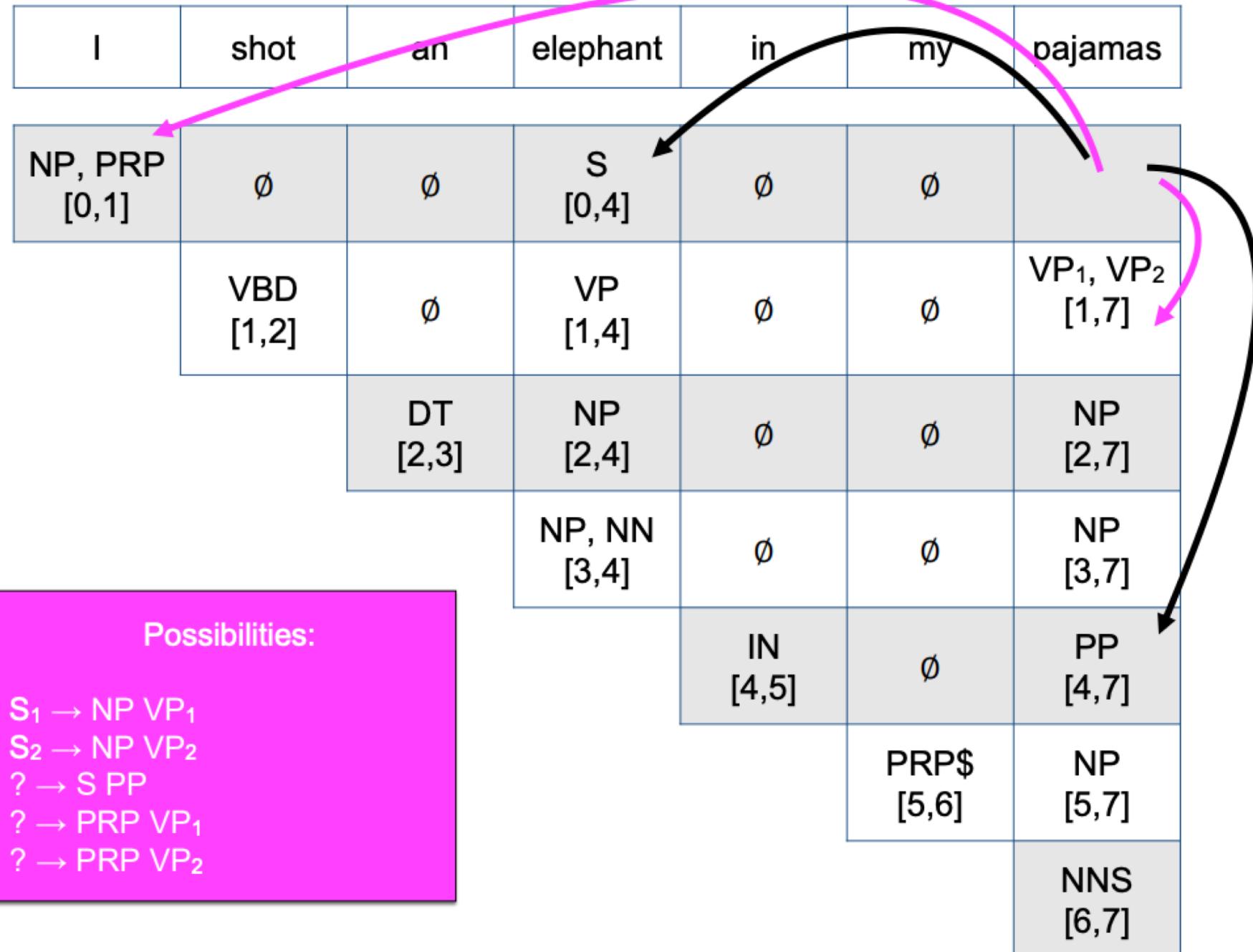
I	shot	an	elephant	in	my	pajamas
---	------	----	----------	----	----	---------

NP, PRP [0,1]	Ø	Ø	S [0,4]	Ø	Ø	
	VBD [1,2]	Ø	VP [1,4]	Ø	Ø	VP ₁ , VP ₂ [1,7]
		DT [2,3]	NP [2,4]	Ø	Ø	NP [2,7]
			NP, NN [3,4]	Ø	Ø	NP [3,7]
				IN [4,5]	Ø	PP [4,7]
					PRP\$ [5,6]	NP [5,7]
						NNS [6,7]

S → NP VP
 VP → VBD NP
 VP → VP PP
 Nominal → Nominal PP
 Nominal → pajamas | elephant | I
 PP → IN NP
 NP → DT NN
 NP → pajamas | elephant | I
 NP → PRP\$ Nominal

VBD → shot
 DT → an | my
 PRP → I
 PRP\$ → my
 IN → in

I shot an elephant in my pajamas



I	shot	an	elephant	in	my	pajamas
---	------	----	----------	----	----	---------

NP, PRP [0,1]	Ø	Ø	S [0,4]	Ø	Ø	S ₁ , S ₂ [0,7]
	VBD [1,2]	Ø	VP [1,4]	Ø	Ø	VP ₁ , VP ₂ [1,7]
		DT [2,3]	NP [2,4]	Ø	Ø	NP [2,7]
			NP, NN [3,4]	Ø	Ø	NP [3,7]
				IN [4,5]	Ø	PP [4,7]
					PRP\$ [5,6]	NP [5,7]
						NNS [6,7]

Success! We've recognized a total of two valid parses

CKY algorithm

```
function CKY-PARSE(words, grammar) returns table
    for j  $\leftarrow$  from 1 to LENGTH(words) do
        for all {A | A  $\rightarrow$  words[j]  $\in$  grammar} do
            table[j - 1, j]  $\leftarrow$  table[j - 1, j]  $\cup$  A
    for i  $\leftarrow$  from j - 2 downto 0 do
        for k  $\leftarrow$  i + 1 to j - 1 do
            for all {A | A  $\rightarrow$  BC  $\in$  grammar and B  $\in$  table[i, k] and C  $\in$  table[k, j]} do
                table[i, j]  $\leftarrow$  table[i, j]  $\cup$  A
```

Figure 12.5 The CKY algorithm.

I	shot	an	elephant	in	my	pajamas
NP, PRP [0,1]	∅	∅	S [0,4]	∅	∅	S ₁ , S ₂ [0,7]
	VBD [1,2]	∅	VP [1,4]	∅	∅	VP ₁ , VP ₂ [1,7]
	DT [2,3]	NP [2,4]		∅	∅	NP [2,7]
		NP, NN [3,4]		∅	∅	NP [3,7]
			IN [4,5]		∅	PP [4,7]
				PRP\$ [5,6]		NP [5,7]
					NNS [6,7]	

Runtime complexity?

CKY

- This use of CKY allows us to:
 - **Recognize** whether a sentence is grammatical in the language defined by the CFG
 - Enumerate all possible parses for a sentence
- But it does not tell us on its own which of those possible parses is most likely



- Under a full grammar, how many valid parses does an average-length sentence have?
 - a) 1-10
 - b) 10-100
 - c) 100-1000
 - d) 1000+

CKY Summary & In Practice

- Dynamic programming parsing algorithms, such as CKY, use a table of partial parses to efficiently parse ambiguous sentences.
- CKY restricts the form of the grammar to Chomsky normal form (CNF)
- In Practice
 - Does not return trees that are consistently with given grammar
 - Conversion to CNF complicates any syntax-driven approach to semantic analysis

PCFG

- Probabilistic context-free grammar: each production is also associated with a probability
- This lets us calculate the probability of a parse for a given sentence; for a given parse tree T for sentence S comprised of n rules from R (each $A \rightarrow \beta$):

$$P(T, S) = \prod_i^n P(\beta | A)$$

PCFG

$$P(T, S) = \prod_i^n P(\beta | A)$$

N Finite set of non-terminal symbols

NP, VP, S

Σ Finite alphabet of terminal symbols

the, dog, a

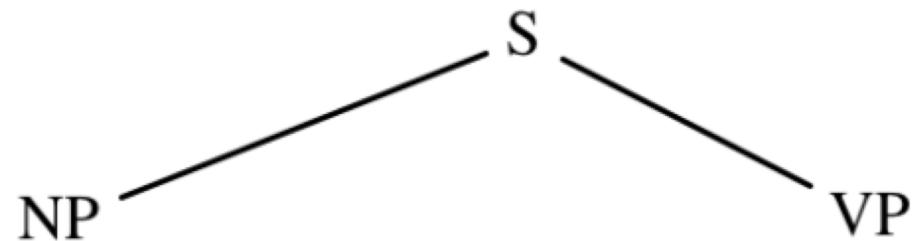
R Set of production rules, each
 $A \rightarrow \beta$ [p]
p = $P(\beta | A)$

$S \rightarrow NP\ VP$
Noun \rightarrow dog

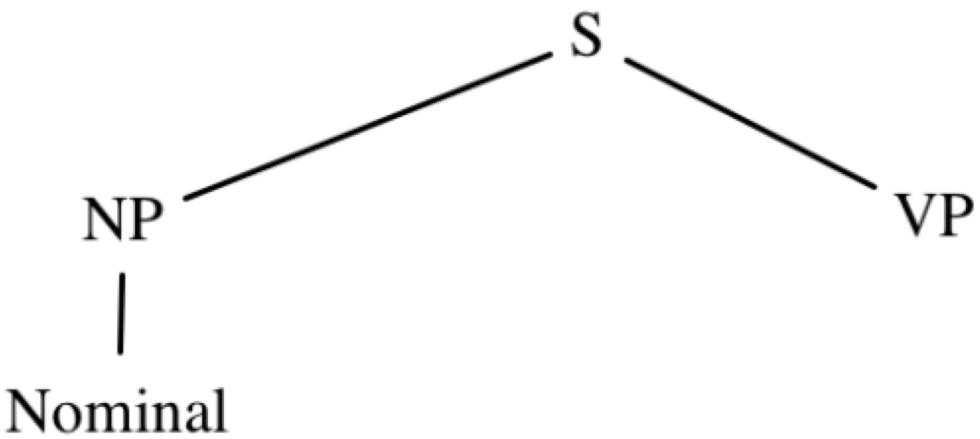
S Start symbol

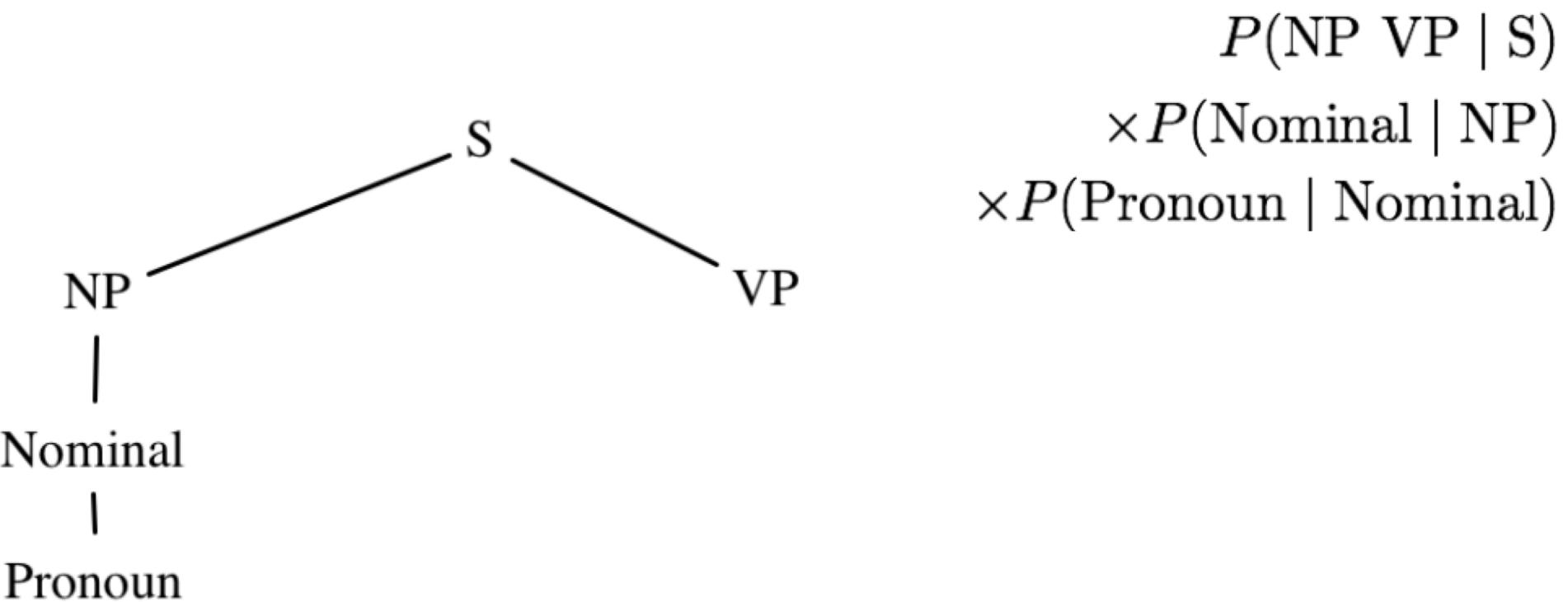
S

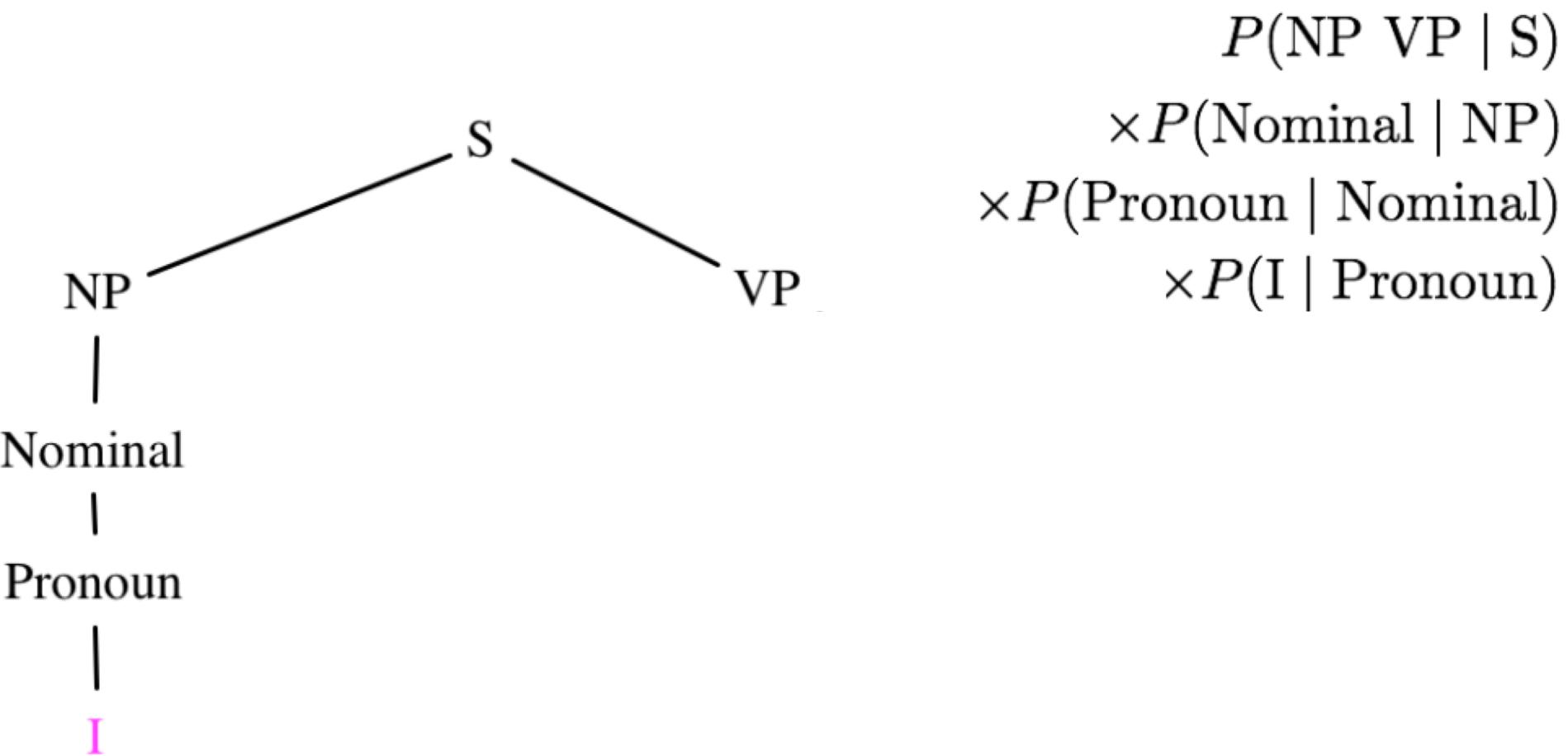
$$P(\text{NP } \text{VP} \mid \text{S})$$

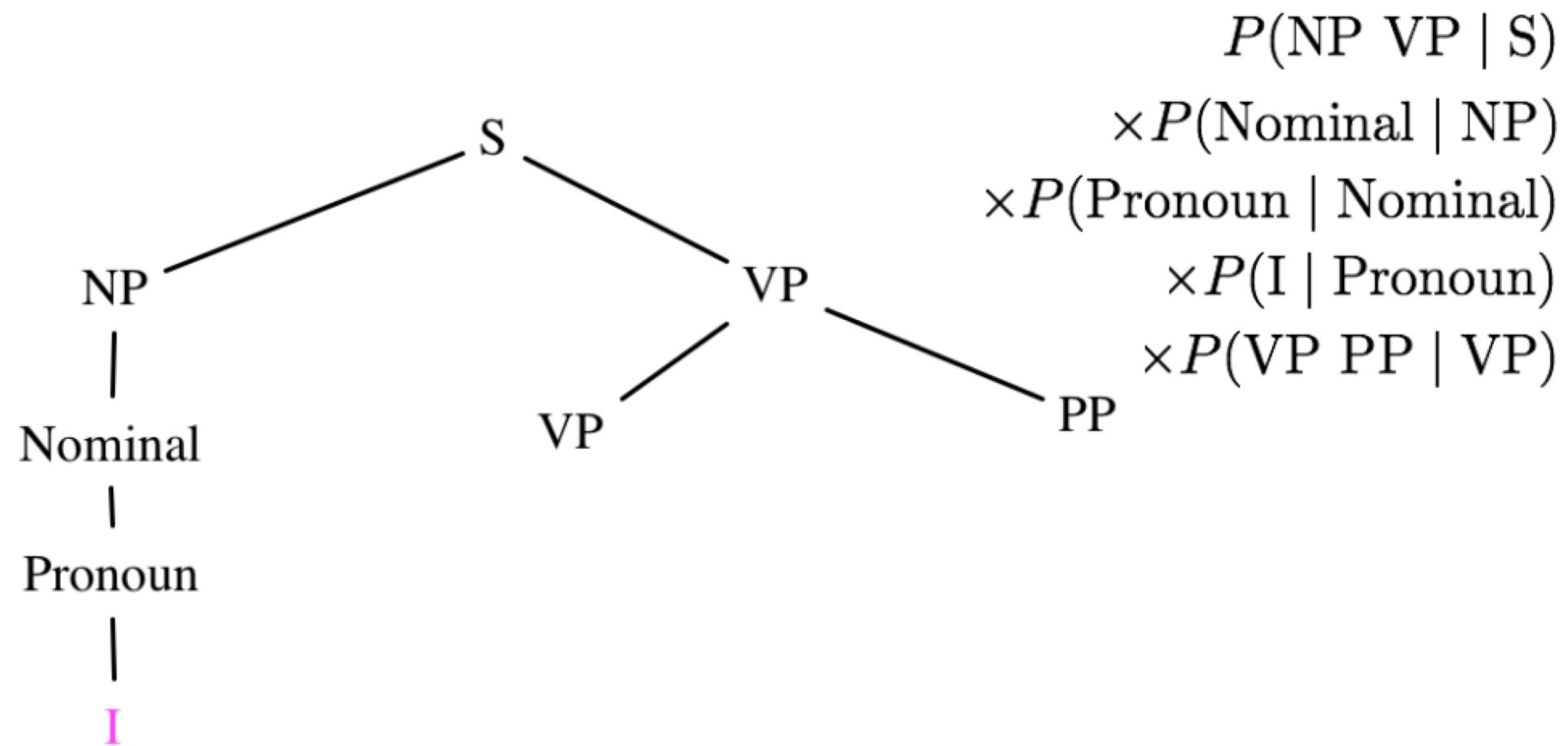


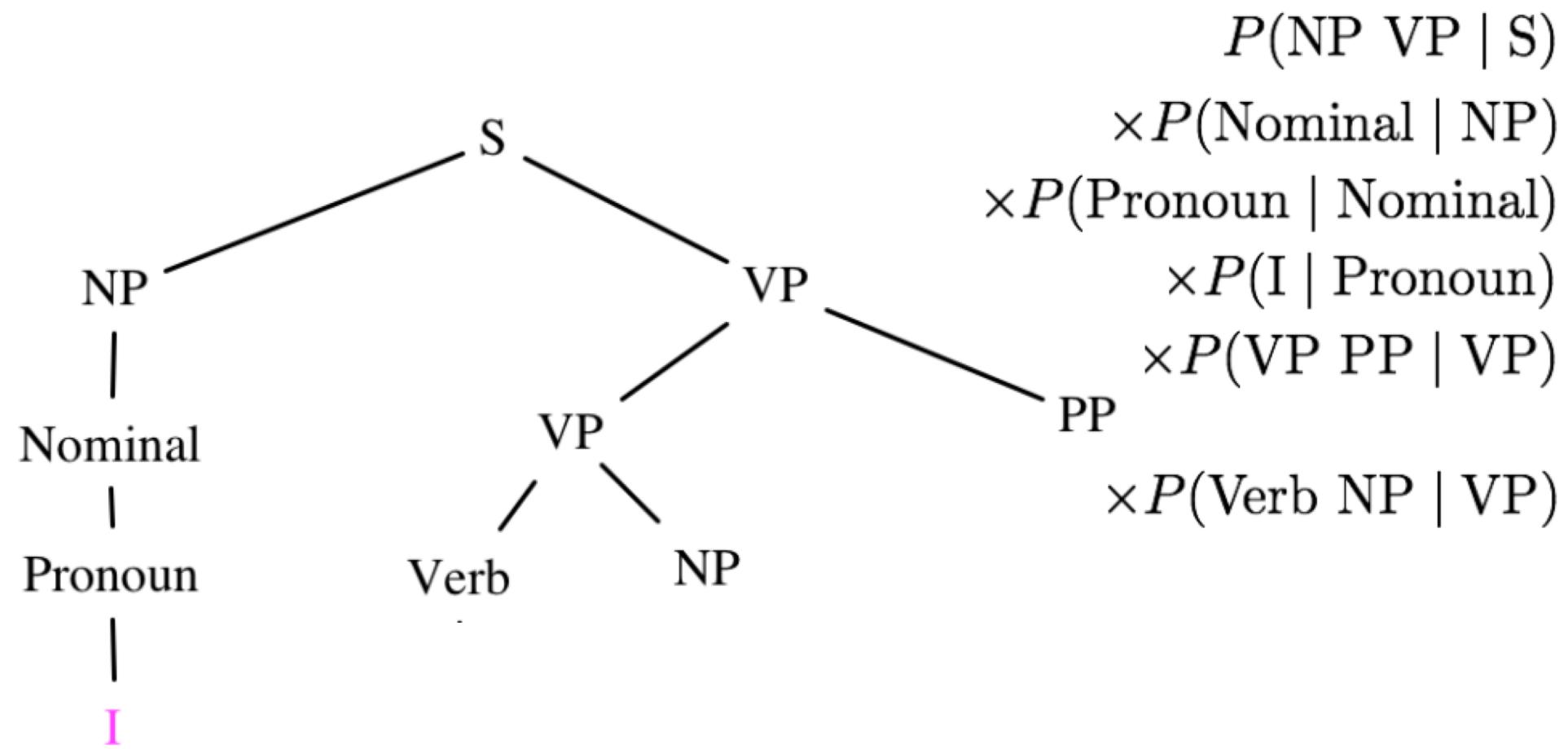
$$P(\text{NP VP} \mid \text{S}) \\ \times P(\text{Nominal} \mid \text{NP})$$

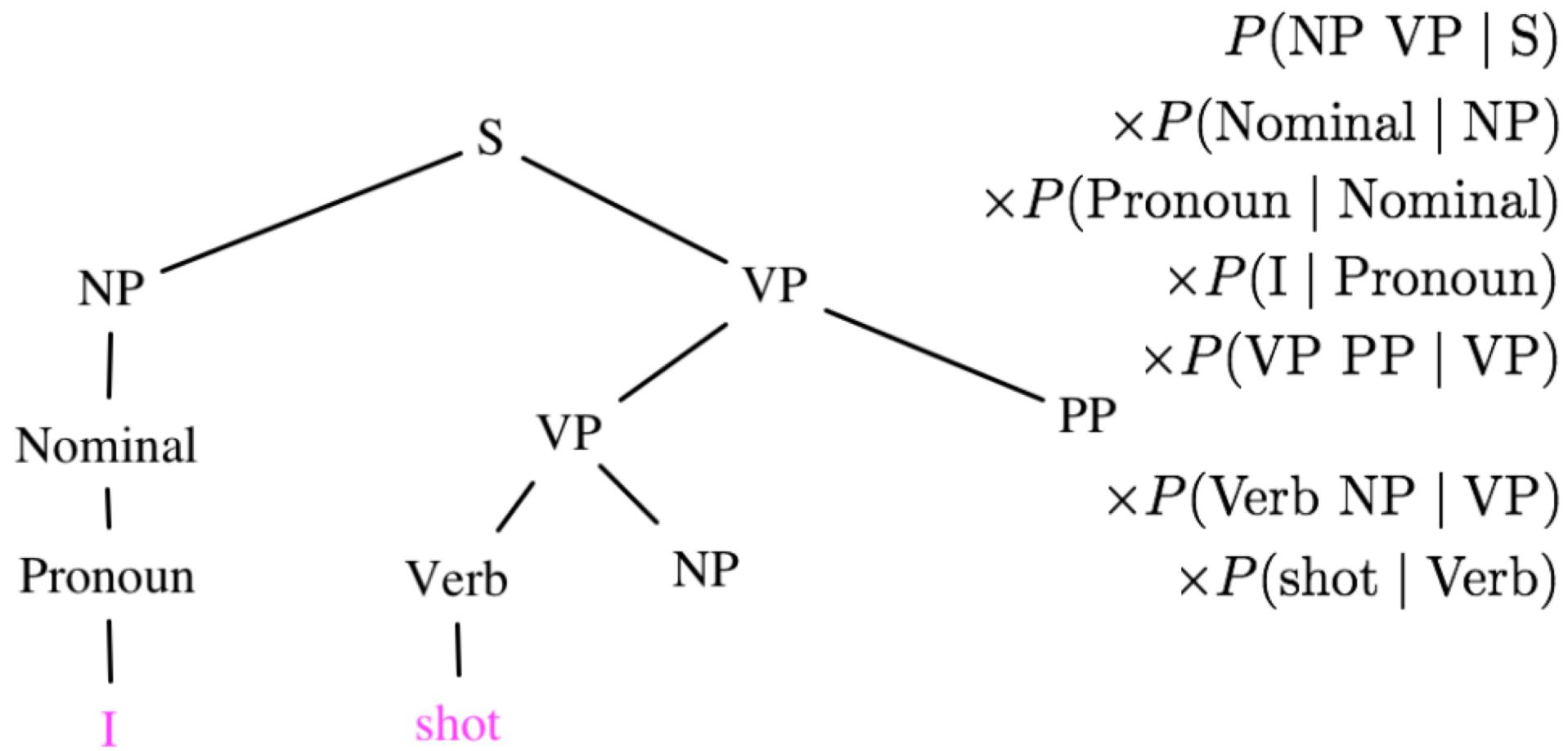


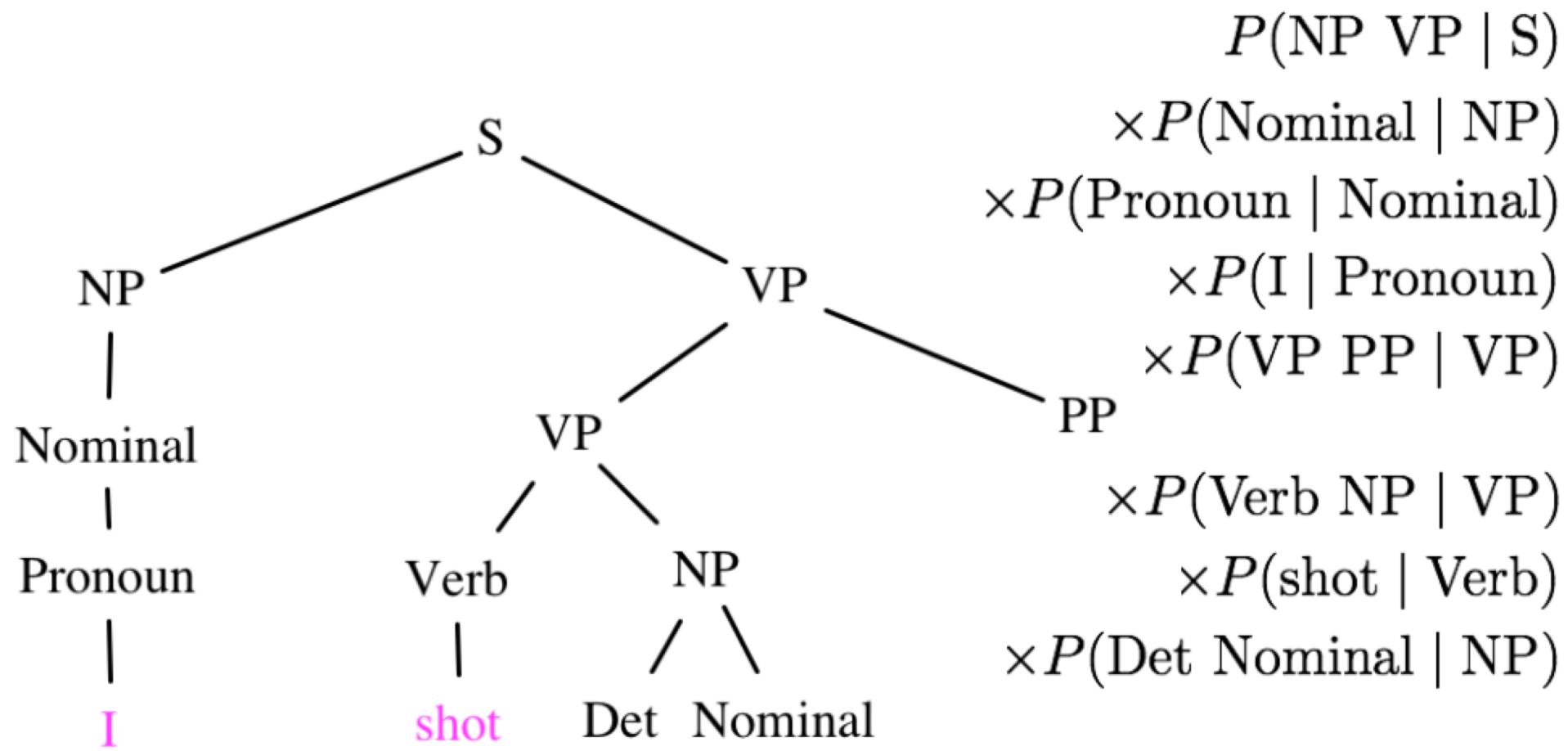


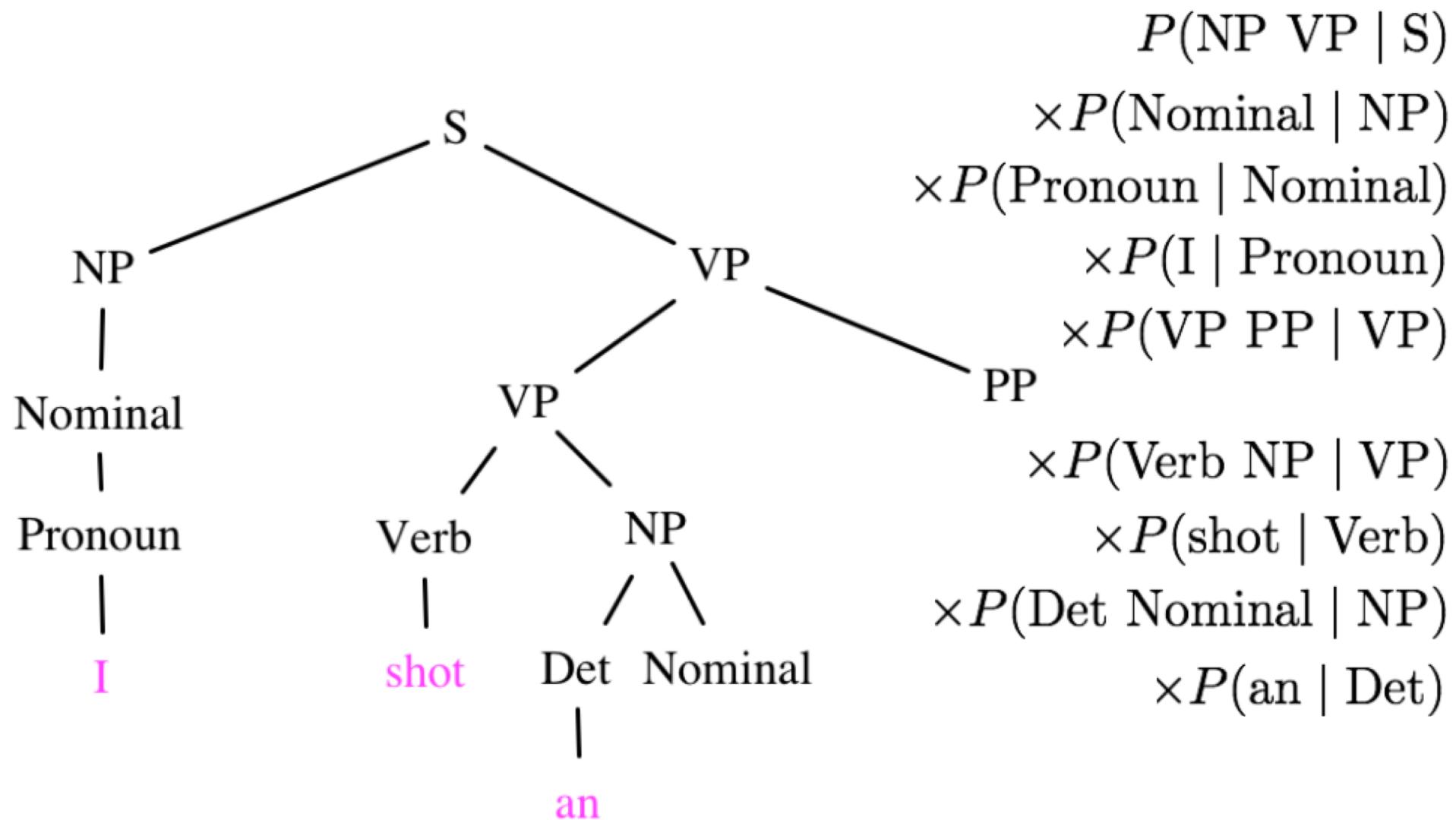


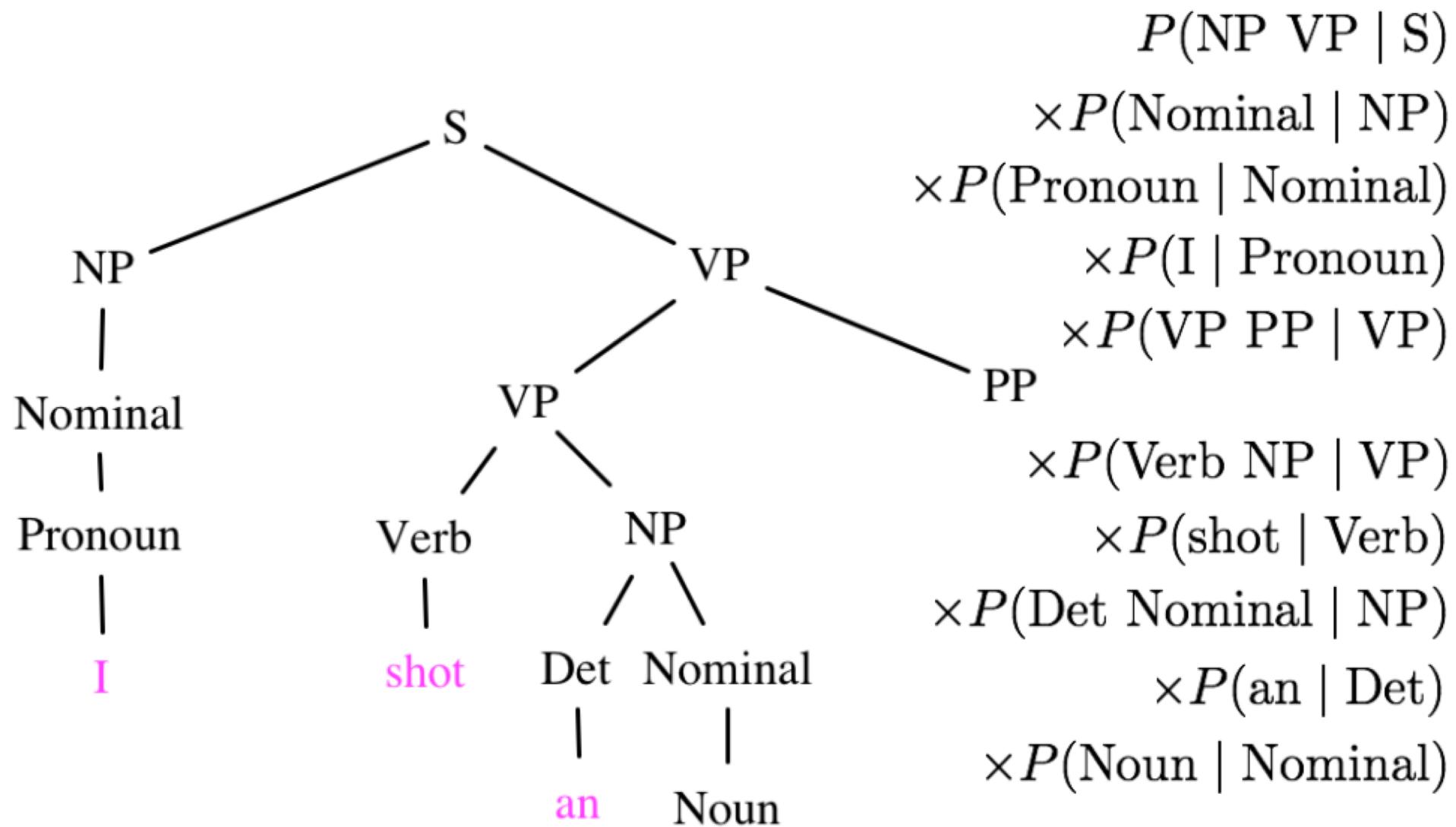


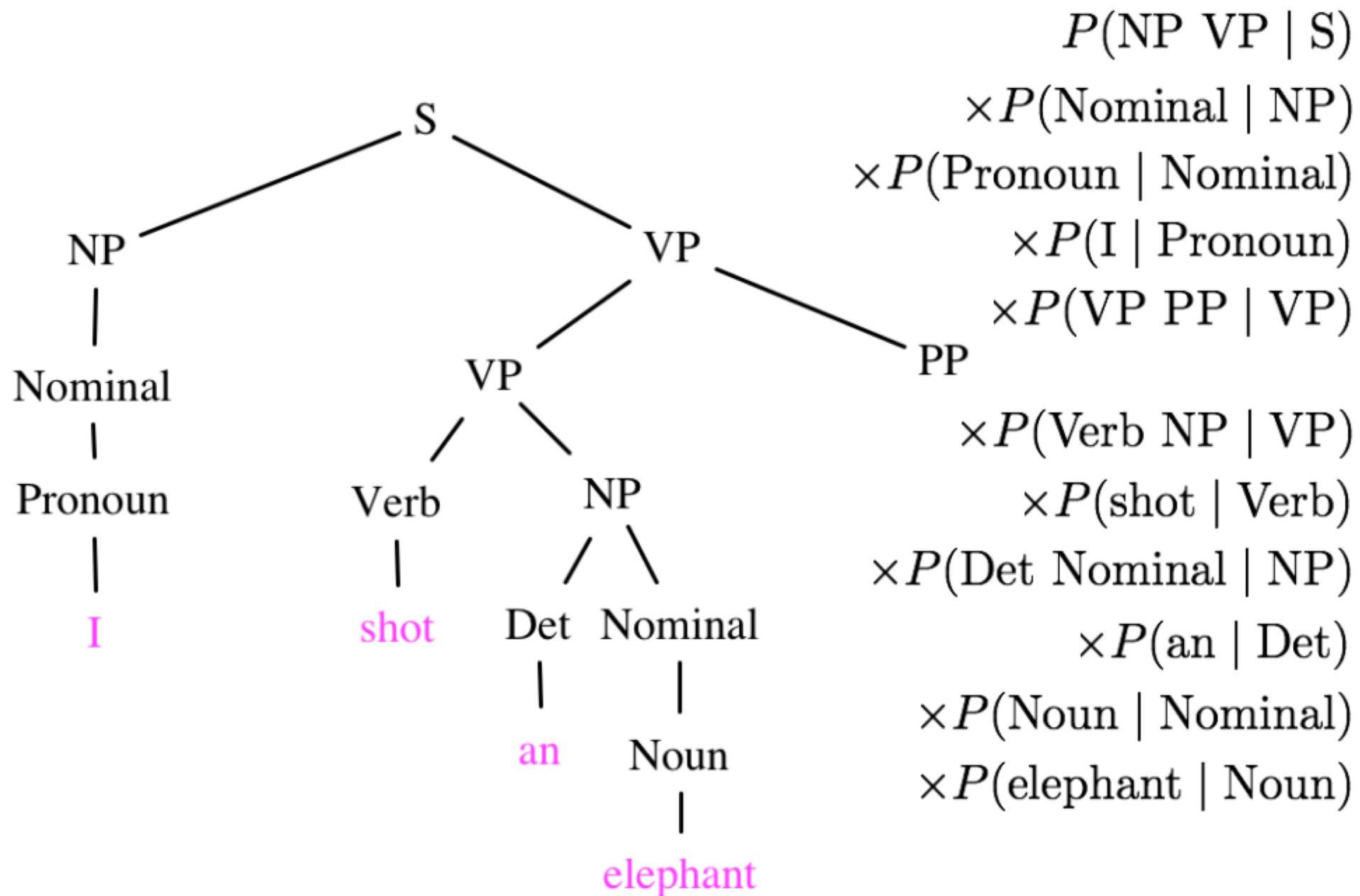


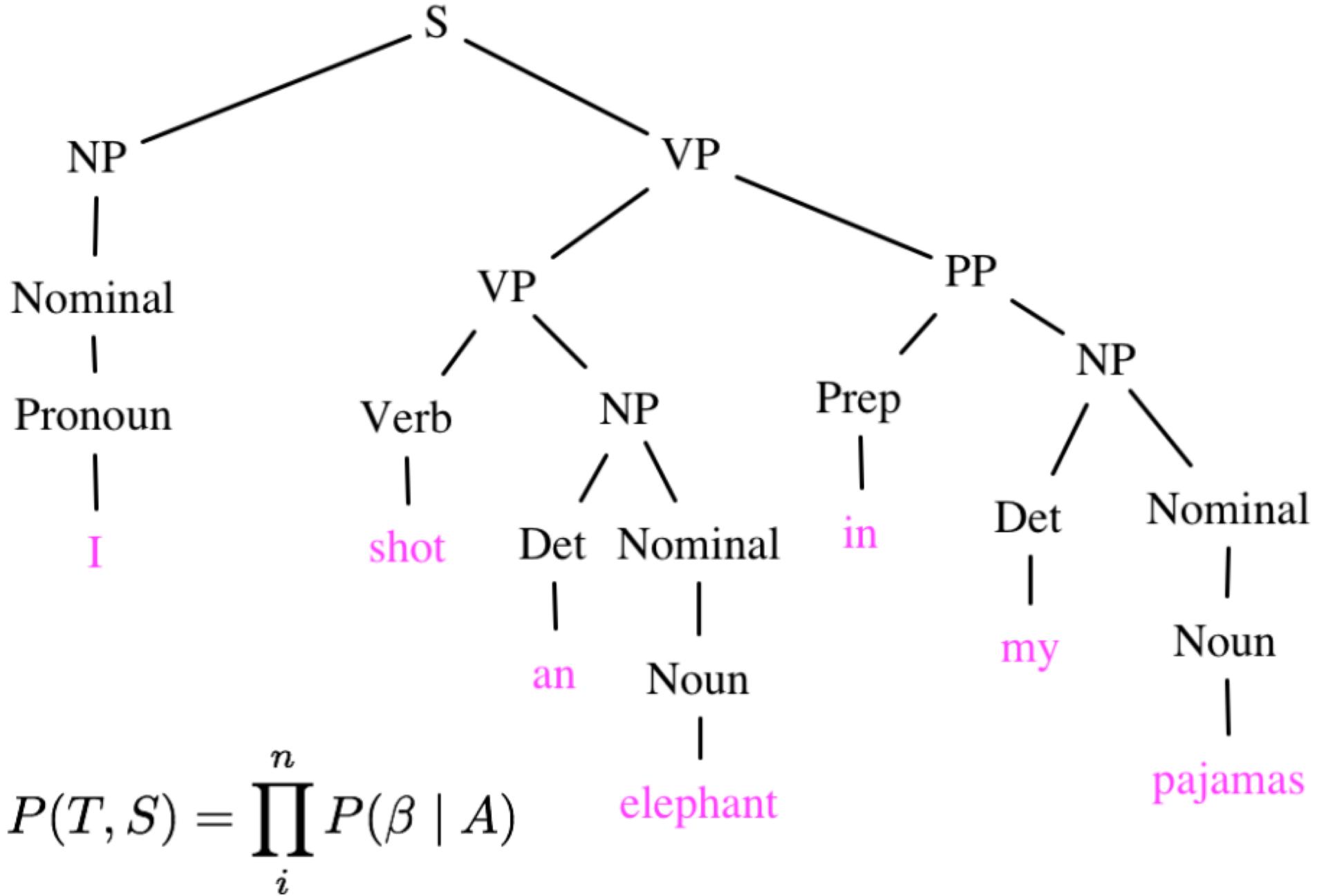












CKY PCFG

- A PCFG gives us a mechanism for assigning scores (here, probabilities) to different parses for the same sentence
- But we often care about is finding **the single best parse** with the highest probability
- In general, when people talk about CKY, they mean the PCFG version, since the CFG version isn't very useful

CKY PCFG

- We calculate the max probability parse using CKY by storing the probability of each phrase within each cell as we built it up.
- In particular, we fill the cell for span $i-j$ and label A with the maximum over splits and rules of

$$table(i, j, A) = P(A \rightarrow BC) \times table(i, k, B) \times table(k, j, C)$$

```

function PROBABILISTIC-CKY(words,grammar) returns most probable parse
                                         and its probability
for  $j \leftarrow$  from 1 to LENGTH(words) do
    for all {  $A \mid A \rightarrow words[j] \in grammar$  }
         $table[j-1, j, A] \leftarrow P(A \rightarrow words[j])$ 
    for  $i \leftarrow$  from  $j-2$  downto 0 do
        for  $k \leftarrow i+1$  to  $j-1$  do
            for all {  $A \mid A \rightarrow BC \in grammar,$ 
                         and  $table[i, k, B] > 0$  and  $table[k, j, C] > 0$  }
                if ( $table[i, j, A] < P(A \rightarrow BC) \times table[i, k, B] \times table[k, j, C]$ ) then
                     $table[i, j, A] \leftarrow P(A \rightarrow BC) \times table[i, k, B] \times table[k, j, C]$ 
                     $back[i, j, A] \leftarrow \{k, B, C\}$ 
    return BUILD-TREE( $back[1, \text{LENGTH}(words), S]$ ),  $table[1, \text{LENGTH}(words), S]$ 

```

VS. CKY

```
function CKY-PARSE(words, grammar) returns table
  for j  $\leftarrow$  from 1 to LENGTH(words) do
    for all {A | A  $\rightarrow$  words[j]  $\in$  grammar}
      table[j - 1, j]  $\leftarrow$  table[j - 1, j]  $\cup$  A
    for i  $\leftarrow$  from j - 2 downto 0 do
      for k  $\leftarrow$  i + 1 to j - 1 do
        for all {A | A  $\rightarrow$  BC  $\in$  grammar and B  $\in$  table[i, k] and C  $\in$  table[k, j]}
          table[i, j]  $\leftarrow$  table[i, j]  $\cup$  A
```

Figure 13.5 The CKY algorithm.

I	shot	an	elephant	in	my	pajamas
PRP:0.04 [0,1]						
	VBD:0.04 [1,2]					
		DT:0.05 [2,3]				
			NN:0.03 [3,4]			
Probability of a terminal (word) given its tag				IN:0.10 [4,5]		
$P(A \rightarrow \beta)$				PRP\$:0.1 2 [5,6]		
				NNS:0.01 [6,7]		

I	shot	an	elephant	in	my	pajamas
---	------	----	----------	----	----	---------

PRP:0.04 [0,1]	∅	∅				
	VBD:0.04 [1,2]	∅				
		DT:0.05 [2,3]	NP:0.0001 5 [2,4]			
			NN:0.03 [3,4]			
				IN:0.10 [4,5]		
					PRP\$:0.12 [5,6]	
						NNS:0.01 [6,7]

$$table(2, 4, NP) = P(NP \rightarrow DT\ NN) \times table(2, 3, DT) \times table(3, 4, NN)$$

I	shot	an	elephant	in	my	pajamas
PRP:0.04 [0,1]	∅	∅				
	VBD:0.04 [1,2]	∅	VP:0.0000 006 [1,4]			
		DT:0.05 [2,3]	NP:0.0001 5 [2,4]			
			NN:0.03 [3,4]			
				IN:0.10 [4,5]		
					PRP\$:0.12 [5,6]	
						NNS:0.01 [6,7]

We just calculated this value
and can use it now

$$table(1, 4, VP) = P(VP \rightarrow VBD\ NP) \times table(1, 2, VBD) \times \text{table}(2, 4, NP)$$

I	shot	an	elephant	in	my	pajamas
---	------	----	----------	----	----	---------

PRP:0.04 [0,1]	Ø	Ø	S: 0.0000000 048 [0,4]			
	VBD:0.04 [1,2]	Ø	VP:0.0000 006 [1,4]			
		DT:0.05 [2,3]	NP:0.0001 5 [2,4]			
			NN:0.03 [3,4]			
				IN:0.10 [4,5]		
					PRP\$:0.12 [5,6]	
						NNS:0.01 [6,7]

We just calculated this value
and can use it now

$$table(0, 4, S) = P(S \rightarrow NP\ VP) \times table(0, 1, NP) \times \text{table}(1, 4, VP)$$

I	shot	an	elephant	in	my	pajamas
PRP:0.04 [0,1]	∅	∅	S: 0.0000000 048 [0,4]			
	VBD:0.04 [1,2]	∅	VP:0.0000 006 [1,4]			
		DT:0.05 [2,3]	NP:0.0001 5 [2,4]			
			NN:0.03 [3,4]			
				IN:0.10 [4,5]		
					PRP\$:0.12 [5,6]	
						NNS:0.01 [6,7]

Note these values are getting very small! Better to add in log space

I	shot	an	elephant	in	my	pajamas
PRP: -3.21 [0,1]	∅	∅	S: -19.2 [0,4]			
	VBD: -3.21 [1,2]	∅	VP: -14.3 [1,4]			
		DT: -3.0 [2,3]	NP: -8.8 [2,4]			
			NN: -3.5 [3,4]			
				IN: -2.3 [4,5]		
					PRP\$: - 2.12 [5,6]	
						NNS: -4.6 [6,7]

Note these values are getting very small! Better to add in log space

I	shot	an	elephant	in	my	pajamas
PRP: -3.21 [0,1]	∅	∅	S: -19.2 [0,4]	∅	∅	
VBD: -3.21 [1,2]	∅	∅	VP: -14.3 [1,4]	∅	∅	VP ₁ , VP ₂ [1,7]
DT: -3.0 [2,3]	NP: -8.8 [2,4]	∅	∅	∅	NP: -24.7 [2,7]	
NN: -3.5 [3,4]	∅	∅	∅	∅	NP: -19.4 [3,7]	
IN: -2.3 [4,5]	PP: -13.6 [4,7]	∅	∅	∅	PP: -13.6 [4,7]	
PRP\$: - 2.12 [5,6]	NP: -9.0 [5,7]	∅	∅	∅	NP: -9.0 [5,7]	
NNS: -4.6 [6,7]						

For any phrase type spanning $[i,j]$, we only need to keep the max probability given the assumptions of a PCFG

I	shot	an	elephant	in	my	pajamas
PRP: -3.21 [0,1]	∅	∅	S: -19.2 [0,4]	∅	∅	
	VBD: -3.21 [1,2]	∅	VP: -14.3 [1,4]	∅	∅	VP: -30.2 [1,7]
		DT: -3.0 [2,3]	NP: -8.8 [2,4]	∅	∅	NP: -24.7 [2,7]
			NN: -3.5 [3,4]	∅	∅	NP: -19.4 [3,7]
				IN: -2.3 [4,5]	∅	PP: -13.6 [4,7]
					PRP\$: - 2.12 [5,6]	NP: -9.0 [5,7]
						NNS: -4.6 [6,7]

For any phrase type spanning [i,j], we only need to keep the max probability given the assumptions of a PCFG

I	shot	an	elephant	in	my	pajamas
---	------	----	----------	----	----	---------

PRP: -3.21 [0,1]	\emptyset	\emptyset	S: -19.2 [0,4]	\emptyset	\emptyset	S: -35.7 [0,7]
	VBD: -3.21 [1,2]	\emptyset	VP: -14.3 [1,4]	\emptyset	\emptyset	VP: -30.2 [1,7]
		DT: -3.0 [2,3]	NP: -8.8 [2,4]	\emptyset	\emptyset	NP: -24.7 [2,7]
			NN: -3.5 [3,4]	\emptyset	\emptyset	NP: -19.4 [3,7]
				IN: -2.3 [4,5]	\emptyset	PP: -13.6 [4,7]
					PRP\$: - 2.12 [5,6]	NP: -9.0 [5,7]
						NNS: -4.6 [6,7]

For any phrase type spanning [i,j], we only need to keep the max probability given the assumptions of a PCFG

