

CKY Parsing

Ling571

Deep Processing Approaches to NLP

January 11, 2017

Roadmap

- Motivation:
 - Inefficiencies of parsing-as-search
- Strategy: Dynamic Programming
- Chomsky Normal Form
 - Weak and strong equivalence
- CKY parsing algorithm

Bottom-Up Parsing

- Try to find all trees that span the input
 - Start with input string
 - Book that flight.
- Use all productions with current subtree(s) on RHS
 - E.g., $N \rightarrow \rightarrow \text{Book}$; $V \rightarrow \rightarrow \text{Book}$
- Stop when spanned by S (or no more rules apply)

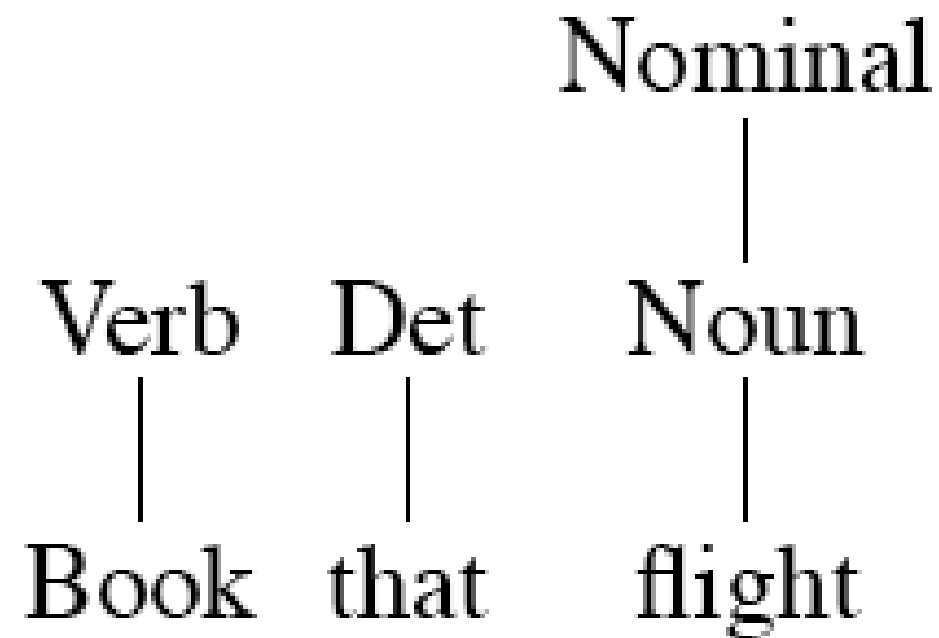
Bottom-Up Search

Book that flight

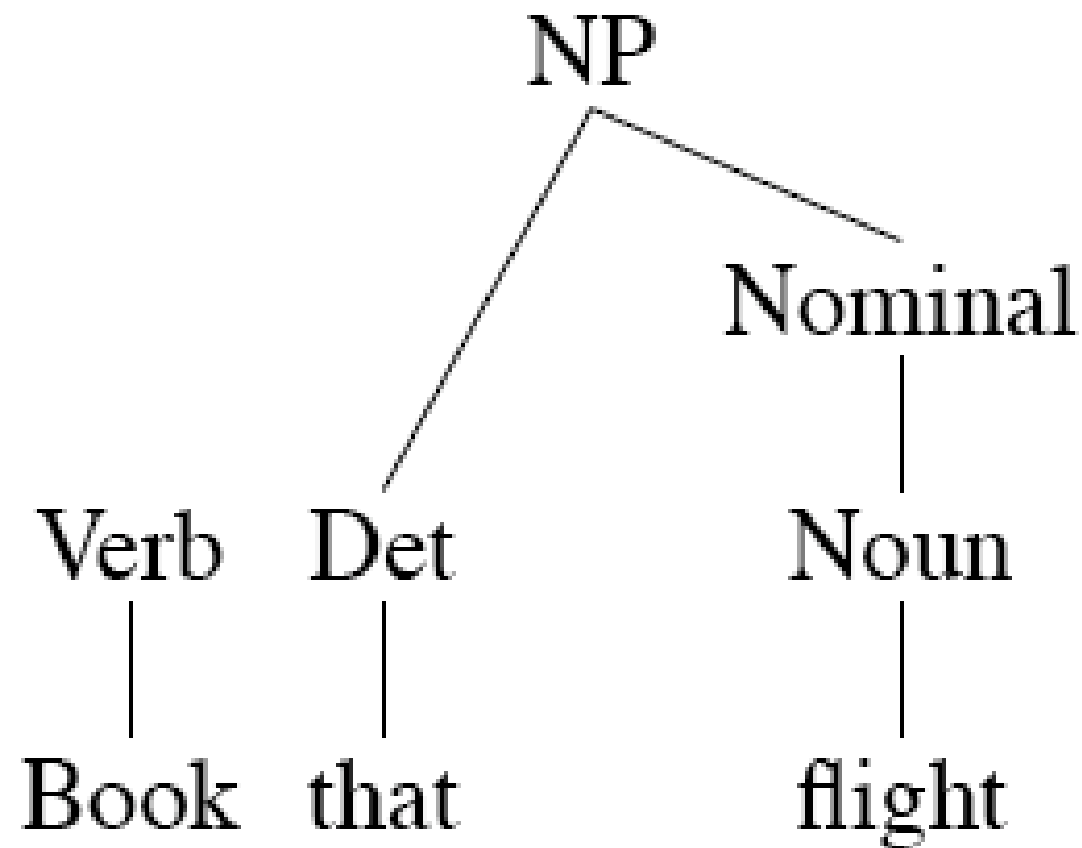
Bottom-Up Search

Verb	Det	Noun
Book	that	flight

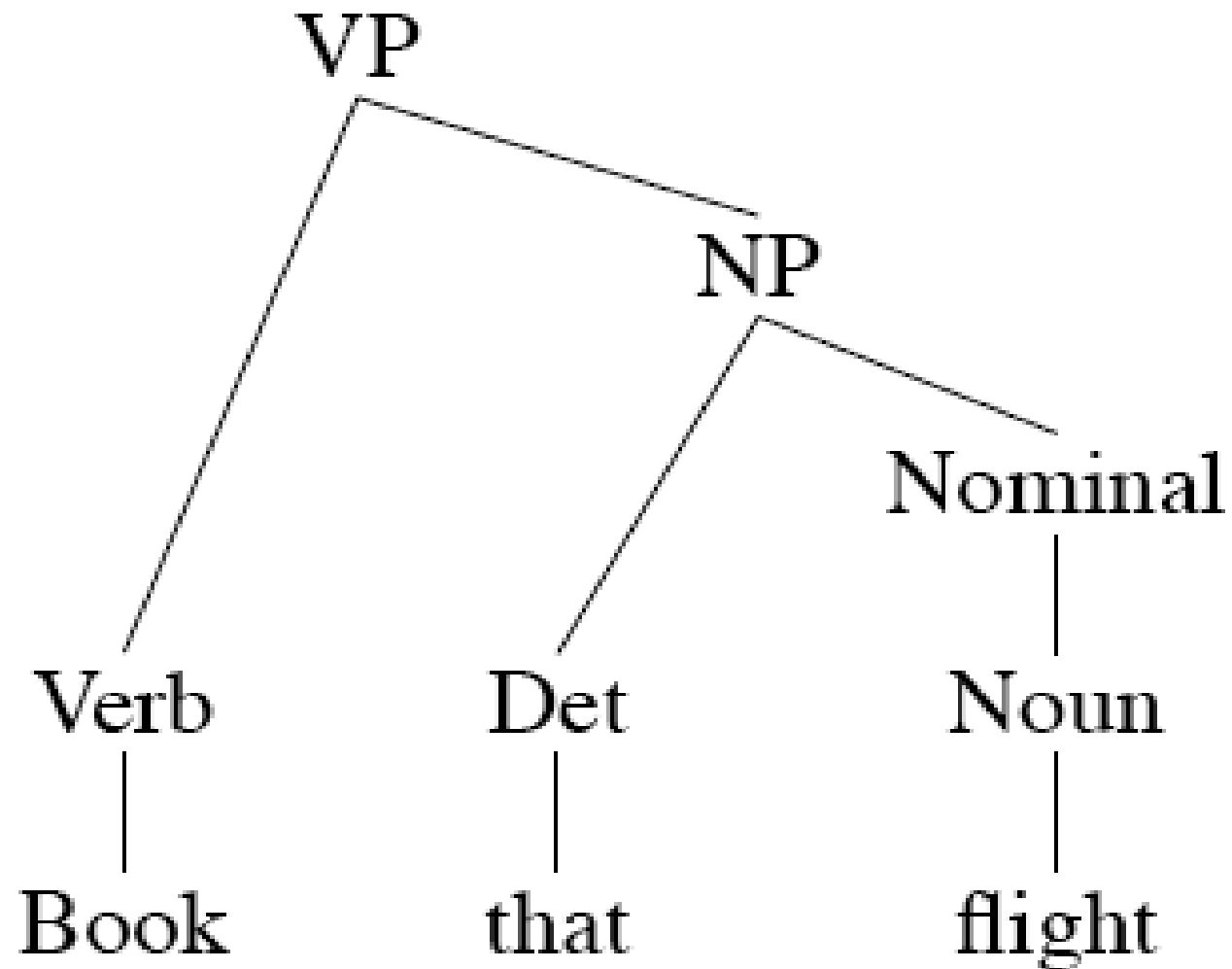
Bottom-Up Search



Bottom-Up Search



Bottom-Up Search



Veib
 |
 Book

Det
 |
 that

Noun
 |
 flight

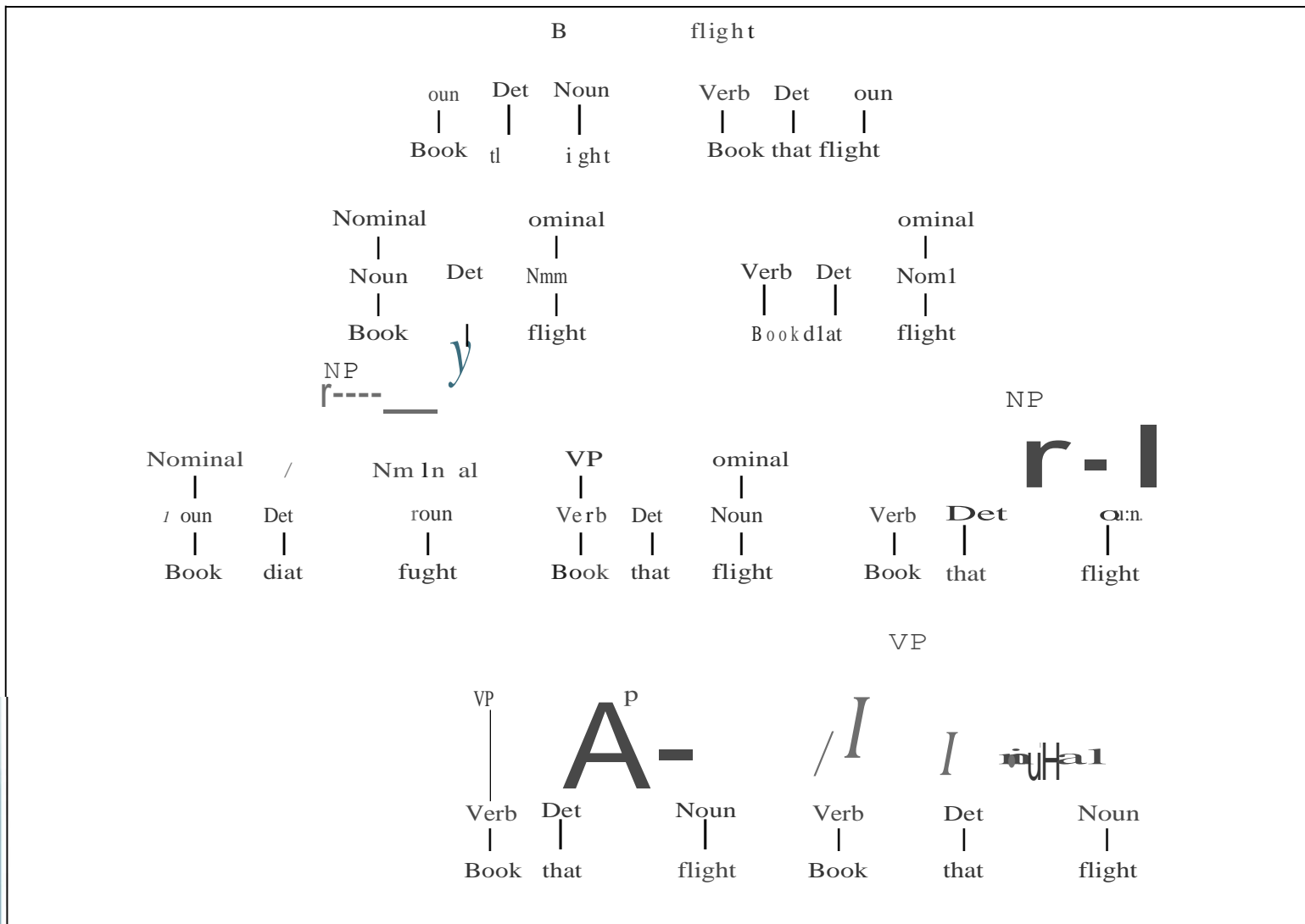
Verb
 |
 Book

Det
 |
 that

roun
 |
 flight

The diagram shows a syntax tree for the sentence "Book that flight Book that flight". The root node is S, which branches into NP and VP. The first NP branches into Det ("the") and N ("book"). The VP branches into V ("is") and NP. The second NP branches into Det ("that") and N ("flight"). The third NP branches into Det ("the") and N ("book"). The fourth VP branches into V ("is") and NP. The fifth NP branches into Det ("that") and N ("flight"). The sixth NP branches into Det ("the") and N ("flight").

Bottom-Up Search



Pros and Cons of Bottom-Up Search

— Pros:

- Will not explore trees that don't match input
- Recursive rules less problematic
- Useful for incremental/ fragment parsing

— Cons:

- Explore subtrees that will not fit full sentences

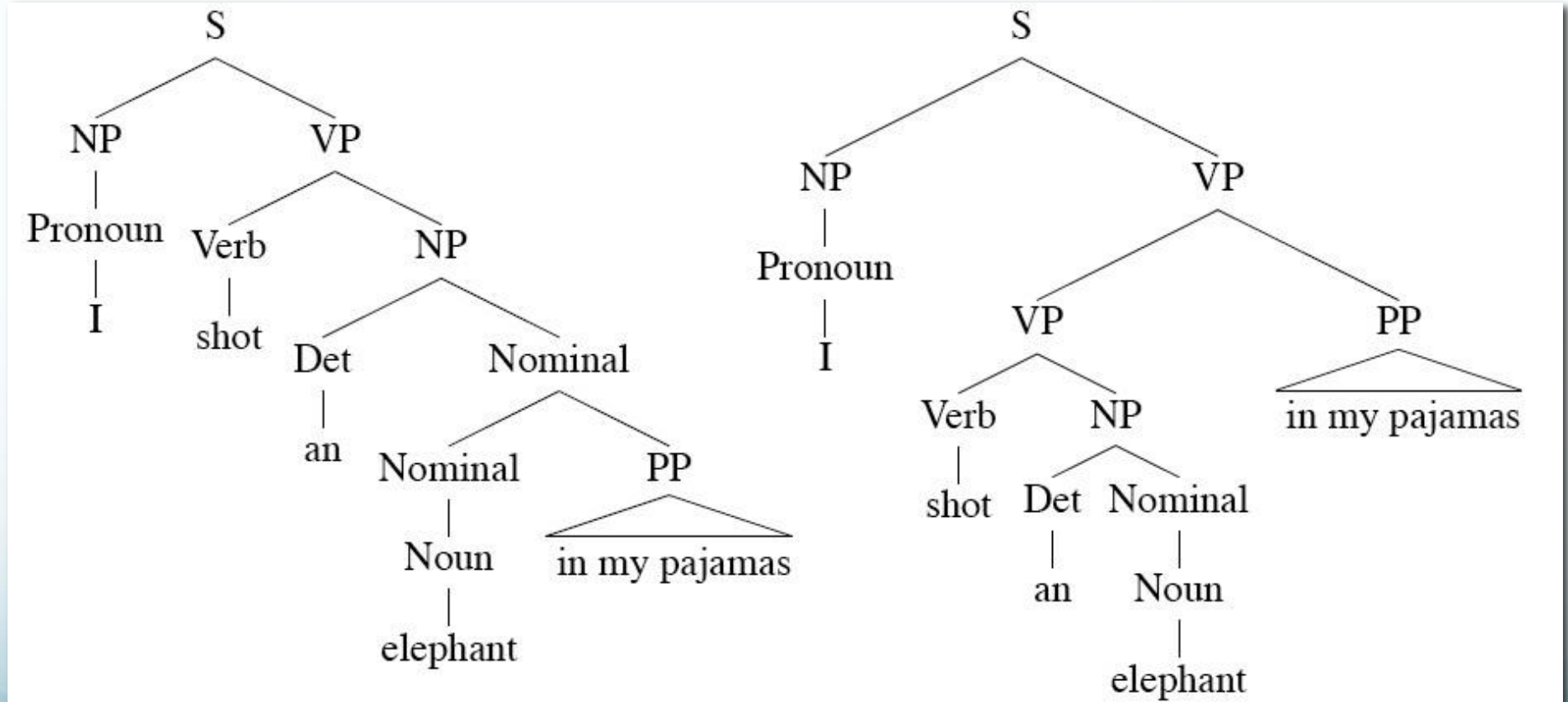
Parsing Challenges

- Ambiguity
- Repeated substructure
- Recursion

Parsing Ambiguity

- Many sources of parse ambiguity
 - Lexical ambiguity
 - Book/N; Book/V
 - Structural ambiguity: Main types:
 - Attachment ambiguity
 - Constituent can attach in multiple places
 - *I shot an elephant in my pyjamas.*
 - Coordination ambiguity
 - Different constituents can be conjoined
 - *Old men and women*

Ambiguity



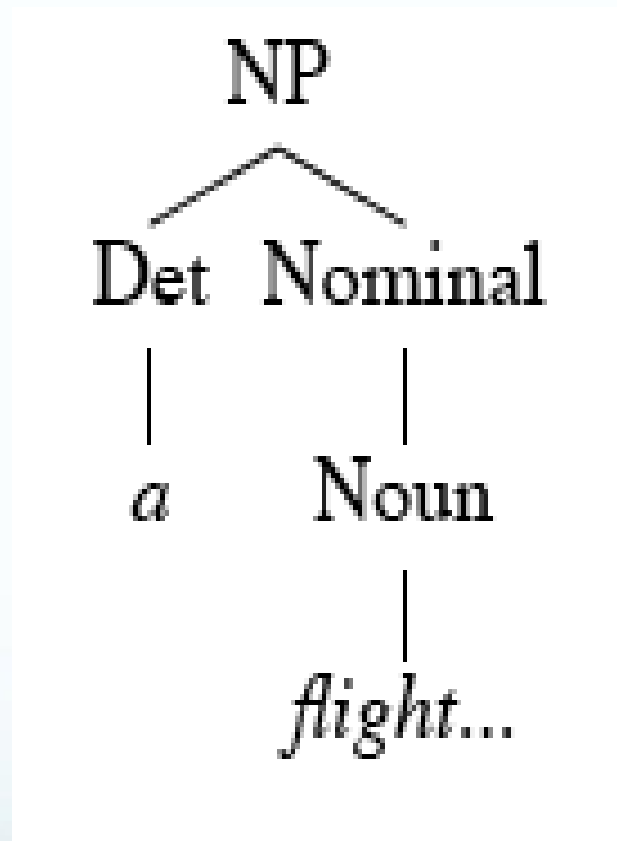
Disambiguation

- Global ambiguity:
 - Multiple complete alternative parses
 - Need strategy to select correct one
 - Approaches exploit other information
 - Statistical
 - Some prepositional structs more likely to attach high/low
 - Some phrases more likely, e.g., (old (men and women))
 - Semantic
 - Pragmatic
 - E.g., elephants and pyjamas
 - Alternatively, keep all
- Local ambiguity:
 - Ambiguity in subtree, resolved globally

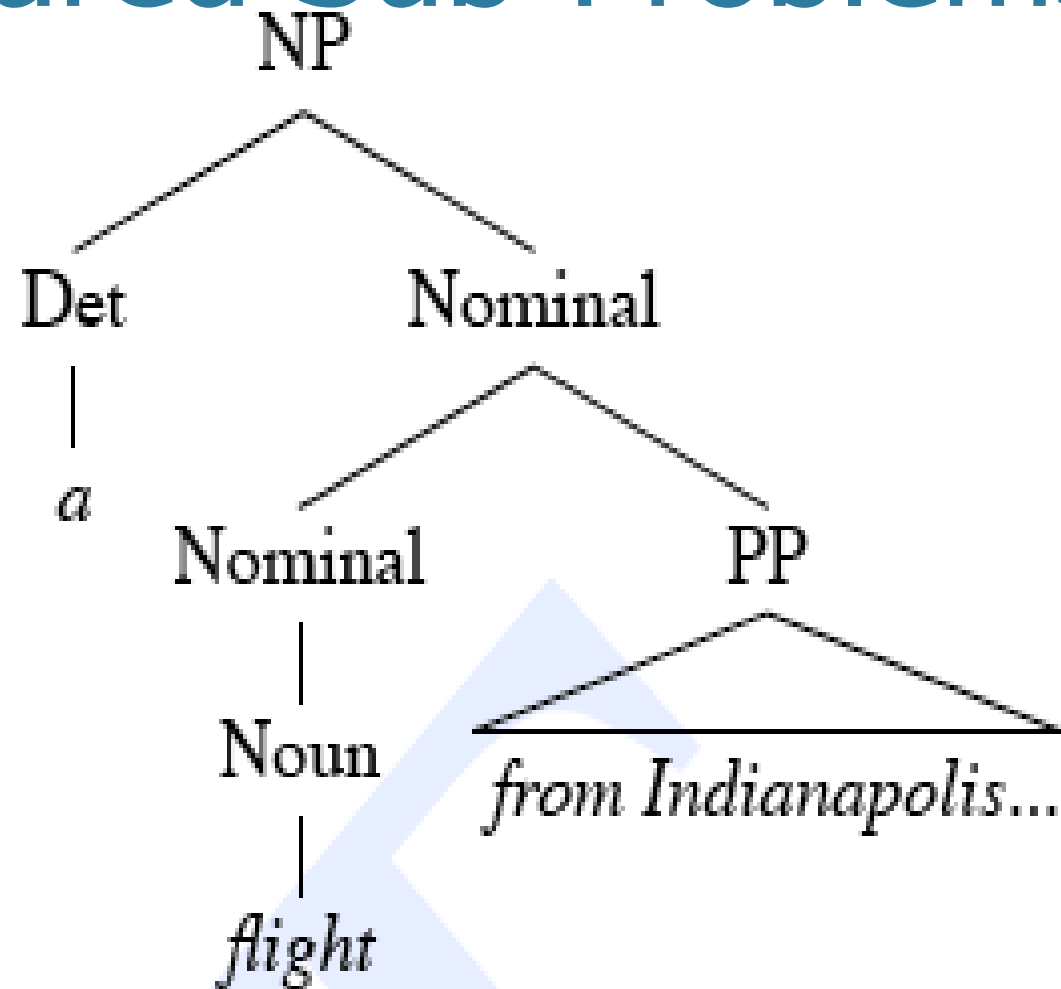
Repeated Work

- Top-down and bottom-up parsing both lead to repeated substructures
 - Globally bad parses can construct good subtrees
 - But overall parse will fail
 - Require reconstruction on other branch
 - No static backtracking strategy can avoid
- Efficient parsing techniques require storage of shared substructure
 - Typically with dynamic programming
- Example: *a flight from Indianapolis to Houston on TWA*

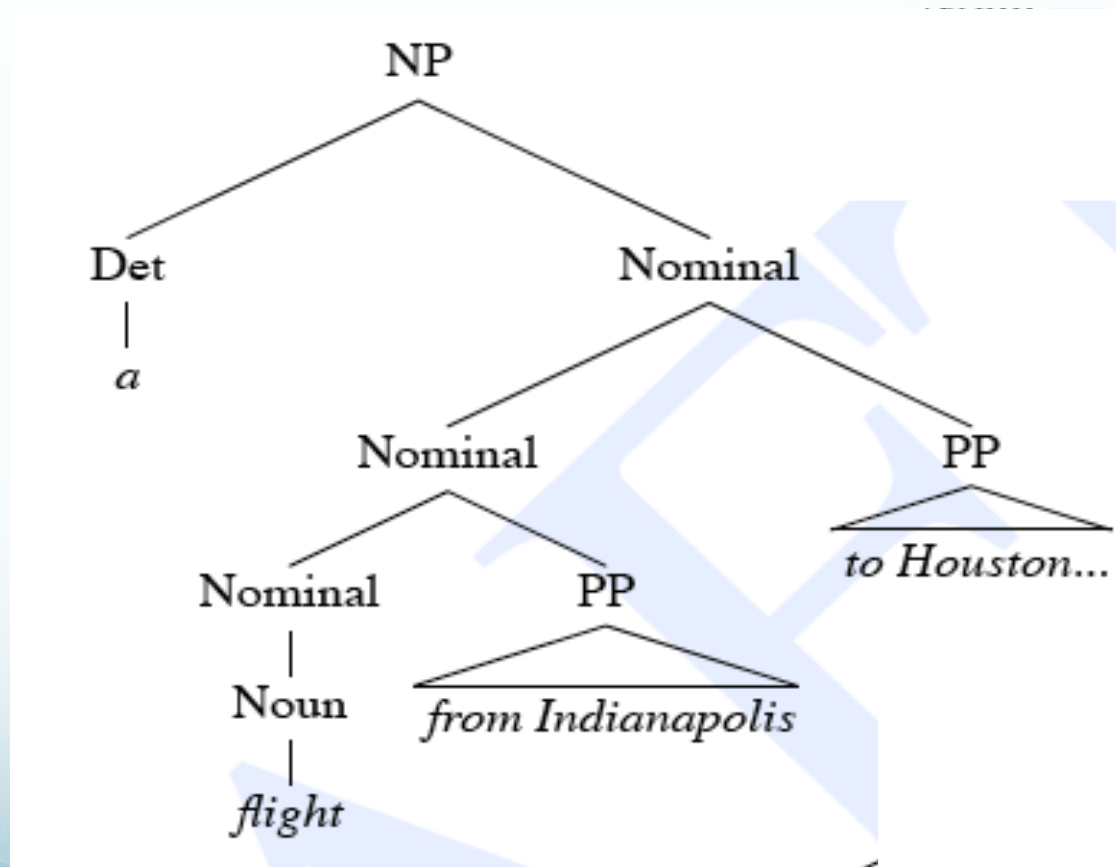
Shared Sub-Problems



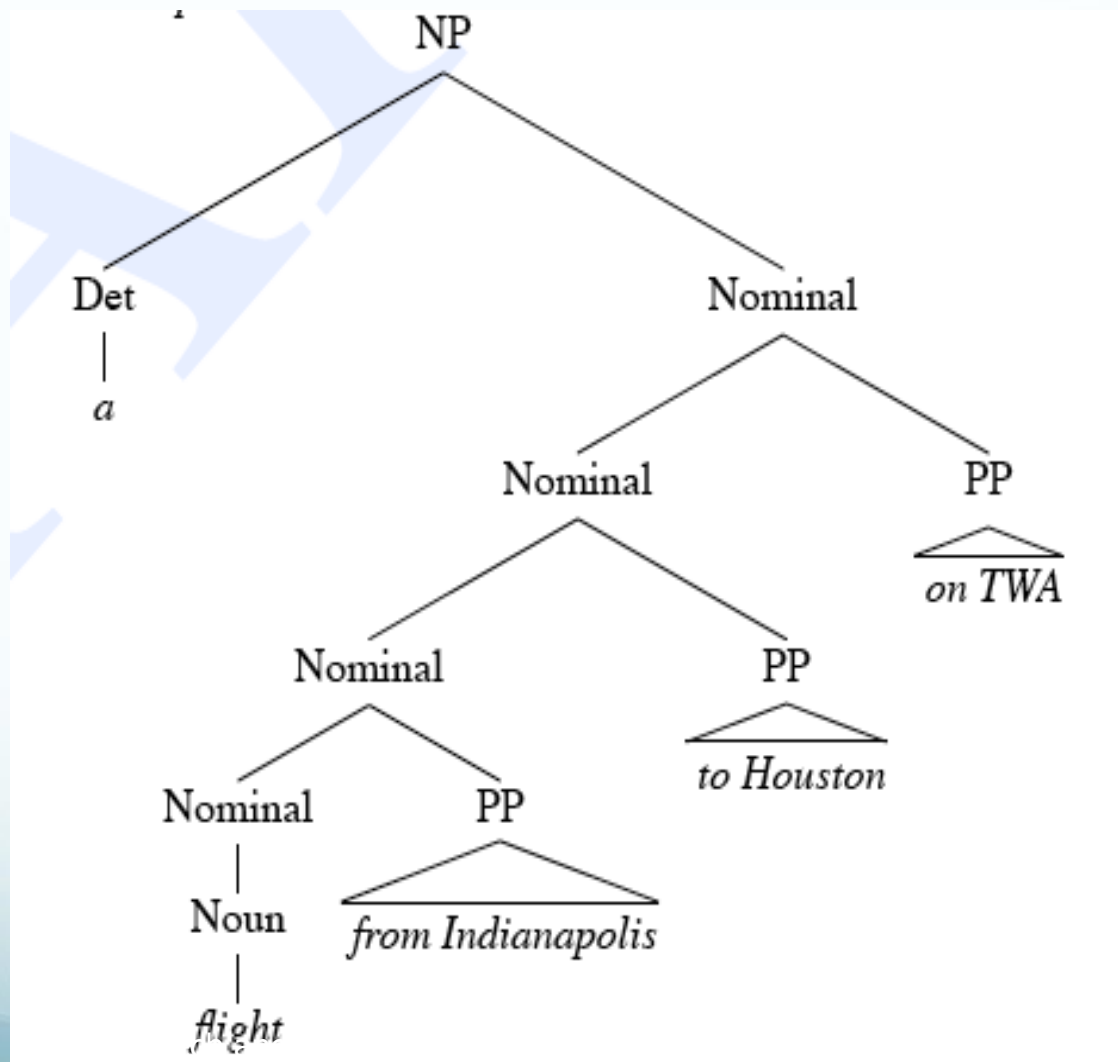
Shared Sub-Problems



Shared Sub-Problems



Shared Sub-Problems



Recursion

- Many grammars have recursive rules
 - E.g., $S \rightarrow S \text{ Conj } S$
- In search approaches, recursion is problematic
 - Can yield infinite searches
 - Esp., top-down

Dynamic Programming

- Challenge: Repeated substructure $\rightarrow\rightarrow$ Repeated work
- Insight:
 - Global parse composed of parse substructures
 - Can record parses of substructures
- Dynamic programming avoids repeated work by tabulating solutions to subproblems
 - Here, stores subtrees

Parsing w/Dynamic Programming

- Avoids repeated work
- Allows implementation of (relatively) efficient parsing algorithms
 - Polynomial time in input length
 - Typically cubic (n^3) or less
- Several different implementations
 - Cocke-Kasami-Younger (CKY) algorithm
 - Earley algorithm
 - Chart parsing

Chomsky Normal Form (CNF)

- CKY parsing requires grammars in CNF
- Chomsky Normal Form
 - All productions of the form:
 - $A \rightarrow BC$, or
 - $A \rightarrow a$
- However, most of our grammars are not of this form
 - E.g., $S \rightarrow \text{Wh-NP Aux NP VP}$
- Need a general conversion procedure
 - Any arbitrary grammar can be converted to CNF

Grammar Equivalence and Form

- Grammar equivalence
 - Weak: Accept the same language, May produce different analyses
 - Strong: Accept same language, Produce same structure

CNF Conversion

- Three main conditions:

- Hybrid rules:

- INF-VP $\rightarrow\rightarrow$ to VP

- Unit productions:

- $A \rightarrow\rightarrow B$

- Long productions:

- $A \rightarrow\rightarrow B C D$

CNF Conversion

- Hybrid rule conversion:
 - Replace all terminals with dummy non-terminals
 - E.g., INF-VP $\rightarrow\rightarrow$ to VP
 - INF-VP $\rightarrow\rightarrow$ TO VP; TO $\rightarrow\rightarrow$ to
- Unit productions:
 - Rewrite RHS with RHS of all derivable non-unit productions
 - If $A \Rightarrow^* B$ and $B \rightarrow\rightarrow w$, then add $A \rightarrow\rightarrow w$

CNF Conversion

- Long productions:
 - Introduce new non-terminals and spread over rules
 - $S \rightarrow Aux\ NP\ VP$
 - $S \rightarrow X1\ VP; X1 \rightarrow Aux\ NP$
- For all non-conforming rules,
 - Convert terminals to dummy non-terminals
 - Convert unit productions
 - Binarize all resulting rules

2.1 Grammar

S1i in CNF

$S \rightarrow IP VP$

$S \rightarrow Aux NP VP$

$S \rightarrow VP$

$NP \rightarrow Pronoun$

$i\backslash TP \rightarrow Proper-Noun$

$i\backslash TP \rightarrow Det Nominal$

$Nominal \rightarrow Noun$

$Nominal \rightarrow Nominal Noun$

$i\backslash Tominal \rightarrow i\backslash To m i n a l PP$

$VP \rightarrow Verb$

$VP \rightarrow Verb NP$

$VP \rightarrow Verb NP PP$

$VP \rightarrow Verb PP$

$VP \rightarrow VPPP$

$PP \rightarrow Preposition i\backslash TP$

$S \rightarrow i\backslash TP VP$

$S \rightarrow X1 VP$

$X1 \rightarrow Aux TP$

$S \rightarrow book \mid include \mid preje r$

$S \rightarrow Verb NP$

$S \rightarrow X2 PP$

$S \rightarrow Verb PP$

$S \rightarrow VPPP$

$NP \rightarrow I \mid she \mid rne$

$i\backslash TP \rightarrow TWA \mid Houston$

$i\backslash TP \rightarrow Det i\backslash Tominal$

$Nominal \rightarrow book \mid flight \mid meal \mid money$

$Nominal \rightarrow Nominal Noun$

$i\backslash Tominal \rightarrow Nominal PP$

$VP \rightarrow book \mid include \mid prefer$

$VP \rightarrow Verb NP$

$VP \rightarrow X2 PP$

$X2 \rightarrow Verb NP$

$VP \rightarrow Verb PP$

$VP \rightarrow VPPP$

$PP \rightarrow Preposition i\backslash TP$

CKY Parsing

- Cocke-Kasami-Younger parsing algorithm:
 - (Relatively) efficient bottom-up parsing algorithm based on tabulating substring parses to avoid repeated work
- Approach:
 - Use a CNF grammar
 - Build an $(n+1) \times (n+1)$ matrix to store subtrees
 - Upper triangular portion
 - Incrementally build parse spanning whole input string

Dynamic Programming in CKY

- Key idea:
 - For a parse spanning substring $[i,j]$, there exists some k such there are parses spanning $[i,k]$ and $[k,j]$
 - We can construct parses for whole sentence by building up from these stored partial parses
- So,
 - To have a rule $A \rightarrow B C$ in $[i,j]$,
 - We must have B in $[i,k]$ and C in $[k,j]$, for some $i < k < j$
 - CNF grammar forces this for all $j > i+1$