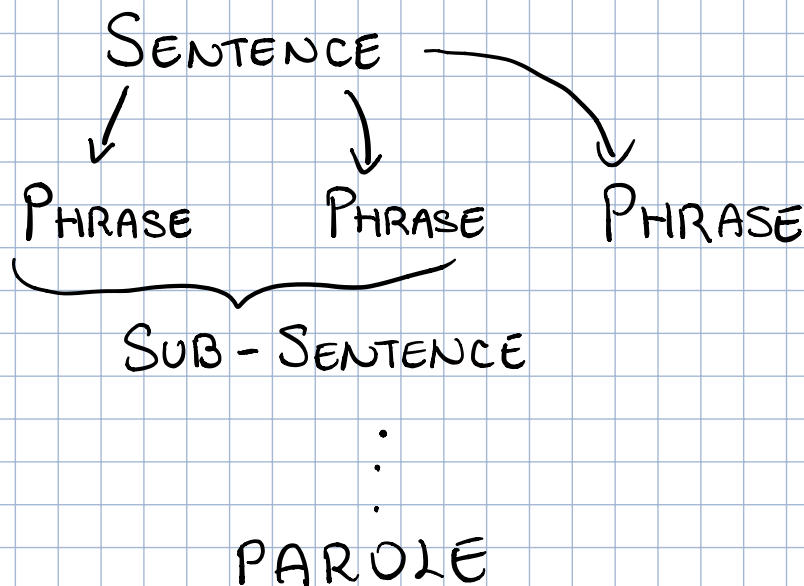


STRUTTURA GRAMMATICALE

L'oggetto minimale di una grammatica non le SENTENCES, le parti che vengono chiamate non le PAROLE, che vengono organizzate in PHRASES che a loro volta possono formare SUB-SENTENCES. Otteniamo quindi una struttura altamente RICORSIVA.



NOUN PHRASE

Governate dal nome. In italiano abbiamo le seguenti regole

NP → ART NOME

NP → ART AD3 NOME

NP → ART NOME AD3P

Ad esempio abbiamo la seguente cosa

Se libro utile allo studio
NP → ART NOME ADSP (FRASE AGGETTIVALE)

PREPOSITIONAL PHRASE

Una PP è una NP a cui ci mettiamo davanti
una preposizione.

PP → PREP NP

Una PP può anche essere dentro una NP
o una ADSP

ADSP → ADP PP

VERBAL PHRASE

VP → VX NP

VX → V

VP → VX NP PP

VX → Aux V

VP → VX NP PP PP

VP → VX PP

SENTENCE

Una tipica frase è data da

$$S \rightarrow NP VP$$

Notiamo a questo punto notiamo che la grammatica ottenuta non è in CNF, e quindi non possiamo utilizzare il CYK per PARSARE le varie sentences. Necesitiamo di un altro algoritmo di parsing che sia in grado di gestire grammatiche CF generali.

ANNOTAZIONI ~ (15:00 / 1 min)

Una ANNOTAZIONE è un ARCO tra due punti: un punto di INIZIO e un punto di FINE.

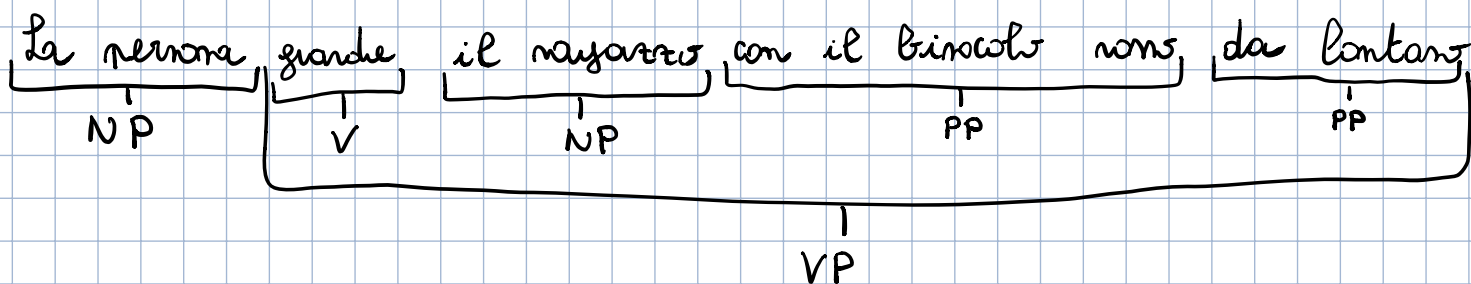
Annotare una frase ci permette di costruire un GRAFO i cui nodi sono gli SPAZI bianchi tra le parole e NON le parole stesse.

Ad una annotazione è associata una ETICHETTA, che può essere strutturata e prende il nome di FEATURE STRUCTURE.

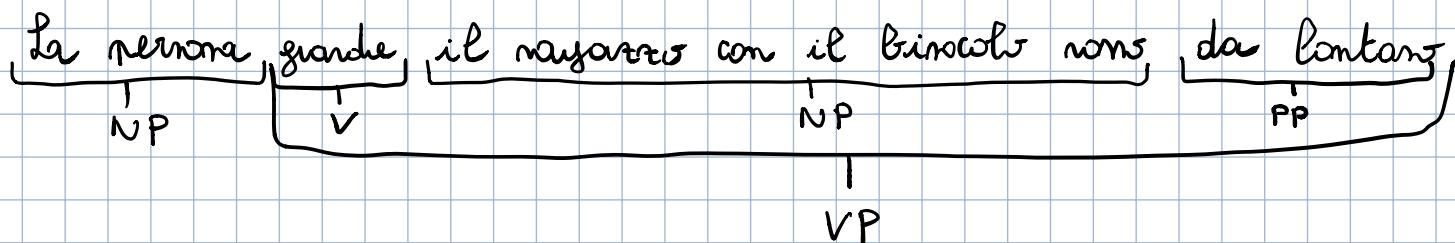
Consideriamo la frase

"La persona guarda il ragazzo con il binocolo non da lontano"

Poniamo annotare la frase in almeno i seguenti due modi:



In questa interpretazione è la persona che utilizza il binocolo non per guardare il ragazzo.



In quest'altro interpretazione è il ragazzo ad avere il binocolo non.

Un eventuale algoritmo di PARSING deve tenere in considerazione la potenziale AMBIGUITÀ della grammatica.

EARLEY ALGORITHM

The early algorithm parses the input string by constructing a CHART TABLE that is heavily used in a DYNAMIC PROGRAMMING way.

The early algorithm uses a CURSOR, denoted by a dot (\cdot), which is used to determine how much of a given rule in the grammar has been MATCHED in that particular step.

In particular, given a NT X and two strings of terminals/non-terminals, the notation

$$X \rightarrow \alpha \cdot \beta$$

is used to represent the fact that α has already been parsed and that β has yet to be parsed.

A combination of the following MAIN OPERATIONS are used by the EARLEY PARSER:

- PREDICT

Used to determine all NON-TERMINALS that will be expanded in the next operation.

- COMPLETE

Used to complete a scan operation, it moves the DOT (•) in all the production that have been MATCHED.

- SCAN

Used to check the input string in order to match grammar's rules.

The main aspects of the EARLEY PARSERS are thus the following :

- DYNAMIC PROGRAMMING (CHART MATRIX)
- ANNOTATIONS
- NOTATION ($X \rightarrow \alpha \cdot \beta$)

We will now see an example of the EARLEY parser. The example has associated AUDIO that explains how the chart was filled.

EARLEY PARSER EXAMPLE

(TAKEN FROM YOUTUBE VIDEO "CYK ALGORITHM BY DEEBA KANNAN")

GRAMMAR G


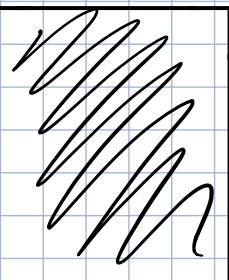



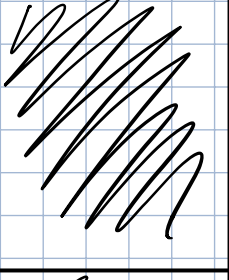

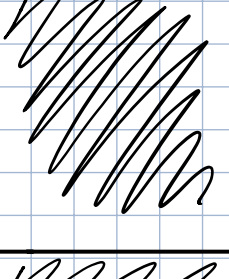

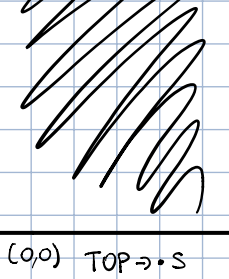
$S \rightarrow NP VP$ $P \rightarrow WITH$
 $PP \rightarrow P NP$ $V \rightarrow DRINK$
 $VP \rightarrow V NP$ $N \rightarrow PAPER$
 $VP \rightarrow VP PP$ $N \rightarrow MILK$
 $NP \rightarrow NP PP$ $N \rightarrow CHOCOLATE$
 $NP \rightarrow N$ $N \rightarrow COFFEE$

🎵: SPIEGAZIONE AUDIO!
ESEMPIO

INPUT STRING

PAPER DRINK MILK WITH CHOCOLATE

PARSER CHART

				(4,4) $NP \rightarrow \cdot NP PP$ $NP \rightarrow \cdot N$ $N \rightarrow \cdot CHOCOLATE$	(4,5) $N \rightarrow CHOCOLATE \cdot$ $NP \rightarrow N \cdot$ $NP \rightarrow NP \cdot PP$
			(3,3) $PP \rightarrow \cdot P NP$ $P \rightarrow \cdot WITH$	(3,4) $P \rightarrow WITH \cdot$ $PP \rightarrow P \cdot NP$	(3,5) $PP \rightarrow P NP \cdot$
		(2,2) $NP \rightarrow \cdot NP PP$ $NP \rightarrow \cdot N$ $N \rightarrow \cdot MILK$	(2,3) $N \rightarrow MILK \cdot$ $NP \rightarrow N \cdot$ $NP \rightarrow NP \cdot PP$		(2,5) $NP \rightarrow NP PP \cdot$
	(1,1) $VP \rightarrow \cdot V NP$ $VP \rightarrow \cdot VP PP$ $PP \rightarrow \cdot P NP$ $V \rightarrow \cdot DRINK$	(1,2) $V \rightarrow DRINK \cdot$ $VP \rightarrow V \cdot NP$	(1,3) $VP \rightarrow V NP \cdot$ $P \rightarrow VP \cdot PP$		(1,5) $VP \rightarrow VP PP \cdot$
(0,0) TOP $\rightarrow \cdot S$ $S \rightarrow \cdot NP VP$ $NP \rightarrow \cdot NP PP$ $NP \rightarrow \cdot N$ $N \rightarrow \cdot PAPER$	(0,1) $N \rightarrow PAPER \cdot$ $NP \rightarrow N \cdot$ $NP \rightarrow NP \cdot PP$ $S \rightarrow NP \cdot VP$		(0,3) $S \rightarrow NP VP \cdot$		(0,5) $S \rightarrow NP VP \cdot$

The obtained PARSE TREES are as follows

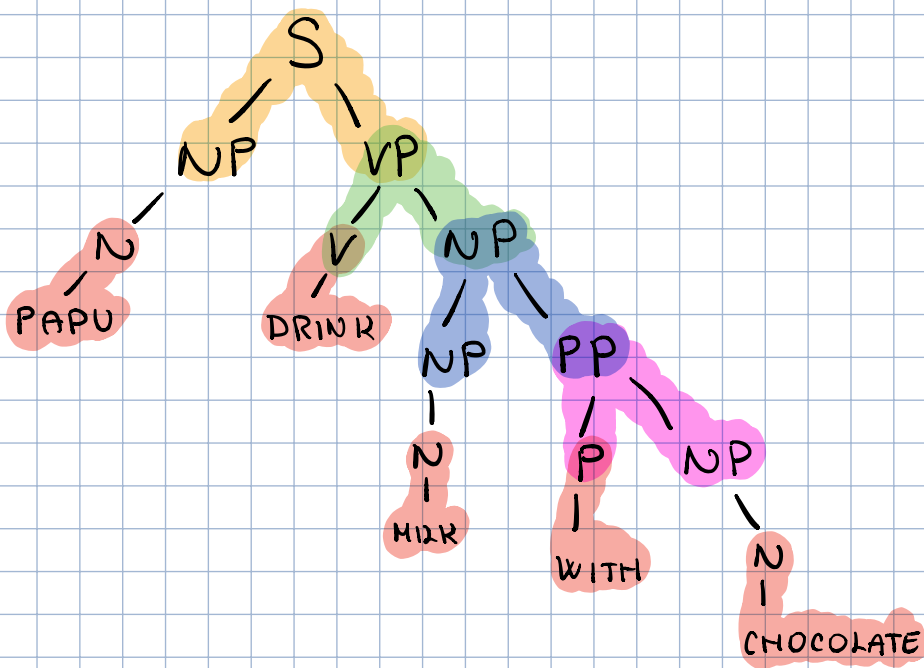
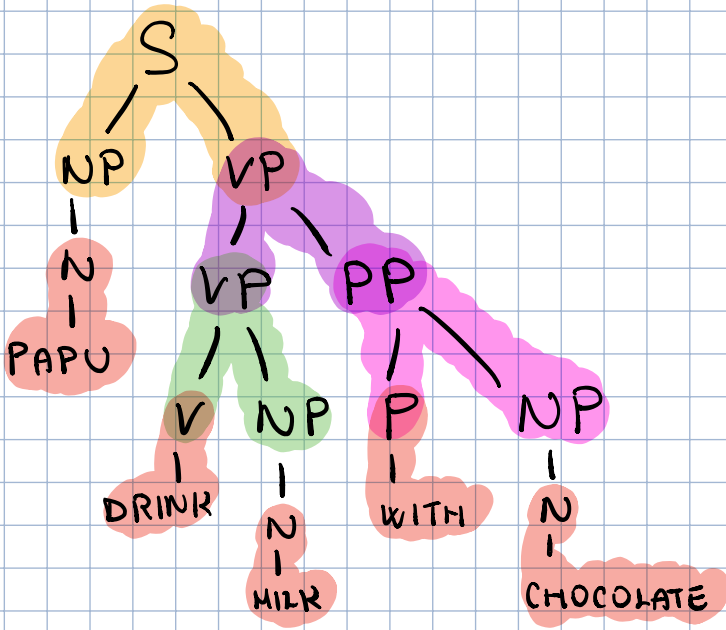


CHART PARSER

~ (03:00/2 min)

Sie il CYK che è EARLEY PARSER si basano sui seguenti concetti fondamentali

- ANNOTAZIONE
- ETICHETTA CON REGOLA
- REGOLA ATTIVA

Per motivi di efficienza è meglio utilizzare in algoritmi BOTTOM-UP. In particolare esiste il BOTTOM-UP CHART PARSING.

A differenza del CYK, nell'EARLEY PARSER durante l'operazione di PREDICT, nella chart table vengono inserite un sacco di REGOLE.

In particolare non ci sono modi per gestire questo, e quindi il fattore $O(n^3)$ nella SPACE COMPLEXITY si fa molto sentire, e rende l'EARLEY PARSER quasi inutilizzabile per grammatiche molto grandi.

TRATTARE LE CONGIUNZIONI COORDINATIVE

Consideriamo la frase

MARIO MANGIA PANE E MARIA MANGIA SALAME

notiamo che la frase può anche essere scritta
come segue

MARIO MANGIA PANE E MARIA SALAME

Come possiamo rappresentare le PARTI SOTTOINTESE
del discorso?

Volendo utilizzare una rappresentazione a COSTITUENTI,
possiamo iniziare utilizzando la seguente regola

$$S \rightarrow S \text{ CC } S$$

Questa regola però ci permette di gestire frasi
semplici senza parti sottintese. In generale
trattare le parti mancanti è un problema molto
difficile.

MARIO MANGIA PANE E BEVE IL VINO

In questa frase
manca il SOGGETTO

Se non trattiamo questo problema al livello della SINTASSI nel successivo livello, quello della SEMANTICA, incapperemo in varie problematiche nel dare la giusta interpretazione alle strutture sintattiche costruite.

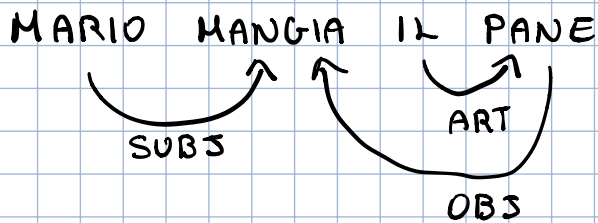
Per trattare queste problematiche poniamo adattare un modulo in grado di inserire dei PLACE-HOLDERS all'interno delle frasi da analizzare che vanno ad ESPlicitARE le parti sottintese delle frasi.

DEPENDENCIES ~ (35:00/2 min)

Oltre alle rappresentazione basate sui SINTAGMI che genera nei i PARSE TREE abbiamo un altro metodo per rappresentare le informazioni sintattiche, che prende il nome di RAPPRESENTAZIONE ALLE DIPENDENZE.

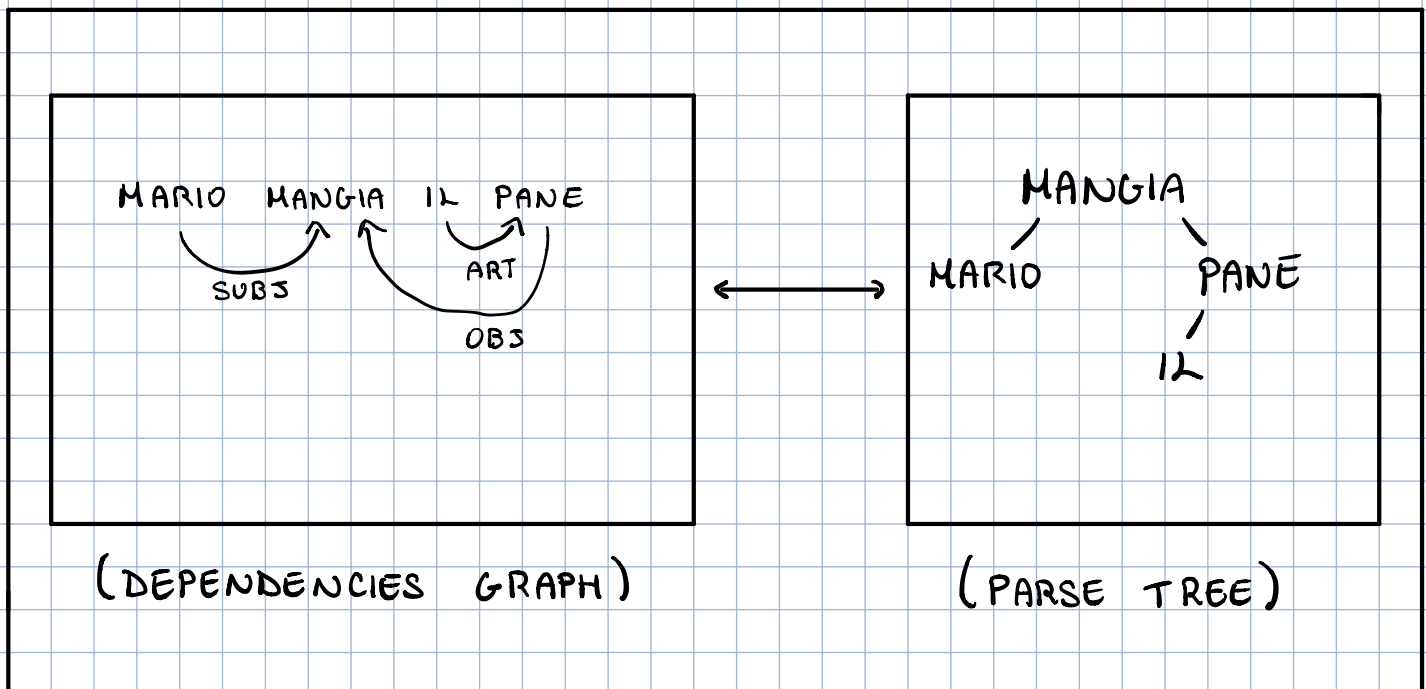
Questa nuova rappresentazione non introduce degli elementi intermedi ma lavora direttamente con le parole, descrivendo come queste ultime sono collegate tra loro tramite un GRAFO ETICHETTATO.

Date una frase come "MARIO MANGIA IL PANE",
abbiamo la seguente rappresentazione



Notiamo che in alcuni casi è possibile
passare da una rappresentazione all'altra.
Per fare questo dobbiamo capire qual è
l'elemento nella rapp. delle dipendenze
che vogliamo far uscire di livello. L'elemento
che facciamo uscire è detto HEAD.

Nel caso di prima abbiamo



Motivato però che con la rapp. alle dipendente possiamo anche far INCROCIARE gli archi. In questo caso i grafi ottenuti vengono detti NON-PROJECTIVES, e non vale più l'equivalenza tra le rappresentazioni.

OSS: Tutti gli alberi presenti nelle UNIVERSAL DEPENDENCIES NON PROJECTIVE.

Per adesso è unico modo per ottenere la rapp. alle dipendente che conosciamo è quello di utilizzare il CYK o l'EARLEY per ottenere il PARSE TREE e poi da questo ottenere il DEPENDENCY TREE/GRAPH. Esistono poi degli algoritmi che costruiscono direttamente la rapp. alle dipendente.

OSS: Inizialmente la rapp. alle dipendente non è stata presa molto in considerazione. Successivamente con la nascita del MACHINE LEARNING (ML) è stata rivista e ripresa in considerazione.