

Natural language syntax: parsing and complexity

Timothée Bernard and Pascal Amsili

Université Paris Cité, Université Sorbonne Nouvelle
timothee.bernard@u-paris.fr, pascal.amsili@ens.fr

Ljubljana, Slovenia – August 7-11, 2023
ESLLI foundational course in Language and Computation

Derivation graph of a grammar

Given a grammar, one can build the directed graph such that

- nodes are words of $(\Sigma \cup N)^*$,
- edges correspond to rewriting.

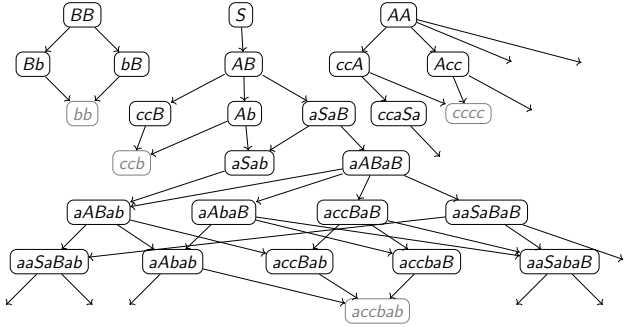
Derivation graph of a grammar

Given a grammar, one can build the directed graph such that

- nodes are words of $(\Sigma \cup N)^*$,
- edges correspond to rewriting.

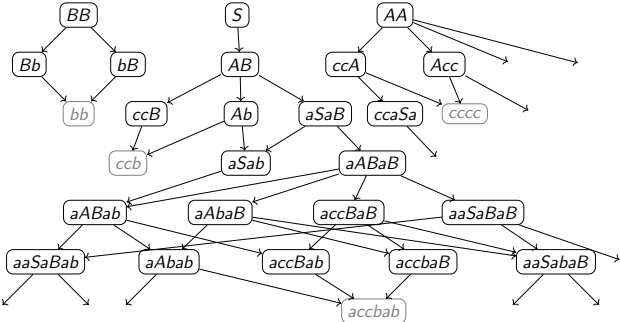
Example:

- $S \rightarrow AB$
- $A \rightarrow cc$
 - $\mid aSa$
- $B \rightarrow b$



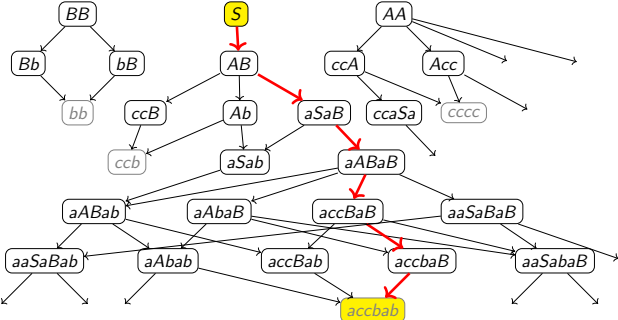
Parsing is a search in the derivation graph

Parsing $w \in \Sigma^*$:



Parsing is a search in the derivation graph

Parsing $w \in \Sigma^*$: finding a path in the graph going from S to w .

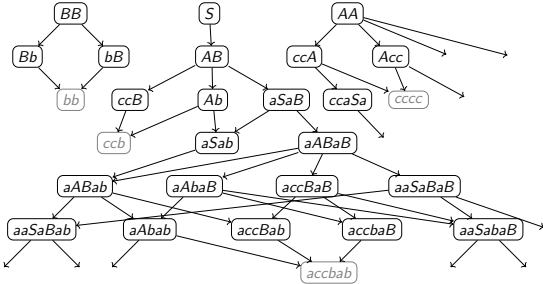


Left derivation

- **Left derivation:** always rewrite the leftmost non-terminal symbol. Examples:
 - $S \Rightarrow AB \Rightarrow ccB \Rightarrow ccb$
 - $S \Rightarrow AB \Rightarrow Ab \Rightarrow ccb$
- In a CFG, every syntactic structure is associated with a single left derivation, and vice versa.

Parsing can focus on left derivations

- Without loss of generality, parsing can be defined as a search for left derivation(s) [or right].

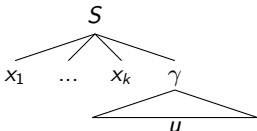


Pruning the graph with the prefix property

Suppose that $S \xRightarrow{*} x_1 \dots x_k \gamma$ with $x_1 \dots x_k \in \Sigma^*$, and $\gamma \in (\Sigma \cup N)^*$.

Any word w s.t. $S \xRightarrow{*} x_1 \dots x_k \gamma \xRightarrow{*} w$ has $x_1 \dots x_k$ as a prefix:

$$\exists u \in \Sigma^* \text{ s.t. } w = x_1 \dots x_k u.$$

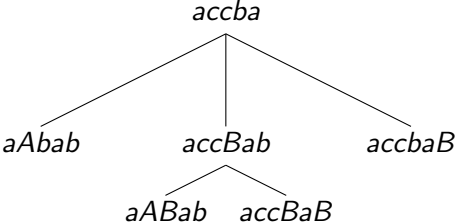


A top-down derivation can be stopped as soon as it contains a non-empty prefix of letters that does not match the query.

Bottom-up approaches: example (I)

At each step:
 “Which parts of w are identical to the right-hand side of a rule?”

$S \rightarrow AB$
 $A \rightarrow Saa$
 | cc
 $B \rightarrow b$



A parsing is a a sequence of *reductions*
 (“rewriting inversions”)

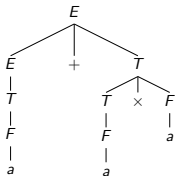
Example

$$E \rightarrow E + T$$

$$E \rightarrow T$$

$$T \rightarrow T \times F$$

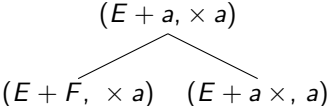
$$T \rightarrow F$$

$$F \rightarrow a$$


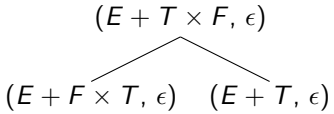
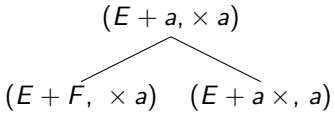
stack	buffer	action
ϵ	$a + a \times a$	shift
a	$+ a \times a$	reduce ($F \rightarrow a$)
F	$+ a \times a$	reduce ($T \rightarrow F$)
T	$+ a \times a$	reduce ($E \rightarrow T$)
E	$+ a \times a$	shift
$E +$	$a \times a$	shift
$E + a$	$\times a$	reduce ($F \rightarrow a$)
$E + F$	$\times a$	reduce ($T \rightarrow F$)
$E + T$	$\times a$	shift
$E + T \times$	a	shift
$E + T \times a$	ϵ	reduce ($F \rightarrow a$)
$E + T \times F$	ϵ	reduce ($T \rightarrow T \times F$)
$E + T$	ϵ	reduce ($E \rightarrow E + T$)
E	ϵ	accept

$$E \rightarrow E + T \rightarrow E + T \times F \rightarrow E + T \times a \rightarrow E + F \times a \rightarrow E + a \times a \rightarrow T + a \times a \rightarrow a + a \times a$$

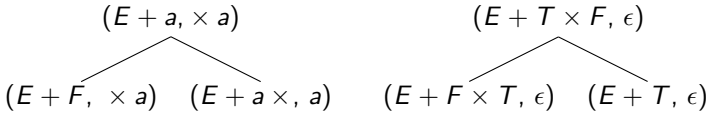
Sources of non-determinism



Sources of non-determinism



Sources of non-determinism



- Choice of the suffix of the stack to reduce.
- Choice of the rule to reduce with.
- Choice between shift and reduce.

A possibly efficient parsing algorithm

A **deterministic shift-reduce parser** can determine at each step and in constant time, based on the content of the stack and the content of the buffer, which action to perform.

- For some CFGs, this is possible.
 - LR(k)* grammars, for the main parts of many programming languages.
- For most CFGs, this is impossible.

Natural language syntax has lots of ambiguities

PP attachment, modifier scope, etc.

- (1) a. Bob saw a passer-by with his telescope.
- b. Wild cats and dogs chase rats.
- c. The men and women from Tirol...

→ Deterministic parsing is not available for natural languages.

Natural language syntax has lots of ambiguities

PP attachment, modifier scope, etc.

- (1) a. Bob saw [a passer-by [with his telescope]].
- b. Wild cats and dogs chase rats.
- c. The men and women from Tirol...

→ Deterministic parsing is not available for natural languages.

Natural language syntax has lots of ambiguities

PP attachment, modifier scope, etc.

- (1) a. Bob saw [a passer-by] [with his telescope].
- b. Wild cats and dogs chase rats.
- c. The men and women from Tirol...

→ Deterministic parsing is not available for natural languages.

Natural language syntax has lots of ambiguities

PP attachment, modifier scope, etc.

- (1)
 - a. Bob saw [a passer-by] [with his telescope].
 - b. [[Wild cats] and dogs] chase rats.
 - c. The men and women from Tirol...

→ Deterministic parsing is not available for natural languages.

Natural language syntax has lots of ambiguities

PP attachment, modifier scope, etc.

- (1)
- a. Bob saw [a passer-by] [with his telescope].
 - b. [Wild [cats and dogs]] chase rats.
 - c. The men and women from Tirol...

→ Deterministic parsing is not available for natural languages.

Natural language syntax has lots of ambiguities

PP attachment, modifier scope, etc.

- (1) a. Bob saw [a passer-by] [with his telescope].
- b. [Wild [cats and dogs]] chase rats.
- c. [[The [men and women]] from Tirol]...

→ Deterministic parsing is not available for natural languages.

Natural language syntax has lots of ambiguities

PP attachment, modifier scope, etc.

- (1)
- a. Bob saw [a passer-by] [with his telescope].
 - b. [Wild [cats and dogs]] chase rats.
 - c. [The men] and [women from Tirol]...

→ Deterministic parsing is not available for natural languages.

Natural language syntax has lots of ambiguities

PP attachment, modifier scope, etc.

- (1) a. Bob saw [a passer-by] [with his telescope].
- b. [Wild [cats and dogs]] chase rats.
- c. [[The men] and [women] from Tirol]...

→ Deterministic parsing is not available for natural languages.

Chart parsing is a possible answer to ambiguity

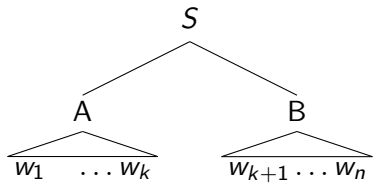
- Idea: decompose the analysis of a word w into the independent analyses of the spans of w .
- The result of these subanalyses are then combined to provide analyses of the whole w .
- Based on a data structure that stores all partial possible parses:
 - computations are done only once, → dynamic programming
 - multiple analyses can be handled.

Two well-known chart parsing algorithms:

- CYK: a bottom-up algorithm that works with grammars in Chomsky Normal Form.
- Earley: a mostly top-down algorithm that works with any CFG.

Factoring the computation

- Given a CFG G and a query w ...
- Suppose $S \overset{*}{\Rightarrow} AB$.
- To answer the question whether $AB \overset{*}{\Rightarrow} w$,
- we may look for a $k \in [1, n]$ (with $n = |w|$) such that:
- $A \overset{*}{\Rightarrow} w_{1:k}$ and $B \overset{*}{\Rightarrow} w_{k+1:n}$.



Constituent chart: cells correspond to spans

5					
4					
3					
2					
1					
j/i	1	2	3	4	5
$w =$	My	sister	likes	Sam	

Convention: $T[i, j] \sim \text{span } w_i \dots w_{j-1}$.

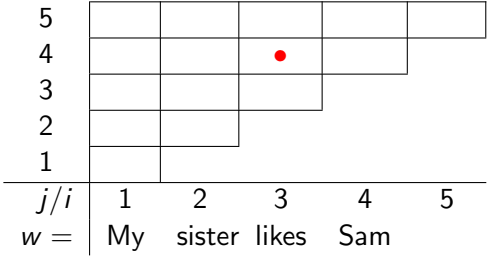
Constituent chart: cells correspond to spans

5					
4					
3					
2		●			
1					
j/i	1	2	3	4	5
$w =$	My	sister	likes	Sam	

Convention: $T[i, j] \sim \text{span } w_i \dots w_{j-1}$.

$T[2, 2] \rightsquigarrow \epsilon$

Constituent chart: cells correspond to spans



Convention: $T[i, j] \sim \text{span } w_i \dots w_{j-1}$.

$T[3, 4] \rightsquigarrow$ likes

Constituent chart: cells correspond to spans

5			●		
4					
3					
2					
1					
j/i	1	2	3	4	5
$w =$	My	sister	likes	Sam	

Convention: $T[i, j] \sim \text{span } w_i \dots w_{j-1}$.

$T[3, 5] \rightsquigarrow \text{likes Sam}$

Constituent chart: cells correspond to spans

5	●				
4					
3					
2					
1					
j/i	1	2	3	4	5
$w =$	My	sister	likes	Sam	

Convention: $T[i, j] \sim \text{span } w_i \dots w_{j-1}$.

$T[1, 5] \rightsquigarrow w$

Constituent chart: information stored in the cells

- CYK: Non-terminal symbols are stored in the chart.
- Earley: ...

5					
4					
3					
2					
1					
j/i	1	2	3	4	5
$w =$	My	sister	likes	Sam	

Convention: $T[i, j] = \{A \in N \mid A \xrightarrow{*} w_{i:j-1}\}$

Constituent chart: information stored in the cells

- CYK: Non-terminal symbols are stored in the chart.
- Earley: ...

5						
4						
3		N				
2						
1						
	j/i	1	2	3	4	5
	$w =$	My	sister	likes	Sam	

Convention: $T[i, j] = \{A \in N \mid A \xrightarrow{*} w_{i:j-1}\}$

$N \rightarrow \text{sister} \in G$

Constituent chart: information stored in the cells

- CYK: Non-terminal symbols are stored in the chart.
- Earley: ...

5					
4					
3		<i>N</i>			
2	<i>Det</i>				
1					
<i>j/i</i>	1	2	3	4	5
<i>w =</i>	My	sister	likes	Sam	

Convention: $T[i, j] = \{A \in N \mid A \xrightarrow{*} w_{i:j-1}\}$

N → sister ∈ *G*

Det → my ∈ *G*

Constituent chart: information stored in the cells

- CYK: Non-terminal symbols are stored in the chart.
- Earley: ...

5					
4					
3	<i>NP</i>	<i>N</i>			
2	<i>Det</i>				
1					
<i>j/i</i>	1	2	3	4	5
<i>w =</i>	My	sister	likes	Sam	

Convention: $T[i, j] = \{A \in N \mid A \xrightarrow{*} w_{i:j-1}\}$

- $N \rightarrow \text{sister} \in G$
- $Det \rightarrow \text{my} \in G$
- $NP \rightarrow Det N \in G$

Constituent chart: a complete example

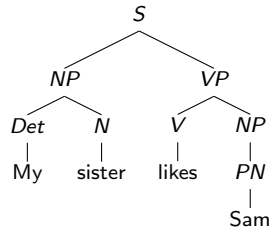
- $S \rightarrow NP VP$
- $NP \rightarrow Det N$
- $NP \rightarrow PN$
- $VP \rightarrow V NP$
- $VP \rightarrow V$
- $Det \rightarrow my \mid the$
- $N \rightarrow sister \mid moon$
- $V \rightarrow likes \mid knows$
- $PN \rightarrow Sam \mid Joan$

5	S	\emptyset	VP	$\{PN, NP\}$	\emptyset
4	S	\emptyset	$\{V, VP\}$	\emptyset	
3	NP	N	\emptyset		
2	Det	\emptyset			
1	\emptyset				
j/i	1	2	3	4	5
$w =$	My	sister	likes	Sam	

- There can be multiple non-terminal symbols in a cell.
- $w \in \Sigma^* \in \mathcal{L}(G)$ iff $S \in T[1, |w| + 1]$

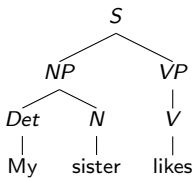
Decoding of the constituent chart

5	<i>S</i>	∅	<i>VP</i>	{ <i>PN, NP</i> }	∅
4	<i>S</i>	∅	{ <i>V, VP</i> }	∅	
3	<i>NP</i>	<i>N</i>	∅		
2	<i>Det</i>	∅			
1	∅				
<i>j/i</i>	1	2	3	4	5
<i>w =</i>	My	sister	likes	Sam	



Decoding of the constituent chart

5	<i>S</i>	∅	<i>VP</i>	{ <i>PN, NP</i> }	∅
4	<i>S</i>	∅	{ <i>V, VP</i> }	∅	
3	<i>NP</i>	<i>N</i>	∅		
2	<i>Det</i>	∅			
1	∅				
<i>j/i</i>	1	2	3	4	5
<i>w =</i>	My	sister	likes	Sam	



Chomsky Normal Form

A grammar is said to be in Chomsky Normal Form (CNF) iff all its rules are of the following form:

- $A \rightarrow BC$, with $A, B, C \in N$,
- or $A \rightarrow a$, with $a \in \Sigma$ and $A \in N$.

- Any CF grammar is weakly equivalent to a CNF grammar.
- Grammars in CNF have no ϵ -rule (apart maybe from a non recursive axiom).
- Grammars in CNF have no non productive cycle.

Filling the chart with a CNF grammar (I)

5					
4					
3					
2					
1					
<i>j/i</i>	1	2	3	4	5
<i>w =</i>	My	sister	likes	Sam	

- S* → *NP VP*
- NP* → *Det N*
- VP* → *V PN*
- Det* → my
- N* → sister
- V* → likes
- PN* → Sam

Filling the chart with a CNF grammar (I)

5					\emptyset
4				\emptyset	
3			\emptyset		
2		\emptyset			
1	\emptyset				
j/i	1	2	3	4	5
$w =$	My	sister	likes	Sam	

- $S \rightarrow NP VP$
- $NP \rightarrow Det N$
- $VP \rightarrow V PN$
- $Det \rightarrow my$
- $N \rightarrow sister$
- $V \rightarrow likes$
- $PN \rightarrow Sam$

- Spans of length 0 (diagonal) are never generated in a CNF:
 $\forall i \in [1, n], T[i, i] = \emptyset.$

Filling the chart with a CNF grammar (I)

5				<i>PN</i>	\emptyset
4			<i>V</i>	\emptyset	
3		<i>N</i>	\emptyset		
2	<i>Det</i>	\emptyset			
1	\emptyset				
<i>j/i</i>	1	2	3	4	5
<i>w =</i>	My	sister	likes	Sam	

- S* → *NP VP*
- NP* → *Det N*
- VP* → *V PN*
- Det* → my
- N* → sister
- V* → likes
- PN* → Sam

- Spans of length 0 (diagonal) are never generated in a CNF:
 $\forall i \in [1, n], T[i, i] = \emptyset.$
- Spans of length 1 are generated by lexical rules:
 $\forall i \in [1, n], T[i, i + 1] = \{A \in N \mid A \rightarrow w_i \in P\}.$

Filling the chart with a CNF grammar (II)

5				<i>PN</i>	\emptyset
4			<i>V</i>	\emptyset	
3		<i>N</i>	\emptyset		
2	<i>Det</i>	\emptyset			
1	\emptyset				
<i>j/i</i>	1	2	3	4	5
<i>w =</i>	My	sister	likes	Sam	

- S* → *NP VP*
- NP* → *Det N*
- VP* → *V PN*
- Det* → my
- N* → sister
- V* → likes
- PN* → Sam

- Spans of length ≥ 2 :

$$T[i, j] = \{A \in N \mid \exists k \in [i + 1, j - 1],$$

$$\exists B \in T[i, k],$$

$$\exists C \in T[k, j],$$

$$A \rightarrow BC \in P\}$$

Filling the chart with a CNF grammar (II)

5				<i>PN</i>	\emptyset
4			<i>V</i>	\emptyset	
3	<i>NP</i>	<i>N</i>	\emptyset		
2	<i>Det</i>	\emptyset			
1	\emptyset				
j/i	1	2	3	4	5
$w =$	My	sister	likes	Sam	

- S* → *NP VP*
- NP* → *Det N*
- VP* → *V PN*
- Det* → my
- N* → sister
- V* → likes
- PN* → Sam

- Spans of length ≥ 2 :

$$\begin{aligned}
 T[i, j] = & \{ A \in N \mid \exists k \in [i + 1, j - 1], \\
 & \exists B \in T[i, k], \\
 & \exists C \in T[k, j], \\
 & A \rightarrow BC \in P \}
 \end{aligned}$$

Filling the chart with a CNF grammar (II)

5			<i>VP</i>	<i>PN</i>	\emptyset
4			<i>V</i>	\emptyset	
3	<i>NP</i>	<i>N</i>	\emptyset		
2	<i>Det</i>	\emptyset			
1	\emptyset				
<i>j/i</i>	1	2	3	4	5
<i>w =</i>	My	sister	likes	Sam	

- S* → *NP VP*
- NP* → *Det N*
- VP* → *V PN*
- Det* → my
- N* → sister
- V* → likes
- PN* → Sam

- Spans of length ≥ 2 :

$$T[i, j] = \{ A \in N \mid \exists k \in [i + 1, j - 1],$$

$$\exists B \in T[i, k],$$

$$\exists C \in T[k, j],$$

$$A \rightarrow BC \in P \}$$

Filling the chart with a CNF grammar (II)

5	<i>S</i>		<i>VP</i>	<i>PN</i>	\emptyset
4			<i>V</i>	\emptyset	
3	<i>NP</i>	<i>N</i>	\emptyset		
2	<i>Det</i>	\emptyset			
1	\emptyset				
<i>j/i</i>	1	2	3	4	5
<i>w =</i>	My	sister	likes	Sam	

- S* → *NP VP*
- NP* → *Det N*
- VP* → *V PN*
- Det* → my
- N* → sister
- V* → likes
- PN* → Sam

- Spans of length ≥ 2 :

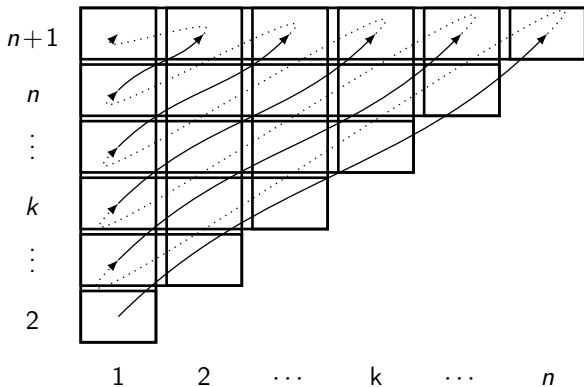
$$T[i, j] = \{A \in N \mid \exists k \in [i + 1, j - 1],$$

$$\exists B \in T[i, k],$$

$$\exists C \in T[k, j],$$

$$A \rightarrow BC \in P\}$$

Diagonal strategy



CYK: Algorithm

```

// Input:   $u \in \Sigma^*$ 
// Output: the constituent chart of  $u$ 
Function CYK-diagonal( $u$ )
┌    $T :=$  empty chart( $u$ );
  │   // First diagonal (unary cases)
  │   for  $i := 1$  to  $|u|$  do
  │   │   foreach  $(A \rightarrow u_i) \in P$  do
  │   │   │    $T[i, i + 1].add(A)$ ;
  │   │
  │   // Other diagonals (binary cases)
  │   for  $l := 2$  to  $|u|$  do           // loop on the length of the span
  │   │   for  $i := 1$  to  $|u| + 1 - l$  do // loop on the beginning of the span
  │   │   │    $j := i + l$ ;                // end of the span
  │   │   │   for  $k := i + 1$  to  $j - 1$  do // loop on the splitting point
  │   │   │   │   foreach  $(A \rightarrow BC) \in P$  do
  │   │   │   │   │   if  $B \in T[i, k]$  and  $C \in T[k, j]$  then
  │   │   │   │   │   │    $T[i, j].add(A)$ ;
  │   │   │   │
  │   │   │   return  $T$ ;
└
```

CYK Summary

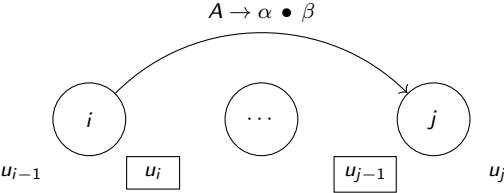
- Time complexity: $O(n^3)$
- Additional information can be stored for decoding the chart into trees.
- Efficient algorithm but requires transformation into CNF.
- Can be adapted for CCG, TAG, probabilistic CFG...

Earley Algorithm

- Works with any CFG (no transformation required).
- For CYK, it was possible to store non-terminals in the chart
($A \in T[i, j]$ iff $A \xrightarrow{*} w_{i:j-1}$).
- For Earley parsing, the chart will contain **dotted rules**:
($A \rightarrow \alpha \bullet \beta$) $\in T[i, j]$ iff $\alpha \xrightarrow{*} w_{i:j-1}$.
- Successful analysis: $\exists \alpha, (S \rightarrow \alpha \bullet) \in T[1, |w| + 1]$.

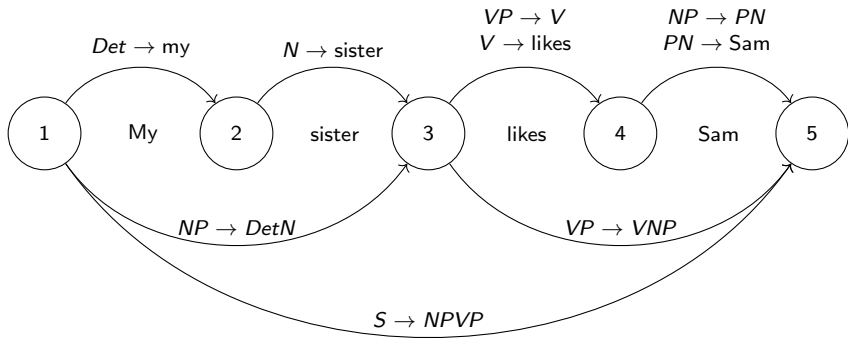
Earley items

- The information that $(A \rightarrow \alpha \bullet \beta) \in T[i, j]$
 - is an **(Earley) item**
 - and is written “ $(A \rightarrow \alpha \bullet \beta, i, j)$ ”.
- Graphical representation:



- Interpretation:
 - one is trying to recognise A starting from u_i ;
 - so far, one has recognised α up to u_{j-1} (included).

Another view on the chart



5	S	∅	VP	{PN, NP}	∅
4	S	∅	{V, VP}	∅	
3	NP	N	∅		
2	Det	∅			
1	∅				
j/i	1	2	3	4	5
w =	My	sister	likes	Sam	

Use of dotted rules

- The use of dotted rules makes it relatively easy to generalise the idea behind CYK to non-binary rules.

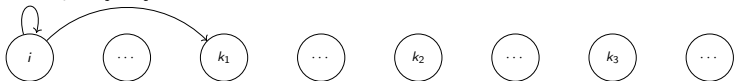
Use of dotted rules

- The use of dotted rules makes it relatively easy to generalise the idea behind CYK to non-binary rules.
- Example:

Use of dotted rules

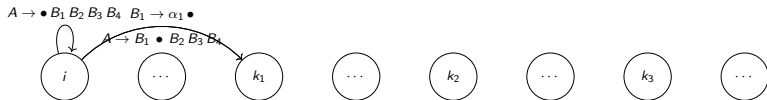
- The use of dotted rules makes it relatively easy to generalise the idea behind CYK to non-binary rules.
- Example:

$A \rightarrow \bullet B_1 B_2 B_3 B_4 \quad B_1 \rightarrow \alpha_1 \bullet$



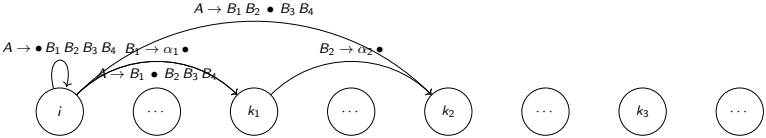
Use of dotted rules

- The use of dotted rules makes it relatively easy to generalise the idea behind CYK to non-binary rules.
- Example:



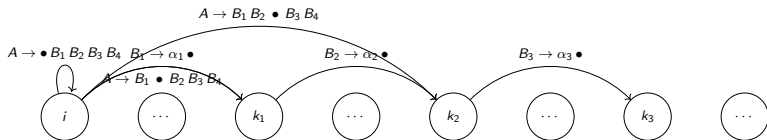
Use of dotted rules

- The use of dotted rules makes it relatively easy to generalise the idea behind CYK to non-binary rules.
- Example:



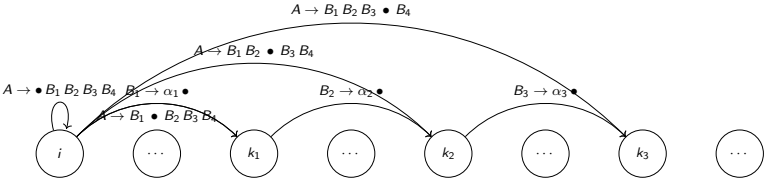
Use of dotted rules

- The use of dotted rules makes it relatively easy to generalise the idea behind CYK to non-binary rules.
- Example:



Use of dotted rules

- The use of dotted rules makes it relatively easy to generalise the idea behind CYK to non-binary rules.
- Example:



Vocabulary

- Inactive item: $(A \rightarrow \alpha \bullet, i, j)$.
- Active item: item that is not inactive.
- Initial item: $(A \rightarrow \bullet \alpha, i, j)$.

Fundamental operation

comp (“complete”)

Fundamental operation

comp (“complete”)

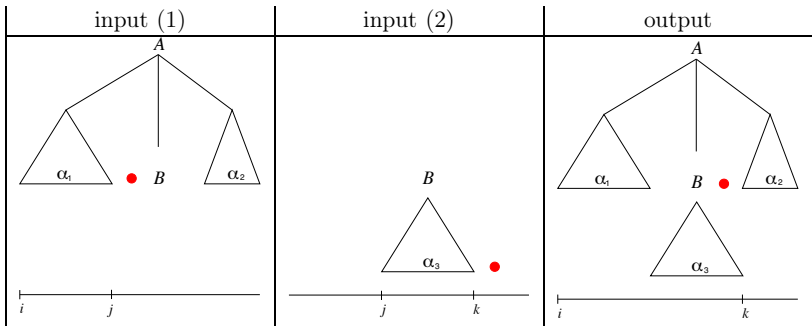
- Input: $(A \rightarrow \alpha_1 \bullet B \alpha_2, i, j)$ and $(B \rightarrow \beta \bullet, j, k)$

Fundamental operation

comp (“complete”)

- Input: $(A \rightarrow \alpha_1 \bullet B \alpha_2, i, j)$ and $(B \rightarrow \beta \bullet, j, k)$
- Output: $(A \rightarrow \alpha_1 B \bullet \alpha_2, i, k)$

Another view on comp



Other essential operation

```
scan
```


Other essential operation

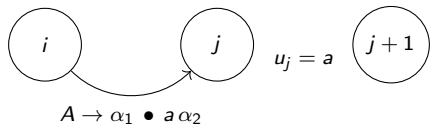
scan

- Input: $(A \rightarrow \alpha_1 \bullet a \alpha_2, i, j)$ provided that $u_j = a$
- Output: $(A \rightarrow \alpha_1 a \bullet \alpha_2, i, j + 1)$

Other essential operation

scan

- Input: $(A \rightarrow \alpha_1 \bullet a \alpha_2, i, j)$ provided that $u_j = a$
- Output: $(A \rightarrow \alpha_1 a \bullet \alpha_2, i, j + 1)$



- comp and scan “advance” existing items.

- comp and scan “advance” existing items.
- How/when are initial items introduced?

- `comp` and `scan` “advance” existing items.
- How/when are initial items introduced?
- → Several versions (i.e. strategies) of the algorithm.

First strategy

- The chart is initialised with all possible initial items (i.e. $(A \rightarrow \bullet \alpha, i, i)$).
- \rightarrow bottom-up parsing (not unlike CYK).

First strategy

Algorithm 1: Simple Earley analysis

```

Function earley-simple( $u$ )
  // Initialisation
   $T :=$  empty chart( $u$ );
  for  $j := 1$  to  $|u| + 1$  do
     $T[j] :=$  ordered_set();
    foreach  $(A \rightarrow \alpha) \in P$  do  $T[j].add((A \rightarrow \bullet \alpha, j))$ ;
  // Main loop
  for  $j := 1$  to  $|u| + 1$  do
     $k := 0$ ;
    while  $k < len(T[j])$  do
       $(A \rightarrow \alpha \bullet \beta, i) := T[j][k]$ ;
      if  $\beta = \epsilon$  then // comp?
         $k' := 0$ ;
        while  $k' < len(T[i])$  do
           $(A' \rightarrow \alpha' \bullet \beta', i') := T[i][k']$ ;
          if  $\beta'_1 = A$  then
             $T[j].add((A' \rightarrow \alpha' \beta'_1 \bullet \beta'_{2:|\beta'|}, i'))$ ;
             $k' += 1$ ;
          else if  $j < |u| + 1$  then // scan?
            if  $\beta_1 = u_j$  then
               $T[j + 1].add((A \rightarrow \alpha \beta_1 \bullet \beta_{2:|\beta|}, i))$ ;
             $k += 1$ ;
    return  $T$ ;

```

First strategy

- Let's analyse *Sabine saw a truck* with a grammar such that

$$P = \left\{ \begin{array}{l} S \rightarrow NP VP, \\ NP \rightarrow DET N \mid PN, \\ VP \rightarrow V \mid VNP, \\ DET \rightarrow the \mid a(n), \\ N \rightarrow truck \mid experiment, \\ PN \rightarrow Sabine \mid Fred \mid Jamy, \\ V \rightarrow saw \mid prepared \end{array} \right\}.$$

Better version

- The original Earley algorithm.
- The only items introduced initially are the ones of the shape $(S \rightarrow \bullet \alpha, 1, 1)$.
- A new operation, `pred` (*predict*), is used to introduce additional initial items.
- `pred` is used to introduce an initial item only if this item may be used to advance an item already introduced.
- \rightarrow bottom-up parsing with top-down information.

Better version

New operation: pred

Better version

Algorithm 2: Earley analysis

Function earley(u)

```

// Initialisation
T := empty chart(u);
for j := 1 to |u| + 1 do T[j] := ordered_set();
foreach ( $S \rightarrow \alpha$ )  $\in P$  do T[1].add( $(S \rightarrow \bullet \alpha, 1)$ );
// Main loop
for j := 1 to |u| + 1 do
  k := 0;
  while k < len(T[j]) do
    ( $A \rightarrow \alpha \bullet \beta, i$ ) :=  $\in T[j][k]$ ;
    if  $\beta = \epsilon$  then // comp?
      k' := 0;
      while k' < len(T[i]) do
        ( $A' \rightarrow \alpha' \bullet \beta', i'$ ) := T[i][k'];
        if  $\beta'_1 = A$  then
          T[j].add( $(A' \rightarrow \alpha' \beta'_1 \bullet \beta'_{2:|\beta'|}, i')$ );
          k' += 1;
        else if  $\beta_1 \in N$  then // pred?
          foreach ( $\beta_1 \rightarrow \gamma$ )  $\in P$  do
            T[j].add( $(\beta_1 \rightarrow \bullet \gamma, j)$ );
          else if j < |u| + 1 then // scan?
            if  $\beta_1 = u_j$  then
              T[j + 1].add( $(A \rightarrow \alpha \beta_1 \bullet \beta_{2:|\beta|}, i)$ );
            k += 1;
  return T;
```
