



# Overview of the course

- Day 1: Formal languages and syntactic complexity.
- Day 2: The complexity of natural language.
- Day 3: Historic algorithms for parsing.
- Day 4: Modern approaches to parsing.
- Day 5: Neural networks and error propagation.



































































































# 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: a complete example

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

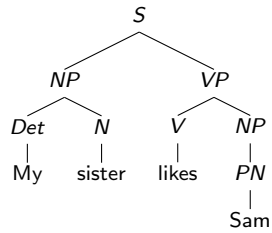
5	$S$	$\emptyset$	$VP$	$\{PN, NP\}$	$\emptyset$
4	$S$	$\emptyset$	$\{V, VP\}$	$\emptyset$	
3	$NP$	$N$	$\emptyset$		
2	$Det$	$\emptyset$			
1	$\emptyset$				
$j/i$	1	2	3	4	5
$w =$	My	sister	likes	Sam	

- 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 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 one of the following forms:

- **(binary rule)**  $A \rightarrow BC$ , with  $A, B, C \in N$ ,
  - **(lexical rule)**  $A \rightarrow a$ , with  $a \in \Sigma$  and  $A \in N$ .
- 
- Any CFG is weakly equivalent to some CFG in CNF.







# 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  $\geq 2$  :
 
$$T[i, j] = \{A \in N \mid \exists k \in [i + 1, j - 1],$$

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

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

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

# 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	$\emptyset$				
<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  $\geq 2$  :
 
$$T[i, j] = \{A \in N \mid \exists k \in [i + 1, j - 1],$$

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

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

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





# 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  $\geq 2$  :
 
$$T[i, j] = \{A \in N \mid \exists k \in [i + 1, j - 1],$$

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

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

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







# 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...

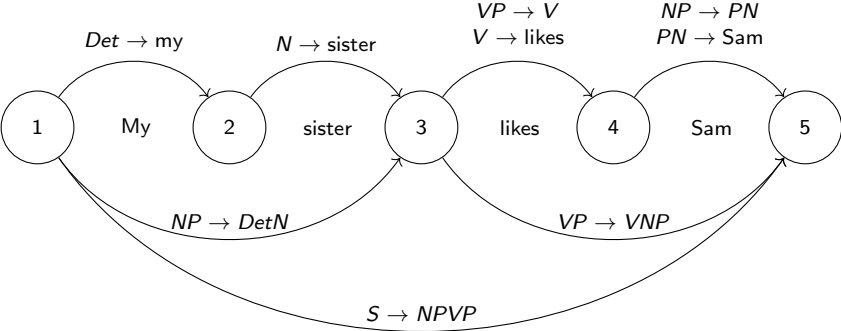
# Earley Algorithm

- 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:  $\exists \alpha, (S \rightarrow \alpha \bullet) \in T[1, |w| + 1]$ .





# Another view on the CYK 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

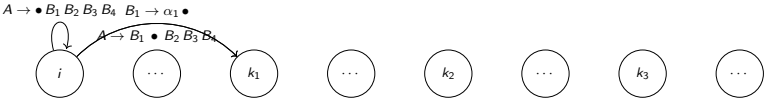
- The use of dotted rules makes it relatively easy to generalise the idea behind CYK to non-binary rules.
- Example:





# Use of dotted rules

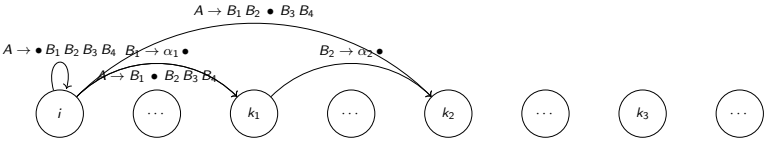
- The use of dotted rules makes it relatively easy to generalise the idea behind CYK to non-binary rules.
- Example:





# Use of dotted rules

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

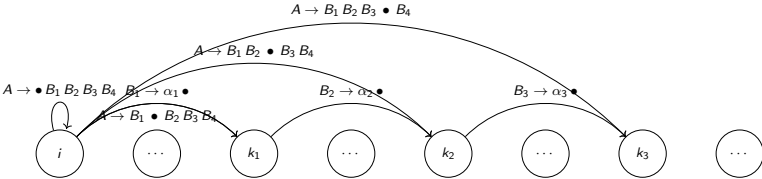






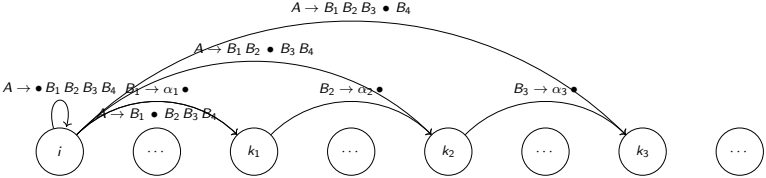
# Use of dotted rules

- 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:



- An Earley item can be interpreted as an hypothesis:  $(A \rightarrow \alpha \bullet \beta, i, j)$  indicates that one is trying to recognise  $A$  starting from  $i$ .

# 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”)





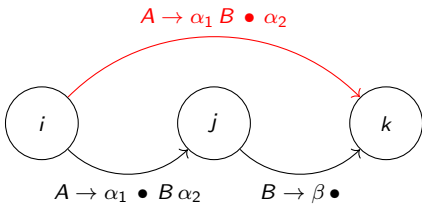




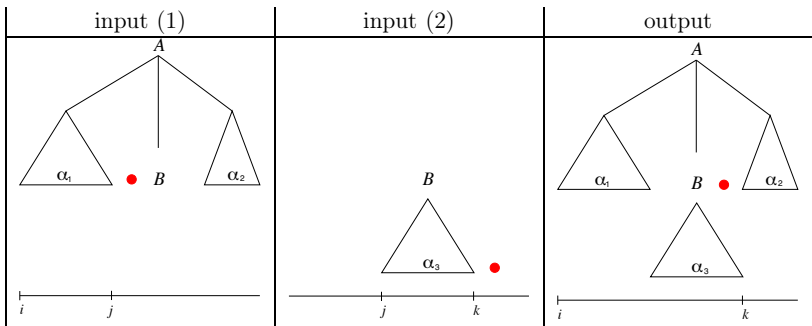
# Fundamental operation

comp ("complete")

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



# Another view on comp





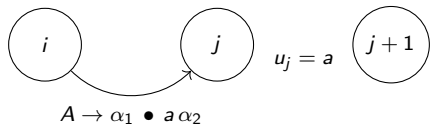




# Other essential operation

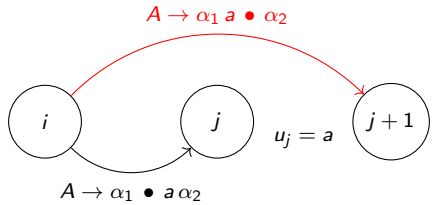
## scan

- 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.









# 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  $T;$ 

```

---





























