

# INTRODUCTION TO NATURAL LANGUAGE PROCESSING

## CHAPTER 11

# Today's Outline

HWs and Project

Review: Earley

Features and Unification

Unification and Parsing

Subcategorization

# Review

Reuse: despite ambiguity and backtracking there are invariants to be taken advantage of to reduce inefficiency.

The Earley algorithm fills a table (the chart) with solutions to sub-problems that:

- does not do repeated work
- does top-down search with bottom-up filtering (sort of)
- solves the left-recursion problem
- solves an exponential problem in less time

# Features and Unification

Recall that there were things that we knew that CFGs could not handle very well.

In particular, agreement and subcategorization could only be handled in a very awkward way.

## Agreement

- *The cat sleeps.*
- *The cats sleep.*

## Subcategorization

- *Cats dream.*
- *Cats like mice.*

# Agreement

What was the problem with the following rules?

$S \rightarrow NP VP$

$NP \rightarrow Det Nominal$

# Agreement

What was the problem with the following rules?

$S \rightarrow NP VP$

$NP \rightarrow Det Nominal$

Answer: Since they don't enforce number and person agreement constraints, they overgenerate and allow the following:

- ★ The cat sleep.
- ★ They sleeps.
- ★ This dogs

# Subcategorization

What was the problem with the following rules?

$VP \rightarrow V$

$VP \rightarrow V NP$

# Subcategorization

What was the problem with the following rules?

$VP \rightarrow V$

$VP \rightarrow V NP$

Answer: Since they don't enforce subcategorization constraints, they allow the following:

- ★ Cats dream mice.
- ★ He took.
- ★ He slept the cot.
- ★ She disappeared the elephant.

## Elegant Solution

We need to constrain the grammar rules to enforce number agreement, subcategorization differences, etc.

We'll do this with feature structures and the constraint-based unification formalism.

## 3sgNP

One way to handle these phenomena in a strictly context-free approach is to encode the constraints into the non-terminal categories and then into the rules. As in the following:

- $S \rightarrow SgS$
- $S \rightarrow PIS$
- $SgS \rightarrow SgNP SgVP$
- $SgNP \rightarrow SgDet SgNom$

and

- $IntransVP \rightarrow IntransVerb$
- $TransVP \rightarrow TransVerb NP$

Why is this inelegant?

# Features

An alternative is to rethink the terminal and non-terminals as complex objects with associated properties (called features) that can be manipulated.

- features take on different values.
- the application of grammar rules is constrained by testing on these features.

This would allow us to code rules as the following:

$$S \rightarrow NP VP$$

*Only if the number of the NP is equal to the number of the VP (that is, the NP and VP agree in number).*

This allows us to have the best of both worlds.

# Features and Feature Structures

We can encode these properties by associating what are called Feature Structures with grammatical constituents.

Features structures are sets of feature-value pairs where:

- the features are atomic symbols and
- the values are either atomic symbols or feature structures.

$$\left[ \begin{array}{ll} \textit{feature}_1 & \textit{value}_1 \\ \textit{feature}_2 & \textit{value}_2 \\ \dots & \\ \textit{feature}_n & \textit{value}_n \end{array} \right]$$

# Example Feature Structures

[ *number sg* ]

[ *number sg*  
*person 3* ]

[ *cat np*  
*number sg*  
*person 3* ]

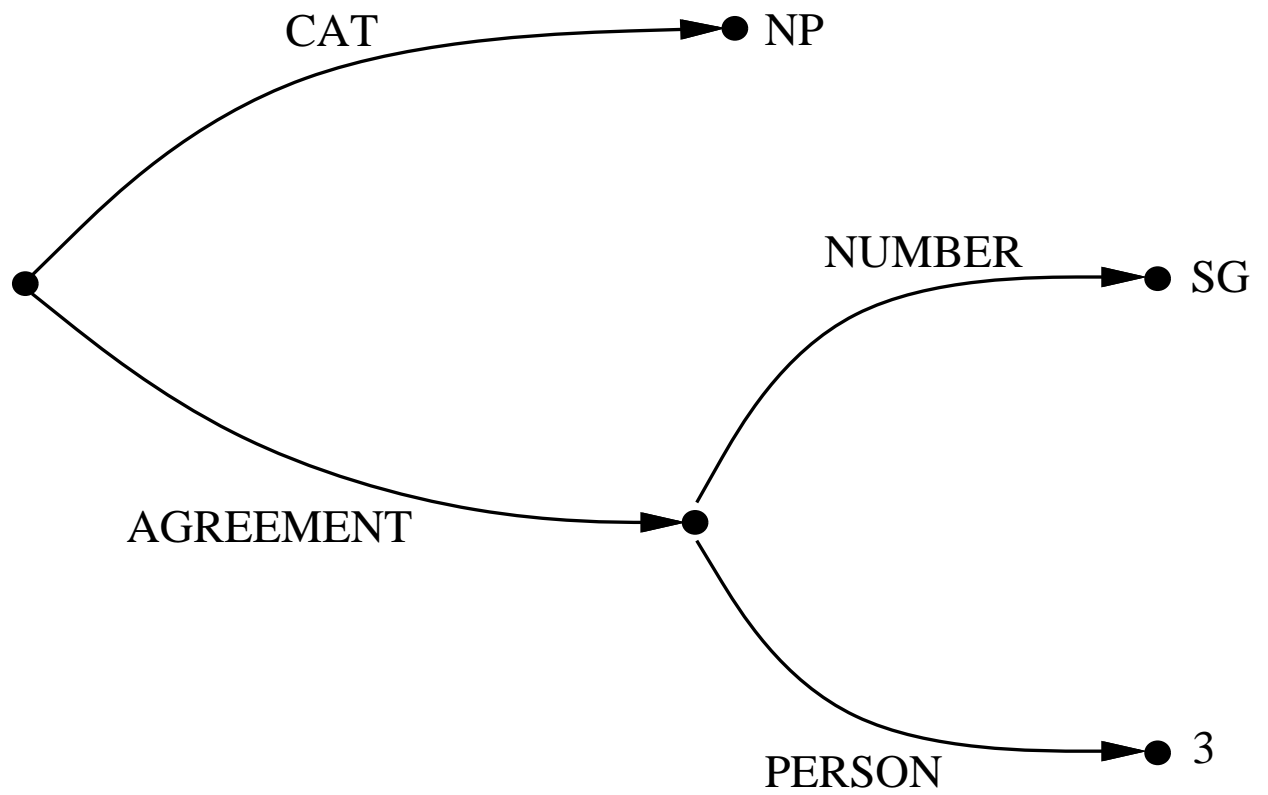
# Bundles of Features

Feature values can be feature structures themselves.

Useful when certain features commonly co-occur, as number and person.

$$\left[ \begin{array}{ll} \textit{cat} & \textit{np} \\ \textit{agreement} & \left[ \begin{array}{ll} \textit{number} & \textit{sg} \\ \textit{person} & 3 \end{array} \right] \end{array} \right]$$

# Feature Structures as DAGs



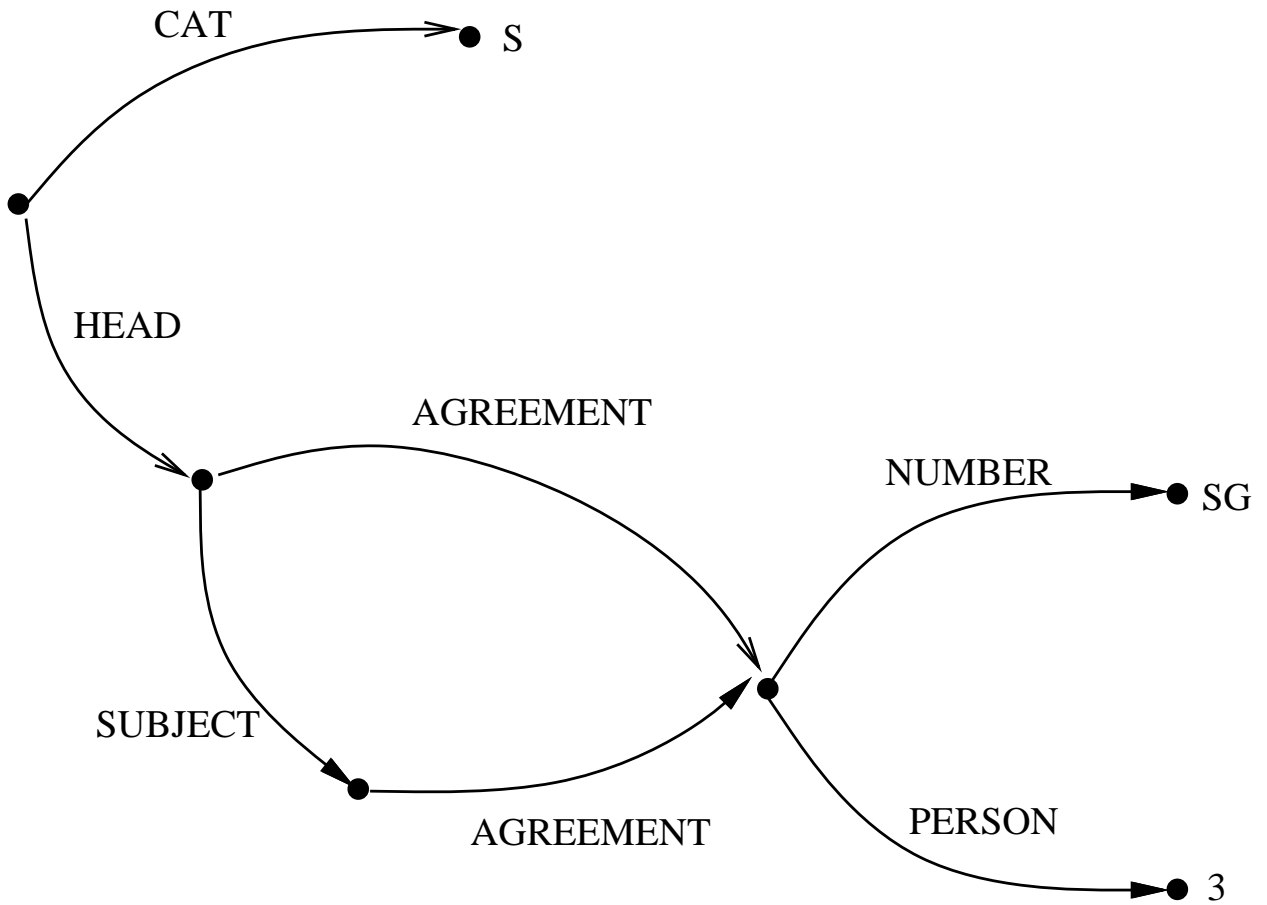
# Reentrant Structures

We'll allow multiple features in a feature structure to share the same values. By this, we mean that they share the same structures not just that they have the same value.

$$\left[ \begin{array}{cc} \textit{cat} & \textit{s} \\ \textit{head} & \left[ \begin{array}{cc} \textit{agreement} & 1 \\ \textit{subject} & \left[ \begin{array}{cc} \textit{agreement} & 1 \end{array} \end{array} \right] \end{array} \right] \end{array} \right]$$

Numerical indices indicate the shared value.

# Reentrant DAGs



# Feature Paths

It will also be useful to talk about paths through feature structures. As in the paths

⟨HEAD AGREEMENT NUMBER⟩

⟨HEAD SUBJECT AGREEMENT NUMBER⟩

$$\left[ \begin{array}{cc} \textit{cat} & \textit{s} \\ \textit{head} & \left[ \begin{array}{cc} \textit{agreement} & 1 \\ \textit{subject} & \left[ \begin{array}{cc} \textit{agreement} & 1 \end{array} \end{array} \right] \end{array} \right] \left[ \begin{array}{cc} \textit{number} & \textit{sg} \\ \textit{person} & 3 \end{array} \right] \end{array} \right]$$

# The Unification Operation

So what do we want to do with these things...

- check the compatibility of two structures
- merge the information in two structures

We can do both with an operation called Unification.

Merging two feature structures produces a new feature structure that is more specific (has more information) than, or is identical to, each of the input feature structures.

# The Unification Operation

We say that two feature structures can be unified if the component features that make them up are compatible.

$[\text{number sg}] \sqcup [\text{number sg}] = [\text{number sg}]$

$[\text{number sg}] \sqcup [\text{number pl}]$  fails!

Structures are compatible if they contain no features that are incompatible.

If so, unification returns the union of all feature/value pairs.

# The Unification Operation

$$[\text{number sg}] \sqcup [\text{number []}] = [\text{number sg}]$$

$$[\text{number sg}] \sqcup [\text{person 3}] = \begin{bmatrix} \textit{number} & \textit{sg} \\ \textit{person} & 3 \end{bmatrix}$$

# The Unification Operation

$$\left[ \begin{array}{l} \textit{agreement} \ 1 \left[ \begin{array}{l} \textit{number} \ \textit{sg} \\ \textit{person} \ 3 \end{array} \right] \\ \textit{subject} \ \left[ \textit{agreement} 1 \right] \end{array} \right]$$

□

$$\left[ \begin{array}{l} \textit{agreement} \left[ \begin{array}{l} \textit{number} \ \textit{sg} \\ \textit{person} \ 3 \end{array} \right] \\ \textit{subject} \ \left[ \textit{agreement} \left[ \begin{array}{l} \textit{number} \ \textit{pl} \\ \textit{person} \ 3 \end{array} \right] \right] \end{array} \right]$$

Fails!

# Properties of Unification

Monotonic: if some description is true of a feature structure, it will still be true after unifying it with another feature structure.

Order independent: given a set of feature structures to unify, we can unify them in any order and we'll get the same result.

# Features, Unification and Grammars

We'll incorporate all this into our grammars in two ways.

- we'll assume that constituents are objects which have feature-structures associated with them
- we'll associate sets of unification constraints with grammar rules that must be satisfied for the rule to be satisfied.

# Unification Constraints

$$\beta_0 \rightarrow \beta_1 \dots \beta_n$$

{ set of constraints }

$\langle \beta_i \text{ feature path} \rangle = \text{atomic value}$

$\langle \beta_i \text{ feature path} \rangle = \langle \beta_j \text{ feature path} \rangle$

# Agreement

NP → Det Nominal

⟨ Det AGREEMENT ⟩ = ⟨ Nominal AGREEMENT ⟩

⟨ NP AGREEMENT ⟩ = ⟨ Nominal AGREEMENT ⟩

Noun → flight

⟨ Noun AGREEMENT NUMBER ⟩ = SG

Noun → flights

⟨ Noun AGREEMENT NUMBER ⟩ = PL

Nominal → Noun

⟨ Nominal AGREEMENT ⟩ = ⟨ Noun AGREEMENT ⟩

Det → this

⟨ Det AGREEMENT NUMBER ⟩ = SG

# Unification and Parsing

OK, let's assume we've augmented our grammar with sets of path-like unification constraints.

What changes do we need to make to a parser to make use of them?

- building feature structures and associate them with subtree
- unifying feature structures as subtrees are created
- blocking ill-formed constituents

# Unification and Earley Parsing

With respect to an Earley-style parser...

- building feature structures (represented as DAGs) and associate them with states in the chart
- unifying feature structures as states are advanced in the chart
- block ill-formed states from entering the chart

# Building Feature Structures

Features of most grammatical categories are copied from head child to parent (e.g., from V to VP, Nom to NP, N to Nom)

- $VP \rightarrow V NP$
- $\langle VP \text{ HEAD} \rangle = \langle V \text{ HEAD} \rangle$

$S \rightarrow NP VP$

$\langle NP \text{ HEAD AGREEMENT} \rangle = \langle VP \text{ HEAD AGREEMENT} \rangle$

$\langle S \text{ HEAD} \rangle = \langle VP \text{ HEAD} \rangle$

corresponds to

$$\begin{bmatrix} s & [head1] \\ np & [head & [agreement2]] \\ vp & [head & 1[agreement2]] \end{bmatrix}$$

## Augmenting States with DAGs

We just add a new field to the representation of the states.

$S \rightarrow \cdot NP VP, [0,0], [], Dag$

# Unifying States and Blocking

We want to unify the DAGs of existing states as they are combined as specified by the grammatical constraints - in the COMPLETER.

```
function EARLEY-PARSE(words, grammar) returns chart
  ENQUEUE( $(\gamma \rightarrow \bullet S, [0, 0], dag_\gamma), chart[0]$ )
  for  $i \leftarrow$  from 0 to LENGTH(words) do
    for each state in chart[i] do
      if INCOMPLETE?(state) and
        NEXT-CAT(state) is not a part of speech then
          PREDICTOR(state)
      elseif INCOMPLETE?(state) and
        NEXT-CAT(state) is a part of speech then
          SCANNER(state)
      else
        COMPLETER(state)
      end
    end
  return(chart)

procedure PREDICTOR( $(A \rightarrow \alpha \bullet B \beta, [i, j], dag_A)$ )
  for each  $(B \rightarrow \gamma)$  in GRAMMAR-RULES-FOR(B, grammar) do
    ENQUEUE( $(B \rightarrow \bullet \gamma, [j, j], dag_B), chart[j]$ )
  end

procedure SCANNER( $(A \rightarrow \alpha \bullet B \beta, [i, j], dag_A)$ )
  if  $B \subset$  PARTS-OF-SPEECH(word[j]) then
    ENQUEUE( $(B \rightarrow word[j], [j, j + 1], dag_B), chart[j + 1]$ )

procedure COMPLETER( $(B \rightarrow \gamma \bullet, [j, k], dag_B)$ )
  for each  $(A \rightarrow \alpha \bullet B \beta, [i, j], dag_A)$  in chart[j] do
    if  $new\_dag \leftarrow$  UNIFY-STATES( $dag_B, dag_A, B$ )  $\neq$  Fails!
      ENQUEUE( $(A \rightarrow \alpha B \bullet \beta, [i, k], new\_dag), chart[k]$ )
    end

procedure UNIFY-STATES(dag1, dag2, cat)
  dag1-cp  $\leftarrow$  COPYDAG(dag1)
  dag2-cp  $\leftarrow$  COPYDAG(dag2)
  UNIFY(FOLLOW-PATH(cat, dag1-cp), FOLLOW-PATH(cat, dag2-cp))

procedure ENQUEUE(state, chart-entry)
  if state is not subsumed by a state in chart-entry then
    PUSH(state, chart-entry)
  end
```

# Modifying Earley

## Completer

- Recall: Completer adds new states to chart by finding states whose dot can be advanced (i.e., category of next constituent matches that of completed constituent)
- Now: Completer will only advance those states if their feature structures unify

Also, new test for whether to enter a state in the chart

- Now DAGs may differ, so check must be more complex
- Don't add states that have DAGs that are more specific than states in chart; is new state subsumed by existing states?

## Example

$NP \rightarrow Det \cdot Nominal[0,1], [S_{Det}], Dag1$

where Dag1 equals

$$\left[ \begin{array}{ll} np & [head1] \\ det & [head \ [agreement2 \ [numbersg]]] \\ nominal & [head \ 1[agreement2]] \end{array} \right]$$

## Example

Nominal  $\rightarrow$  Noun  $\cdot$ , [1,2], [ $S_{Noun}$ ], Dag2

where Dag2 equals

$$\left[ \begin{array}{l} \textit{nominal} \text{ [head1]} \\ \textit{noun} \text{ [head 1[agreement [numbersg]]]} \end{array} \right]$$

and Dag1 as before

$$\left[ \begin{array}{l} \textit{np} \text{ [head1]} \\ \textit{det} \text{ [head [agreement2 [numbersg]]]} \\ \textit{nominal} \text{ [head 1[agreement2]]} \end{array} \right]$$

# Copying

Why the need for all the copying in the procedure UNIFY-STATES?

## Subcategorization Example

Recall: Different verbs take different types of arguments

Some techniques

- atomic subcat symbols
- encoding subcat lists as feature structures

Verbs also subcategorize for subjects, and other POS subcategorize

## Subcat Symbols

VP → Verb

⟨ VP HEAD ⟩ = ⟨ Verb HEAD ⟩

⟨ VP HEAD SUBCAT ⟩ = INTRANS

VP → Verb NP

⟨ VP HEAD ⟩ = ⟨ Verb HEAD ⟩

⟨ VP HEAD SUBCAT ⟩ = TRANS

VP → Verb NP NP

⟨ VP HEAD ⟩ = ⟨ Verb HEAD ⟩

⟨ VP HEAD SUBCAT ⟩ = DITRANS

## Subcat Symbols

Verb → slept

⟨ Verb HEAD SUBCAT ⟩ = INTRANS

Verb → served

⟨ Verb HEAD SUBCAT ⟩ = TRANS

Verb → gave

⟨ Verb HEAD SUBCAT ⟩ = DITRANS

# Subcat Lists

Verb → gave

⟨ Verb HEAD SUBCAT FIRST CAT ⟩ = NP

⟨ Verb HEAD SUBCAT SECOND CAT ⟩ = NP

⟨ Verb HEAD SUBCAT THIRD ⟩ = END

## Subcat Lists

VP  $\rightarrow$  Verb NP NP

$\langle$  VP HEAD  $\rangle = \langle$  Verb HEAD  $\rangle$

$\langle$  VP HEAD SUBCAT FIRST CAT  $\rangle = \langle$  NP CAT  $\rangle$

$\langle$  VP HEAD SUBCAT SECOND CAT  $\rangle = \langle$  NP CAT  $\rangle$

$\langle$  VP HEAD SUBCAT THIRD  $\rangle =$  END

All we're really doing is encoding lists using positional features.

## Summary

Feature structures encode rich information about components of grammar rules

Unification provides a mechanism for merging structures and for comparing them

Feature structures can be quite complex as with subcategorization

Unification parsing can merge or fail, and be incorporated into Earley

## For Next Time

Read parsing paper, project training data, study for exam!