**Parsing**

1. Grammars and parsing

2. Top-down and bottom-up parsing

3. Chart parsers

4. Bottom-up chart parsing

5. The Earley Algorithm

# Syntax

**syntax**: from the Greek *syntaxis*, meaning "setting out together or arrangement."

Refers to the way words are arranged together.

Why worry about syntax?

- The boy ate the frog.

- The frog was eaten by the boy.

- The frog that the boy ate died.

- The boy whom the frog was eaten by died.

# Syntactic Analysis

Key ideas:

- **constituency**: groups of words may behave as a single unit or phrase

- **grammatical relations**: refer to the SUBJECT, OBJECT, INDIRECT OBJECT, etc.

- **subcategorization and dependencies**: refer to certain kinds of relations between words and phrases, e.g. *want* can be followed by an infinitive, but *find* and *work* cannot.

All can be modeled by various kinds of grammars that are based on context-free grammars.
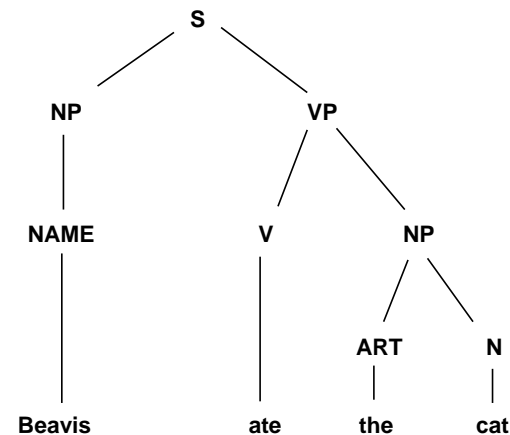
# Grammars and Parsing

Need a **grammar:** a formal specification of the structures allowable in the language.

Need a **parser**: algorithm for assigning syntactic structure to an input sentence.

**Sentence**                                    **Parse Tree**

Beavis ate the cat.

```
                    S
              /            \
           NP               VP
           |             /      \
          NAME          V        NP
           |            |       /    \
                        |     ART     N
                        |      |      |
         Beavis        ate    the    cat
```

# CFG example

CFG's are also called phrase-structure grammars.
Equivalent to Backus-Naur Form (BNF).

1. S → NP VP
2. VP → V NP
3. NP → NAME
4. NP → ART N

5. NAME → Beavis
6. V → ate
7. ART → the
8. N → cat

- CFG's are *powerful* enough to describe most of the structure in natural languages.

- CFG's are *restricted* enough so that efficient parsers can be built.

# CFG's

A context free grammar consists of:

1. a set of non-terminal symbols $N$

2. a set of terminal symbols $\Sigma$ (disjoint from $N$)

3. a set of productions, $P$, each of the form $A \rightarrow \alpha$, where A is a non-terminal and $\alpha$ is a string of symbols from the infinite set of strings $(\Sigma \cup N)^*$

4. a designated start symbol $S$

# Derivations

- If the rule $A \to \beta \in P$, and $\alpha$ and $\gamma$ are strings in the set $(\Sigma \cup N)^*$, then we say that $\alpha A \gamma$ **directly derives** $\alpha \beta \gamma$, or $\alpha A \gamma \Rightarrow \alpha \beta \gamma$

- Let $\alpha_1, \alpha_2, \ldots, \alpha_m$ be strings in $(\Sigma \cup N)^*$, $m > 1$, such that

$$\alpha_1 \Rightarrow \alpha_2, \alpha_2 \Rightarrow \alpha_3, \ldots, \alpha_{m-1} \Rightarrow \alpha_m,$$

then we say that $\alpha_1$ **derives** $\alpha_m$ or $\alpha_1 \overset{*}{\Rightarrow} \alpha_m$

$$L_G$$

The language $L_G$ generated by a grammar $G$ is the set of strings composed of terminal symbols that can be derived from the designated start symbol $S$.

$$L_G = \{w | w \in \Sigma^*, S \overset{*}{\Rightarrow} w\}$$

Parsing: the problem of mapping from a string of words to its parse tree according to a grammar G.

# General Parsing Strategies

| Grammar | Top-Down | Bottom-Up |
|---|---|---|
| 1. S → NP VP | S → NP VP | → NAME ate the cat |
| 2. VP → V NP | → NAME VP | → NAME V the cat |
| 3. NP → NAME | → Beav VP | → NAME V ART cat |
| 4. NP → ART N | → Beav V NP | → NAME V ART N |
| 5. NAME → Beavis | → Beav ate NP | → NP V ART N |
| 6. V → ate | → Beav ate ART N | → NP V NP |
| 7. ART → the | → Beav ate the N | → NP VP |
| 8. N → cat | → Beav ate the cat | → S |