

Parsing: Tokens to Trees

Matt Might

University of Utah

matt.might.net



Today

- Context-free languages
- Context-free grammars
- Notations for grammars
- Syntax trees
- Recursive-descent parsers
- Racket lexers



What is parsing?

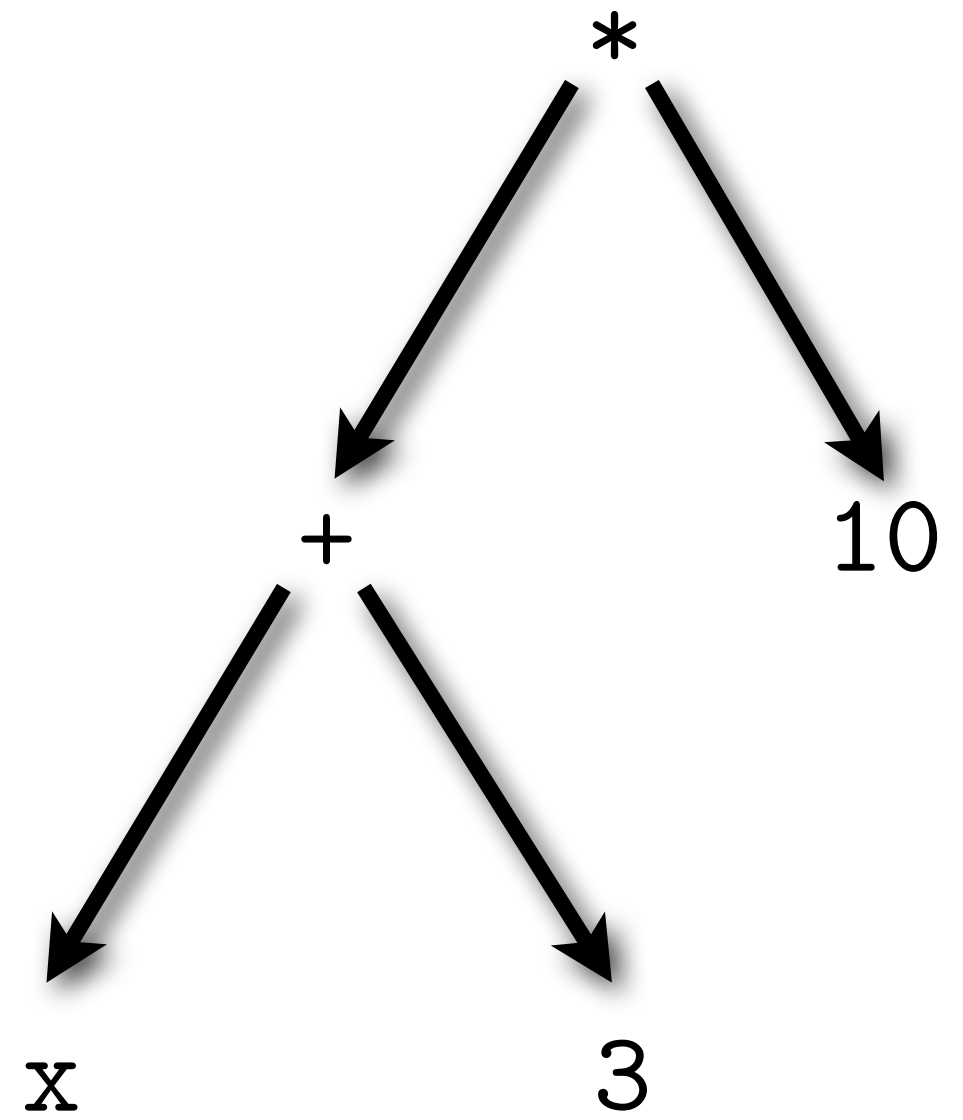
Unstructure to structure.

**A transformation from
a sequence to a tree.**

$$(x + 3) * 10$$

$$(x + 3) * 10$$

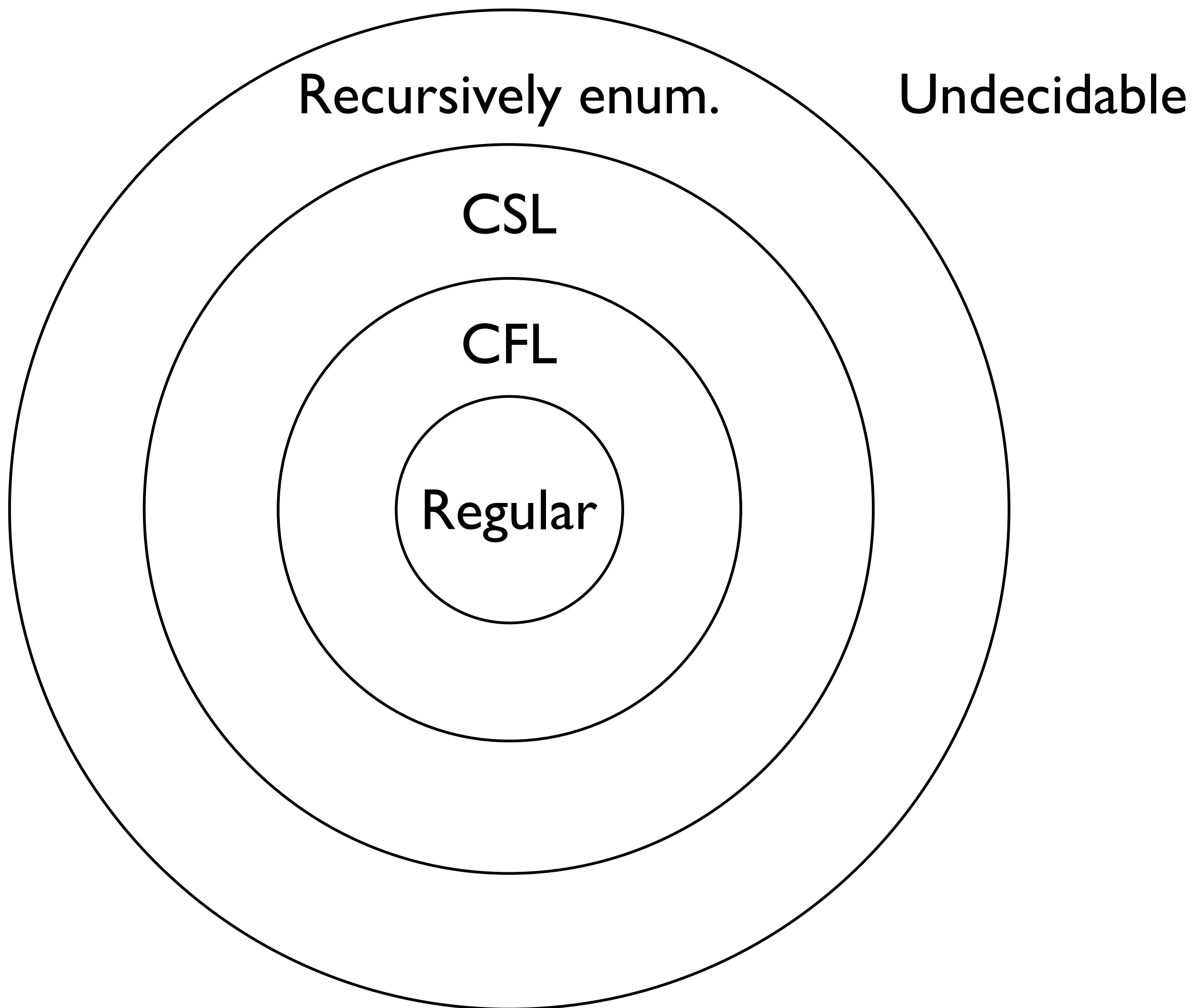
$$(x + 3) * 10$$

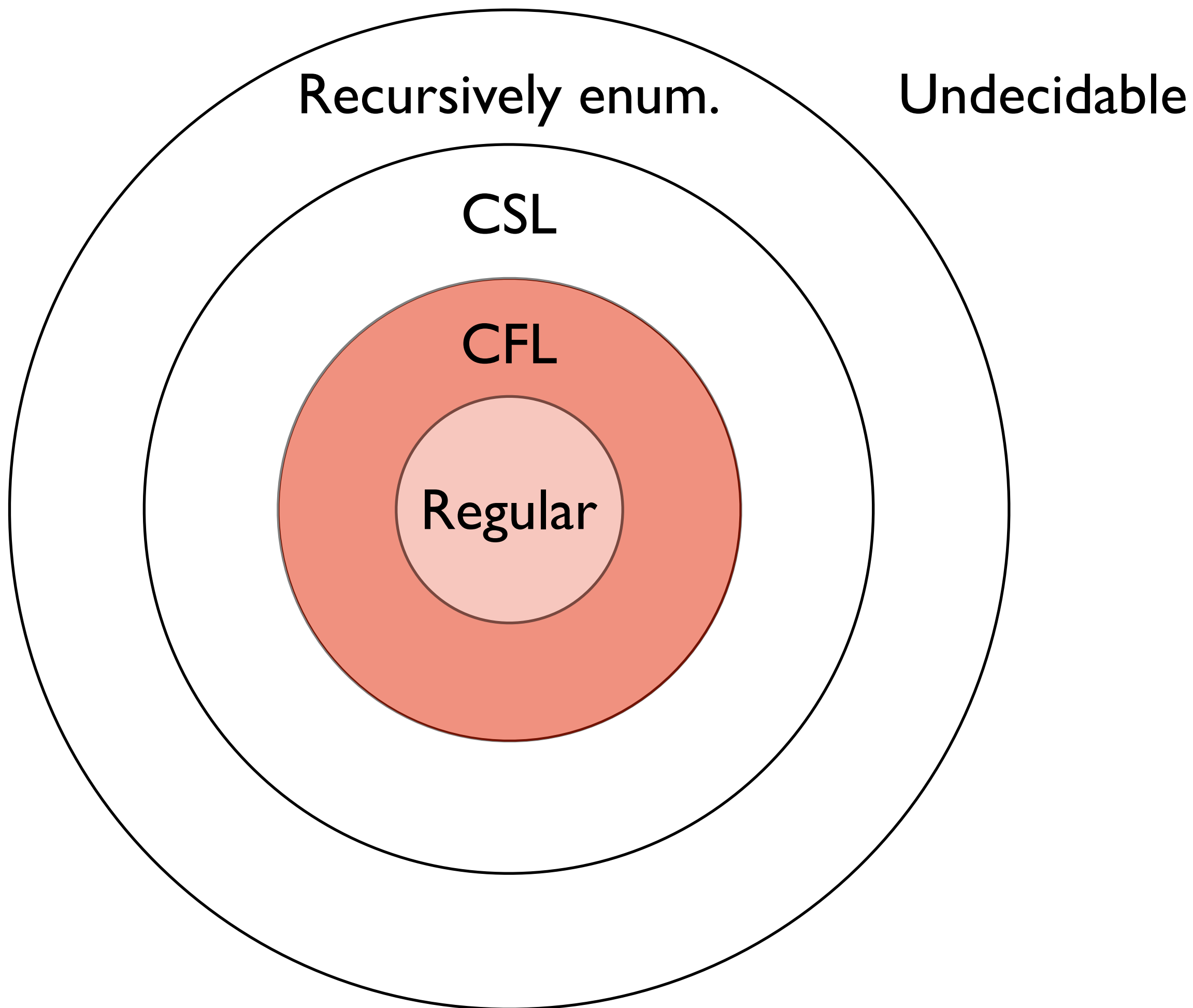


grammars :: tree structure

What's a context-free grammar?

Recursive regular expressions.

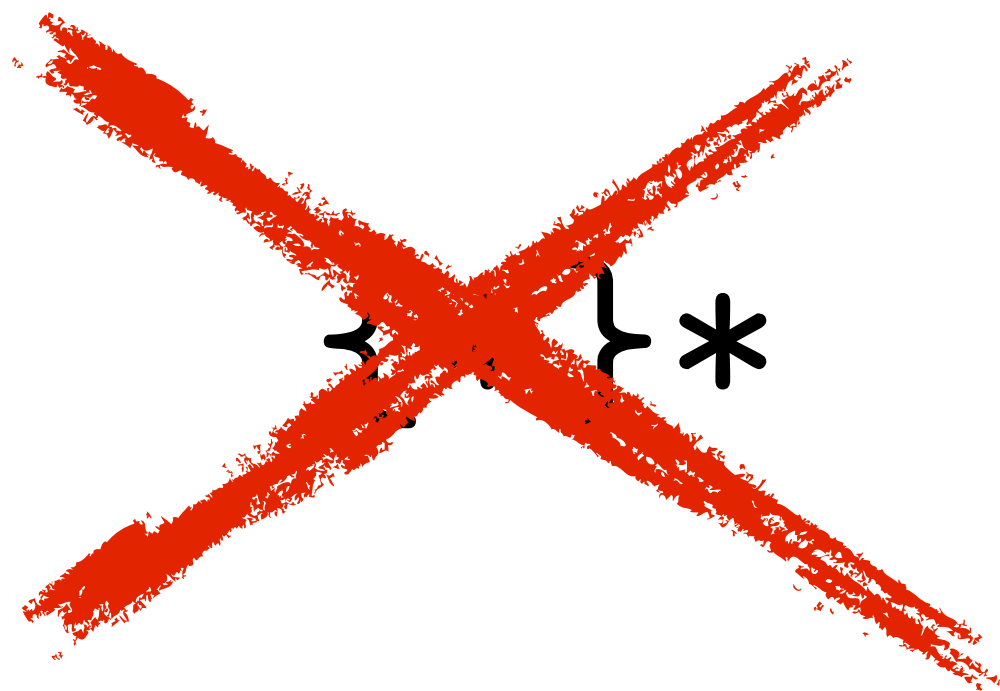




Regex for balanced braces?

$| \{ \} | \{ \{ \} \} | \{ \{ \{ \} \} \} | \dots$

{*}*



$$S ::= \{S\} \mid \varepsilon$$

CFG ingredients

- Concatenation
- Union
- Recursion

Kleene star?

... *A** ...

$\dots B \dots$

$$B :: = AB$$

$$| \quad \varepsilon$$

Context-free grammar

- **Terminals:** Set of atomic symbols
- **Nonterminals:** Phrase types
- **Productions:** Phrase structure rules
- **Start symbol:** Top-level phrase

Notations

- Backus-Naur Form (BNF)
- Extended Backus-Naur Form (EBNF)
- Regular extensions to BNF
- Augmented Backus-Naur Form (ABNF)
- Phrase structure rules
- Set-based language construction

Backus-Naur Form

$$\langle nonterm \rangle ::= exp_1 \dots exp_2$$

Backus-Naur Form

exp is *<nonterm>* or *term*

Backus-Naur Form

term is "literal"

EBNF

exp is *nonterm*

or *term*

or { *exp* }

or [*exp*]

or (*exp*)

or *exp* | *exp*

Regular EBNF

exp is *nonterm*

or *term*

or *exp** or { *exp* }

or *exp*+

or *exp*? or [*exp*]

or (*exp*)

or *exp* | *exp*

ABNF

exp is *nonterm*

or *term*

or *n*m exp*

or *n* exp*

or ** exp*

or *[exp]*

or *(exp)*

or *exp / exp*

Phrase structure rules

$$NT \rightarrow S_1 \dots S_n$$

Phrase structure rules

S is nonterminal or terminal

Set-based languages

$$L' = \epsilon \text{ or } \{\langle \rangle\}$$

$$L' = c \text{ or } \{c\}$$

$$L' = \emptyset \text{ or } \{\}$$

$$L' = L_1 \cup L_2$$

$$L' = L_1 \circ L_2$$

Examples

S-Expression

$\langle S \rangle ::= "(\langle S \rangle^*)" \mid \langle A \rangle$

$\langle A \rangle ::= \text{SYMBOL}$
 $\mid \text{NUMBER}$

Terminal symbols

$\langle S \rangle ::= \text{" ("} \langle S \rangle^* \text{") " } \mid \langle A \rangle$

$\langle A \rangle ::= \text{SYMBOL}$

$\mid \text{NUMBER}$

Non-terminal symbols

$\langle S \rangle ::= "(" \langle S \rangle^* ") " \mid \langle A \rangle$

$\langle A \rangle ::= \text{SYMBOL}$
 $\mid \text{NUMBER}$

Non-terminal symbols

$\langle S \rangle ::= "(\langle S \rangle^*)" \mid \langle A \rangle$
 $\langle A \rangle ::= \text{SYMBOL}$
 $\quad \mid \text{NUMBER}$

Productions

$\langle S \rangle ::= "(\langle S \rangle^*)" \mid \langle A \rangle$

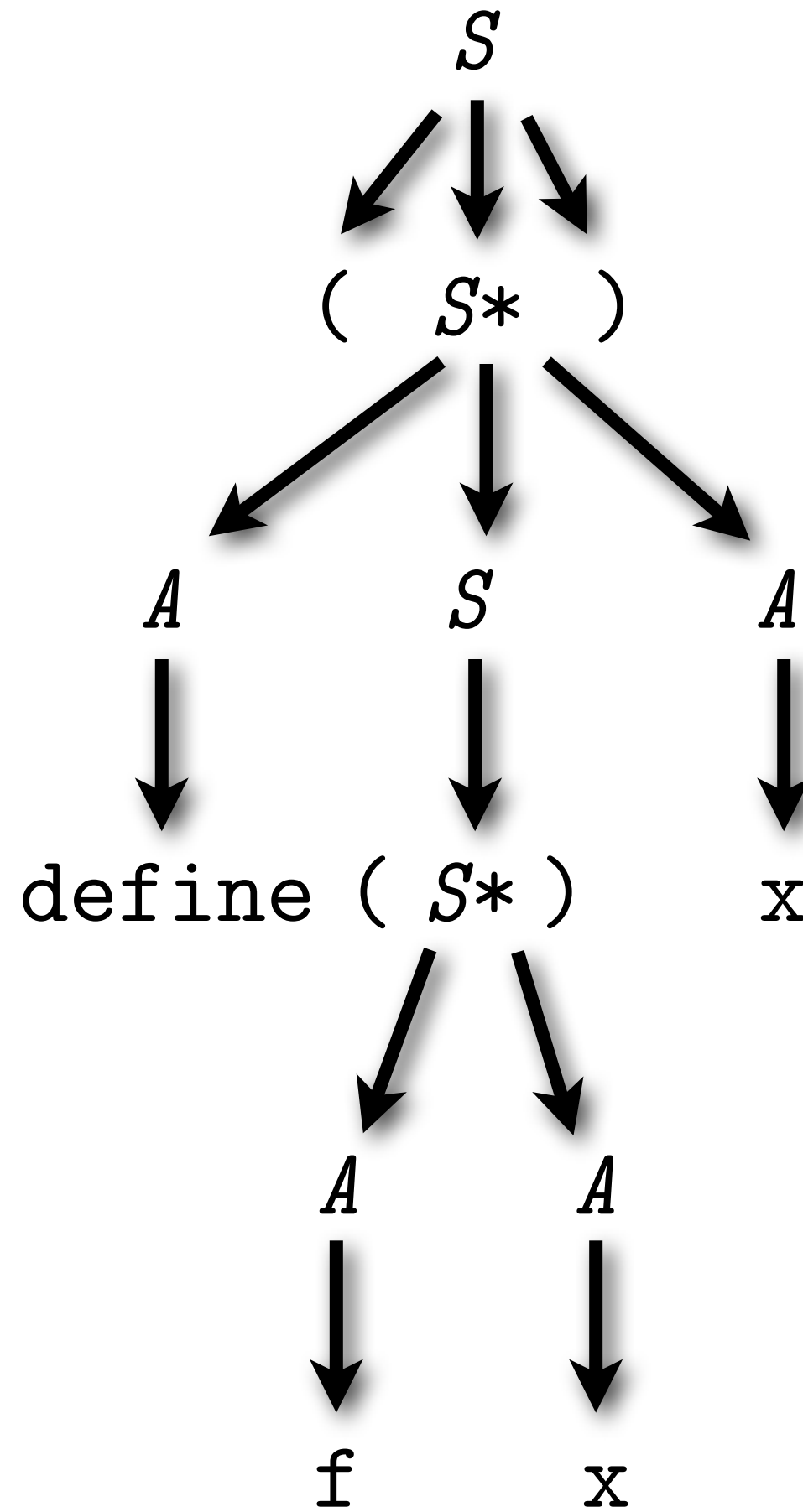
$\langle A \rangle ::= \text{SYMBOL}$
 $\quad \mid \text{NUMBER}$

Start symbol

$\langle S \rangle ::= "(" \langle S \rangle^* ") " \mid \langle A \rangle$

$\langle A \rangle ::= \text{SYMBOL}$
 $\mid \text{NUMBER}$

```
(define (id x) x)
```

Terms

$$E \rightarrow T + E$$

$$E \rightarrow T$$

$$T \rightarrow F \times T$$

$$T \rightarrow F$$

$$F \rightarrow (E)$$

$$F \rightarrow n$$

Terminal symbols

$$E \rightarrow T \boxed{+} E$$

$$E \rightarrow T$$

$$T \rightarrow F \boxed{\times} T$$

$$T \rightarrow F$$

$$F \rightarrow \boxed{(} E \boxed{)}$$

$$F \rightarrow \boxed{n}$$

Nonterminal symbols

$$\boxed{E} \rightarrow T + E$$

$$\boxed{E} \rightarrow T$$

$$\boxed{T} \rightarrow F \times T$$

$$\boxed{T} \rightarrow F$$

$$\boxed{F} \rightarrow (E)$$

$$\boxed{F} \rightarrow n$$

Nonterminal symbols

$$E \rightarrow T + E$$

$$E \rightarrow T$$

$$T \rightarrow F \times T$$

$$T \rightarrow F$$

$$F \rightarrow (E)$$

$$F \rightarrow n$$

Productions

$$E \rightarrow T + E$$

$$E \rightarrow T$$

$$T \rightarrow F \times T$$

$$T \rightarrow F$$

$$F \rightarrow (E)$$

$$F \rightarrow n$$

Start symbol

$$\boxed{E} \rightarrow T + E$$

$$\boxed{E} \rightarrow T$$

$$T \rightarrow F \times T$$

$$T \rightarrow F$$

$$F \rightarrow (E)$$

$$F \rightarrow n$$

Balanced parens

$$P = (\circ P \circ) \circ P \cup \epsilon$$

Terminals

$$P = \boxed{(} \circ P \circ \boxed{)} \circ P \cup \epsilon$$

Nonterminals

$$\boxed{P} = (\circ \boxed{P} \circ) \circ \boxed{P} \cup \epsilon$$

Productions

$$P = (\circ P \circ) \circ P \cup \epsilon$$

Start symbol

$$\boxed{P} = (\circ P \circ) \circ P \cup \epsilon$$

Syntax trees

Syntax trees encode structure

Expressions

$$E \rightarrow T + E$$

$$E \rightarrow T$$

$$T \rightarrow F \times T$$

$$T \rightarrow F$$

$$F \rightarrow (E)$$

$$F \rightarrow n$$

```
abstract class Exp {}
```

```
class Sum extends Exp {  
    Exp left ;  
    Exp right ;  
}
```

```
class Prod extends Exp {  
    Exp left ;  
    Exp right ;  
}
```

```
class Int extends Exp {  
    int value ;  
}
```



```
(struct sum      {left right})  
(struct product {left right})  
(struct int      {value})
```

How to parse
(the hard way)

Recursive descent

- One function per nonterminal
- Allowed to look one token ahead

Recursive descent

- **peek()** : preview next token
- **next()** : consume & return next token
- **eat(token)** : match next token

Recursive descent example

S-Expressions

$\langle S \rangle ::= "(" \langle SL \rangle ") "$

$| \langle A \rangle$

$\langle SL \rangle ::= \langle S \rangle \langle SL \rangle$

$|$

$\langle A \rangle ::= \text{SYMBOL}$

$|$

INTEGER

$|$

$\#t$

$|$

$\#f$

S-Expressions

S

SL

A

S-Expressions

parse S()

parse SL()

parse A()