

Parsing: Tokens to Trees

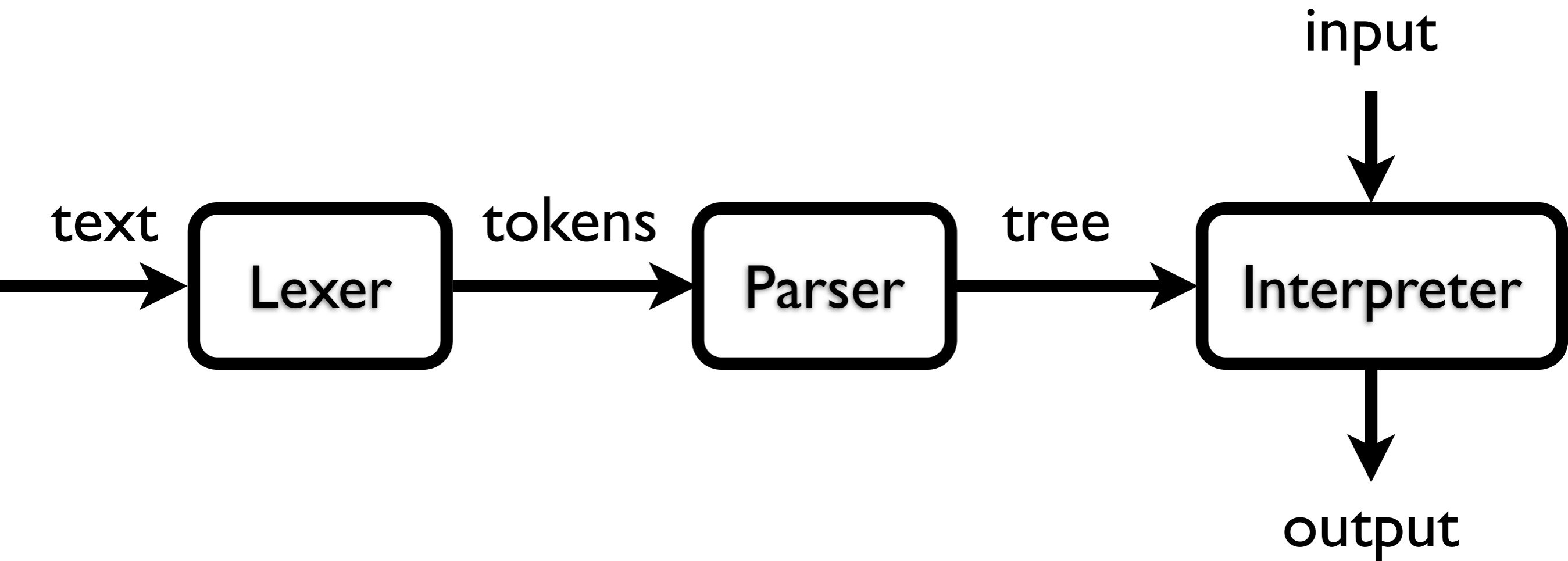
Matt Might
University of Utah
matt.might.net



Today

- Context-free languages
- Context-free grammars
- Recursive-descent parsers
- LALR parser-generators

Interpreters



What is parsing?

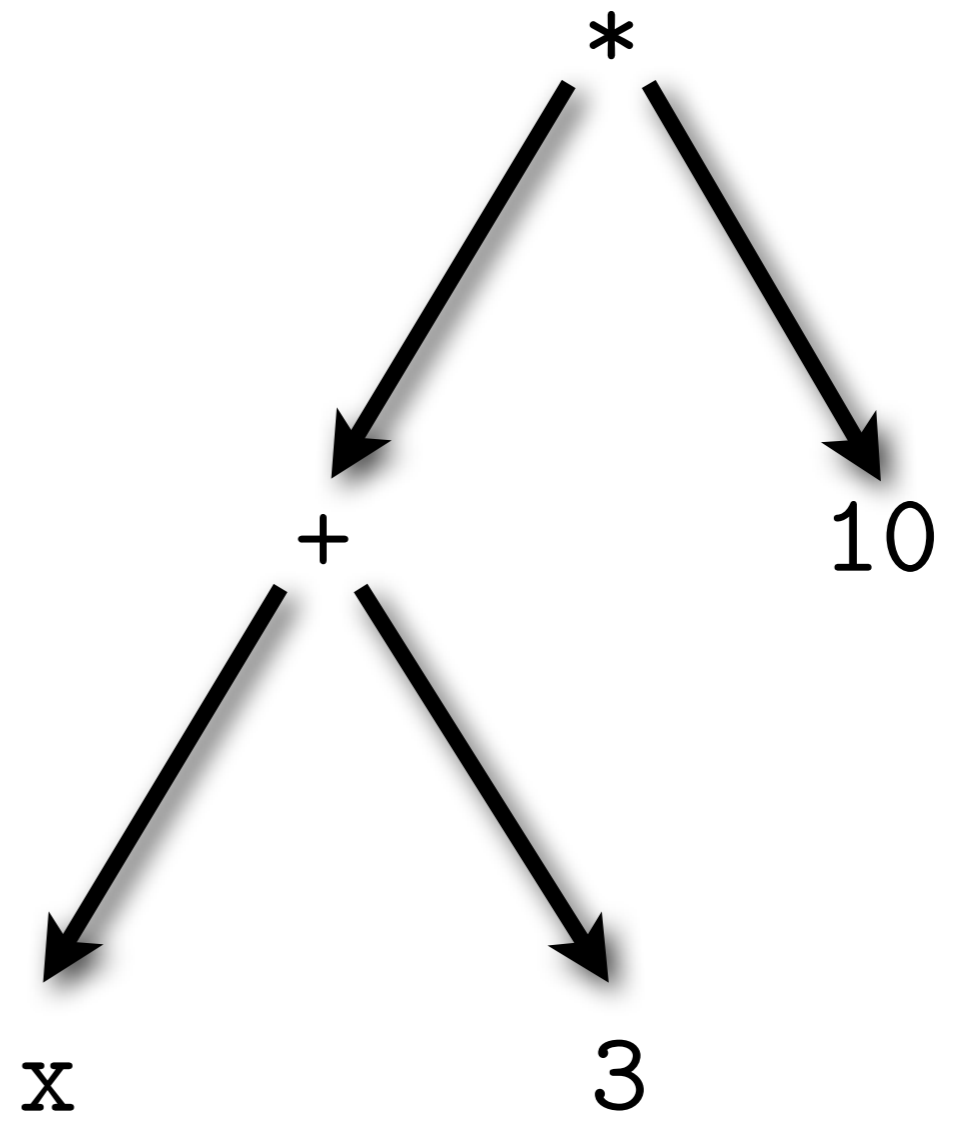
Unstructure to structure.

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

$$(x + 3) * 10$$

$$(x + 3) * 10$$

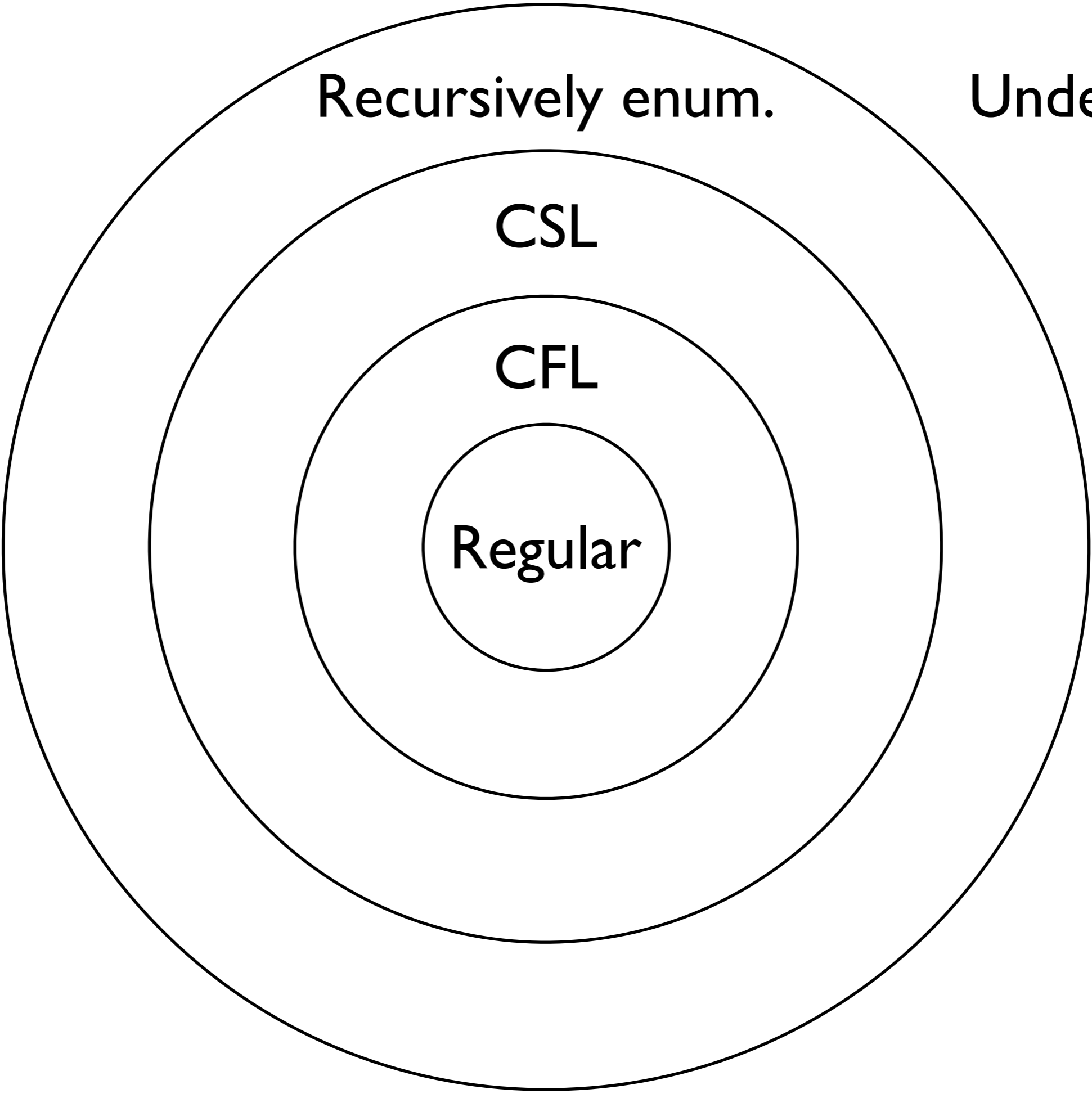
$$(x + 3) * 10$$



grammars :: tree structure

What's a context-free grammar?

Recursive regular expressions.



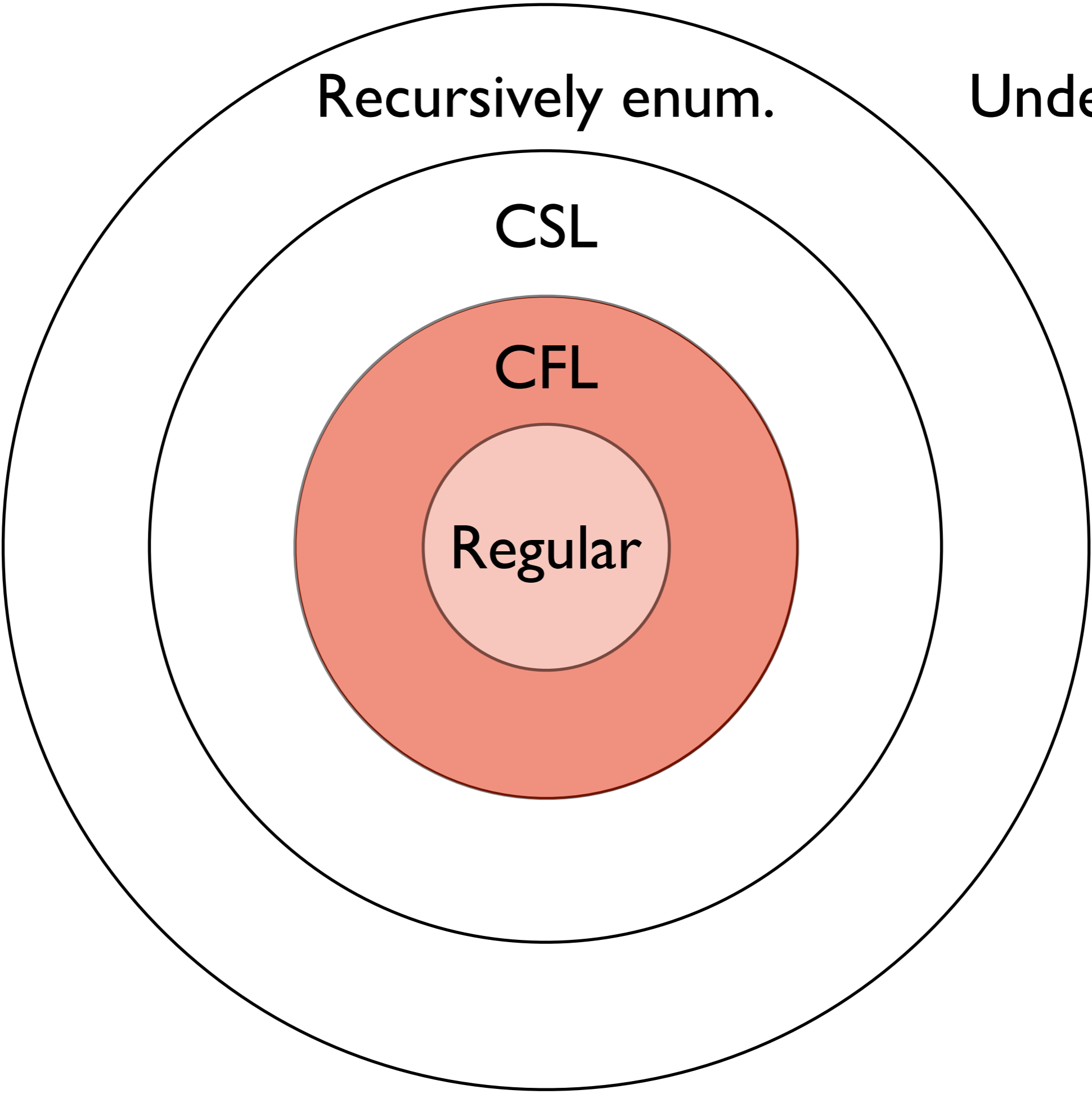
Recursively enum.

Undecidable

CSL

CFL

Regular



Recursively enum.

Undecidable

CSL

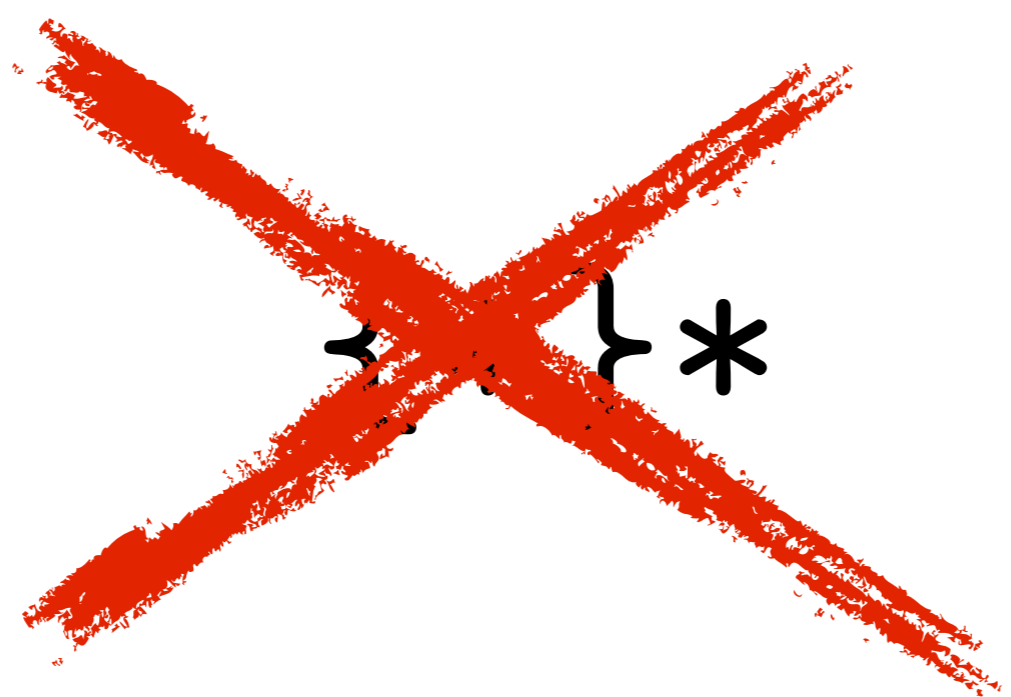
CFL

Regular

Regex for balanced braces?

| {} | {{}} | {{{}}} | ...

{*}*



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

CFG ingredients

- Sequence
- Alternation
- Recursion

Kleene star?

• • • ***A**** • • •

• • • *B* • • •

$$B ::= AB$$
$$| \epsilon$$

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)
- Phrase structure rules
- Set-based language construction

Backus-Naur Form

$\langle \text{nonterm} \rangle ::= \text{exp}_1 \dots \text{exp}_2$

Backus-Naur Form

exp is *<nonterm>* or *term*

EBNF

exp is $\langle nonterm \rangle$

or *term*

or *exp** or { *exp* }

or *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

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

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

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

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

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

Productions

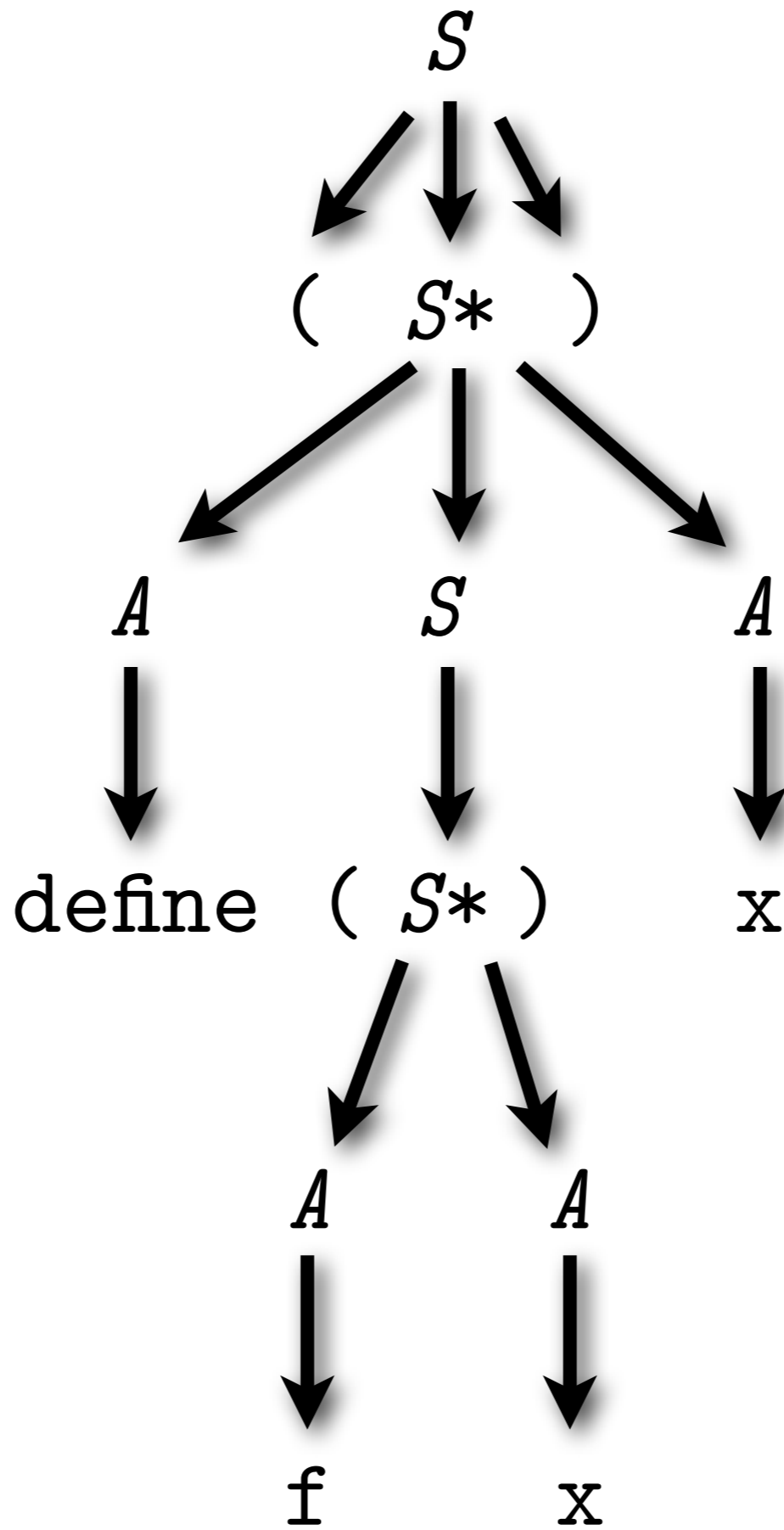
$\langle S \rangle ::= \text{" (" } \langle S \rangle^* \text{ ") " } \mid \langle A \rangle$

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

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

$$E \rightarrow T + E$$

$$E \rightarrow T$$

$$T \rightarrow F \times T$$

$$T \rightarrow F$$

$$F \rightarrow (E)$$

$$F \rightarrow n$$

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

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

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 examples

Expressions

$$E \rightarrow T + E$$

$$E \rightarrow T$$

$$T \rightarrow F \times T$$

$$T \rightarrow F$$

$$F \rightarrow (E)$$

$$F \rightarrow n$$

S-Expressions

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

$| \langle A \rangle$

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

$|$

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

$| \text{INTEGER}$

$| \#t$

$| \#f$