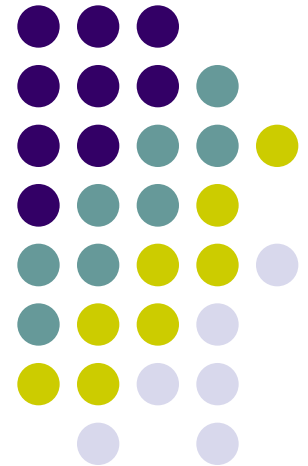
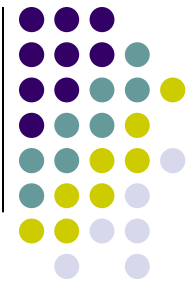


A brief yacc tutorial

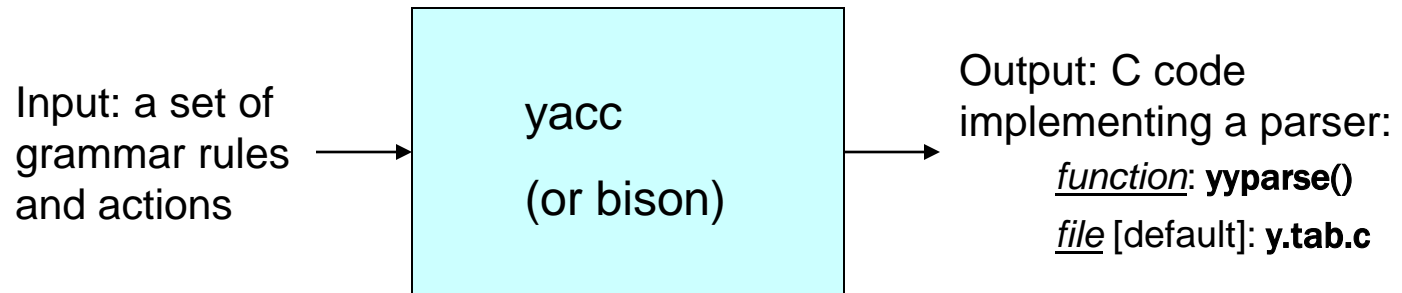


Yacc: Overview

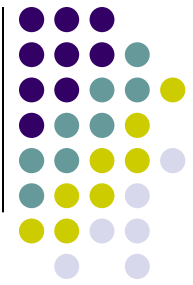


Parser generator:

- ▶ Takes a specification for a context-free grammar.
- ▶ Produces code for a parser.

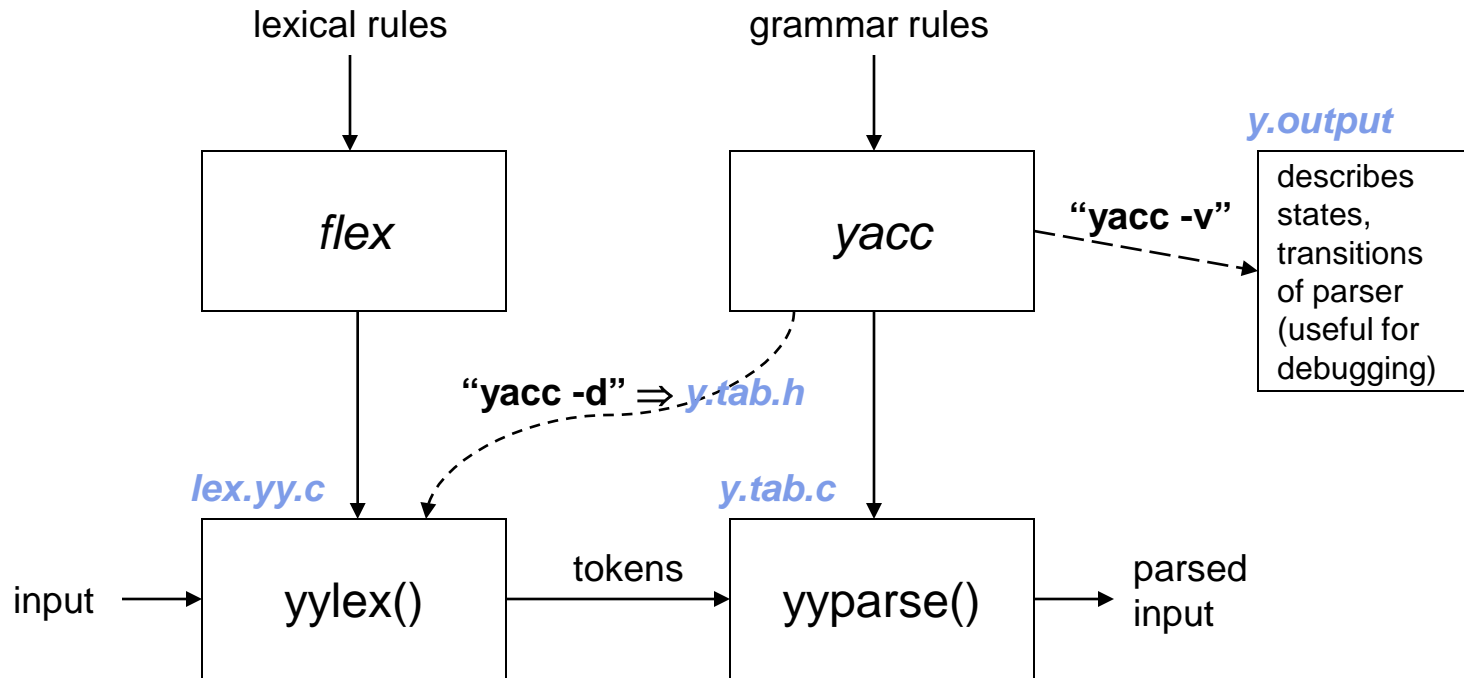


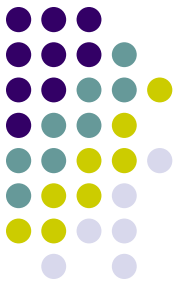
Scanner-Parser interaction



- Parser assumes the existence of a function **'int yylex()'** that implements the scanner.
- Scanner:
 - ▶ return value indicates the type of token found;
 - ▶ other values communicated to the parser using **yytext**, **yyval** (see man pages).
- Yacc determines integer representations for tokens:
 - ▶ Communicated to scanner in file **y.tab.h**
 - use **"yacc -d"** to produce **y.tab.h**
 - ▶ Token encodings:
 - "end of file" represented by '0';
 - a character literal: its ASCII value;
 - other tokens: assigned numbers ≥ 257 .

Using Yacc





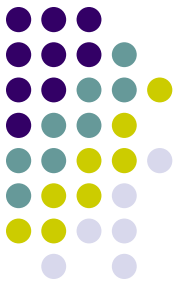
int yyparse()

- Called once from main() [*user-supplied*]
- Repeatedly calls yylex() until done:
 - ▶ On syntax error, calls yyerror() [*user-supplied*]
 - ▶ Returns 0 if all of the input was processed;
 - ▶ Returns 1 if aborting due to syntax error.

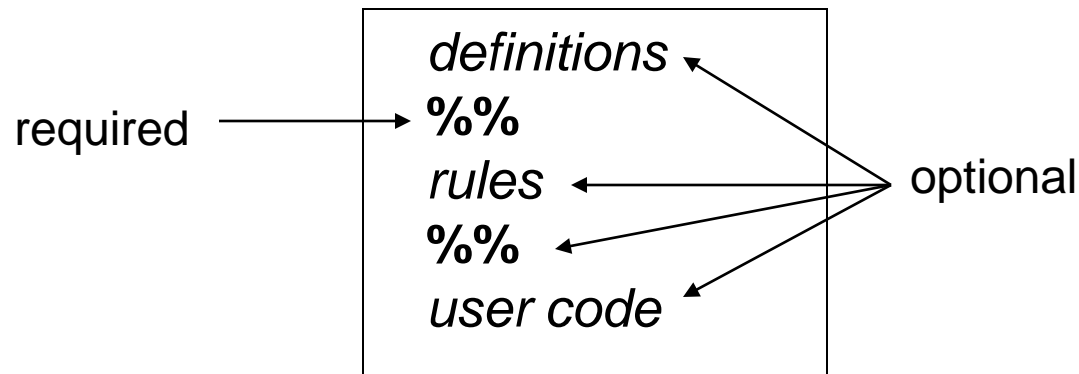
Example:

```
int main() { return yyparse(); }
```

yacc: input format



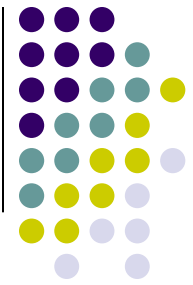
A yacc input file has the following structure:



Shortest possible legal yacc input:

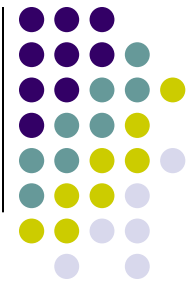
%%

Definitions



- Information about tokens:
 - ▶ token names:
 - declared using '**%token**'
 - single-character tokens don't have to be declared
 - any name not declared as a token is assumed to be a nonterminal.
 - ▶ start symbol of grammar, using '**%start**' [optional]
 - ▶ operator info:
 - precedence, associativity
 - ▶ stuff to be copied verbatim into the output (e.g., declarations, **#includes**): enclosed in **%{ ... }%**

Rules



Grammar production

$$A \rightarrow B_1 B_2 \dots B_m$$
$$A \rightarrow C_1 C_2 \dots C_n$$
$$A \rightarrow D_1 D_2 \dots D_k$$


yacc rule

$$A \rightarrow B_1 B_2 \dots B_m$$
$$/ C_1 C_2 \dots C_n$$
$$/ D_1 D_2 \dots D_k$$

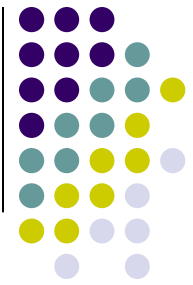
; /* ';' optional, but advised */

- Rule RHS can have arbitrary C code embedded, within { ... }. E.g.:

A : B1 { printf("after B1\n"); x = 0; } B2 { x++; } B3

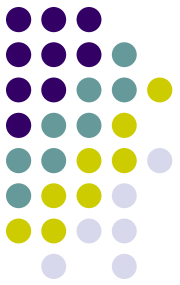
- Left-recursion more efficient than right-recursion:
 - A : A x | ... rather than A : x A | ...

Conflicts



- Conflicts arise when there is more than one way to proceed with parsing.
- Two types:
 - ▶ shift-reduce [default action: *shift*]
 - ▶ reduce-reduce [default: *reduce with the first rule listed*]
- Removing conflicts:
 - ▶ specify operator precedence, associativity;
 - ▶ restructure the grammar
 - use **y.output** to identify reasons for the conflict.

Specifying Operator Properties



- Binary operators: **%left**, **%right**, **%nonassoc**:

`%left '+' '-'`

`%left '*' '/'`

`%right '^'`

{ Operators in the same group
have the same precedence

- Unary operators: **%prec**

- ▶ Changes the precedence of a rule to be that of the token specified. E.g.:

`%left '+' '-'`

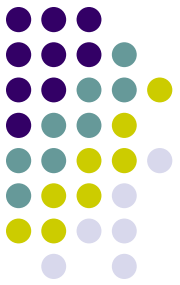
`%left '*' '/'`

`Expr: expr '+' expr`

`| '-' expr %prec '*'`

`| ...`

Specifying Operator Properties



- Binary operators: **%left**, **%right**, **%nonassoc**:

`%left '+' '-'`

`%left '*' '/'`

`%right '^'`

Operators in the same group
have the same precedence

Across groups, precedence
increases going down.

- Unary operators: **%prec**

- Changes the precedence of a rule to be that of the token specified. E.g.:

`%left '+' '-'`

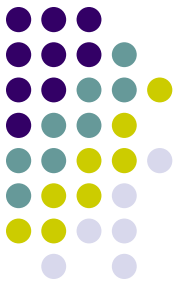
`%left '*' '/'`

Expr: expr '+' expr

| '-' expr %prec '*'

| ...

Specifying Operator Properties



- Binary operators: **%left**, **%right**, **%nonassoc**:

%left '+' '-'

%left '*' '/'

%right '^'

Operators in the same group
have the same precedence

Across groups, precedence
increases going down.

- Unary operators: **%prec**

- Changes the precedence of a rule to be that of the token specified. E.g.:

%left '+' '-'

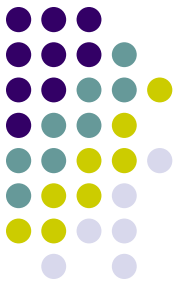
%left '*' '/'

Expr: expr '+' expr

| '-' expr **%prec** '*'

| ...

The rule for unary '-' has the
same (high) precedence as '*'



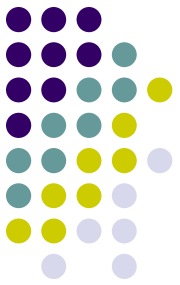
Error Handling

- The “token” ‘error’ is reserved for error handling:
 - ▶ can be used in rules;
 - ▶ suggests places where errors might be detected and recovery can occur.

Example:

```
stmt : IF '(' expr ')' stmt  
      | IF '(' error ')' stmt  
      | FOR ...  
      | ...
```

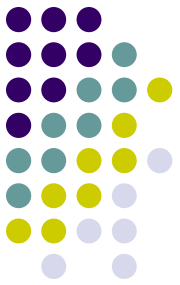
Intended to recover from errors in ‘expr’



Parser Behavior on Errors

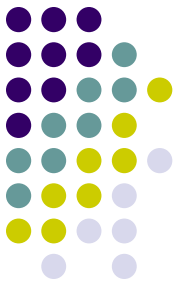
When an error occurs, the parser:

- ▶ pops its stack until it enters a state where the token 'error' is legal;
- ▶ then behaves as if it saw the token 'error'
 - performs the action encountered;
 - resets the lookahead token to the token that caused the error.
- ▶ If no 'error' rules specified, processing halts.



Controlling Error Behavior

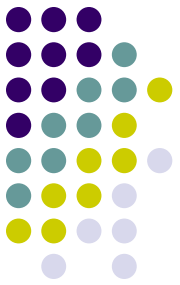
- Parser remains in error state until three tokens successfully read in and shifted
 - ▶ prevents cascaded error messages;
 - ▶ if an error is detected while parser in error state:
 - no error message given;
 - input token causing the error is deleted.
- To force the parser to believe that an error has been fully recovered from:
`yyerrok;`
- To clear the token that caused the error:
`yyclearin;`



Placing 'error' tokens

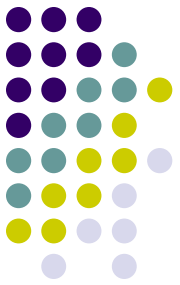
Some guidelines:

- Close to the start symbol of the grammar:
 - ▶ To allow recovery without discarding all input.
- Near terminal symbols:
 - ▶ To allow only a small amount of input to be discarded on an error.
 - ▶ Consider tokens like ')', ';', that follow nonterminals.
- Without introducing conflicts.



Error Messages

- On finding an error, the parser calls a function
`void yyerror(char *s)` /* s points to an error msg */
 - ▶ user-supplied, prints out error message.
- More informative error messages:
 - ▶ `int yychar`: token no. of token causing the error.
 - ▶ user program keeps track of line numbers, as well as any additional info desired.



Adding Semantic Actions

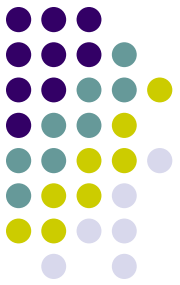
- Semantic actions for a rule are placed in its body:
 - ▶ an action consists of C code enclosed in { ... }
 - ▶ may be placed anywhere in rule RHS

Example:

```
expr : ID { symTbl_lookup(idname); }
```

```
decl : type_name { tval = ... } id_list;
```

Example



```
%{
#include <ctype.h>
%}

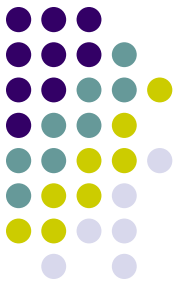
%token DIGIT

%%
line    : expr '\n'          { printf("%d\n", $1); }
        ;
expr    : expr '+' term      { $$ = $1 + $3; }
        | term
        ;
term    : term '*' factor    { $$ = $1 * $3; }
        | factor
        ;
factor  : '(' expr ')'       { $$ = $2; }
        | DIGIT
        ;

%%
```

```
yylex() {
    int c;
    c = getchar();
    if (isdigit(c)) {
        yylval = c - '0';
        return DIGIT;
    }
    return c;
}
```

Example

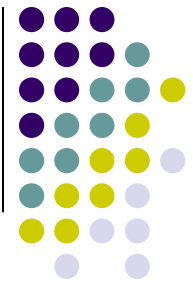


```
%{
#include <ctype.h>
#include <stdio.h>
#define YYSTYPE double /* double type for Yacc stack */
}%
%token NUMBER

%left '+' '-'
%left '*' '/'
%right UMINUS
%%

lines : lines expr '\n'    { printf("%g\n", $2); }
      | lines '\n'
      | /* empty */
      ;
expr  : expr '+' expr      { $$ = $1 + $3; }
      | expr '-' expr      { $$ = $1 - $3; }
      | expr '*' expr      { $$ = $1 * $3; }
      | expr '/' expr      { $$ = $1 / $3; }
      | '(' expr ')'       { $$ = $2; }
      | '-' expr %prec UMINUS { $$ = - $2; }
      | NUMBER
      ;
```

Synthesized Attributes

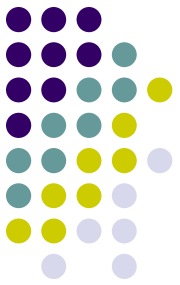


Each nonterminal can “return” a value:

- ▶ The return value for a nonterminal X is “returned” to a rule that has X in its body, e.g.:

$A : \dots X \dots$
 ↑
 value “returned” by X
 $X : \dots$

- ▶ *This is different from the values returned by the scanner to the parser!*



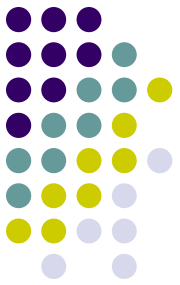
Attribute return values

- To access the value returned by the i^{th} symbol in a rule body, use $\$i$
 - ▶ an action occurring in a rule body counts as a symbol. E.g.:

decl : type { tval = \$1 } id_list { symtbl_install(\$3, tval); }

The diagram illustrates the mapping of symbols to indices in a rule body. Brackets are placed under each symbol: 'type' (1), '{ tval = \$1 }' (2), 'id_list' (3), and 'symtbl_install(\$3, tval);' (4). The index 3 is circled, and a line connects it to the '\$3' symbol within the 'symtbl_install' action, demonstrating how to reference the value of a non-terminal symbol within an action.

- To set the value to be returned by a rule, assign to $$$$
 - ▶ by default, the value of a rule is the value of its first symbol, i.e., \$1.



Declaring Return Value Types

- Default type for nonterminal return values is **int**.
- Need to declare return value types if nonterminal return values can be of other types:

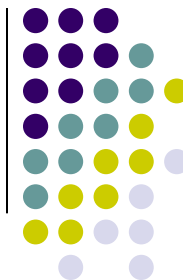
- ▶ Declare the union of the different types that may be returned:

```
%union {  
    symtbl_ptr    st_ptr;  
    id_list_ptr   list_of_ids;  
    tree_node_ptr syntax_tree_ptr;  
    int           value;  
}
```

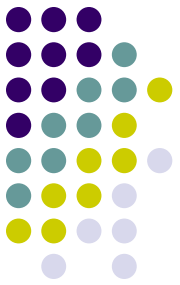
- ▶ Specify which union member a particular grammar symbol will return:

```
%token <value> INTCON, CHARCON; }    terminals  
%type <st_ptr> identifier;           }  
%type <syntax_tree_ptr> expr, stmt; }    nonterminals
```

Conflicts



- A conflict occurs when the parser has multiple possible actions in some state for a given next token.
- Two kinds of conflicts:
 - ▶ *shift-reduce conflict*.
 - The parser can either keep reading more of the input (“shift action”), or it can mimic a derivation step using the input it has read already (“reduce action”).
 - ▶ *reduce-reduce conflict*.
 - There is more than one production that can be used for mimicking a derivation step at that point.



Example of a conflict

Grammar rules:

$S \rightarrow \text{if } (e) S \quad /* 1 */$
 $\quad | \text{if } (e) S \text{ else } S \quad /* 2 */$

Input: **if** (e_1) **if** (e_2) S_2 **else** S_3

Parser state when input token = 'else':

- ▶ Input already seen: **if** (e_1) **if** (e_2) S_2
- ▶ Choices for continuing:

1. keep reading input (“**shift**”):

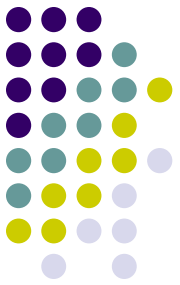
- ‘**else**’ part of innermost if
- eventual parse structure:
 $\text{if } (e_1) \{ \text{if } (e_2) S_2 \text{ else } S_3 \}$

2. mimic derivation step using

$S \rightarrow \text{if } (e) S$ (“**reduce**”):

- ‘**else**’ part of outermost if
- eventual parse structure:
 $\text{if } (e_1) \{ \text{if } (e_2) S_2 \} \text{ else } S_3$

shift-reduce conflict



Handling Conflicts

General approach:

- Iterate as necessary:
 1. Use “yacc -v” to generate the file **y.output**.
 2. Examine **y.output** to find parser states with conflicts.
 3. For each such state, examine the items to figure why the conflict is occurring.
 4. Transform the grammar to eliminate the conflict:

Reason for conflict	Possible grammar transformation
Ambiguity with operators in expressions	Specify associativity, precedence
Error action	Move or eliminate offending error action
Semantic action	Move the offending semantic action
Insufficient lookahead	“expand out” the nonterminal involved
Other	...???...