

Compiler Construction

Chapter 5: Syntax-driven translation

Dittaya Wanvarie

Department of Mathematics and Computer Science
Chulalongkorn University

Second semester, 2024

- 1 Introduction
- 2 Syntax-driven translation
- 3 Attribute grammar
- 4 Modeling the naming environment
- 5 Type systems
- 6 Storage layout

Compiler's tasks

- ➊ Read the source
- ➋ Build intermediate representation: assemble knowledge on the input
 - ▶ Names, types, syntactic structures, etc.
 - ▶ Symbol tables
- ➌ Translate into the target

How to build an IR

- Syntax-driven translation: added actions into a parser
- IR traversal to build another IR

- IR building process
- Name visibility, scope, name binding
- Programming language construct such as variable reference, case statement, heap allocation

We can perform these actions as a side-effect during parsing

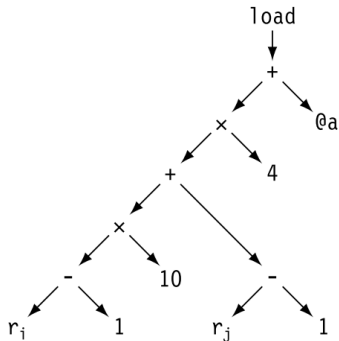
Example: necessary information

Representing $a[i][j]$



Source-Level Tree

As seen in the code text



Low-Level Tree

To generate assembly codes

Reference to x

- Map x to its appropriate runtime instance
 - ▶ Name resolution
- Know type, size, structure, and lifetime of x
 - ▶ E.g. cannot treat a floating point as a memory reference
- Memory location

- 1 Introduction
- 2 Syntax-driven translation
- 3 Attribute grammar
- 4 Modeling the naming environment
- 5 Type systems
- 6 Storage layout

Example: positive integer grammar

Parsing **175** and get its value

1	<i>Number</i>	→	<i>DList</i>
2	<i>DList</i>	→	<i>DList</i> digit
3			digit

(a) The Positive Integer Grammar

State	Action		Goto
	eof	digit	<i>DList</i>
0		s 2	1
1	acc	s 3	
2	r 3	r 3	
3	r 2	r 2	

(b) Its *Action* and *Goto* Tables

■ **FIGURE 5.1** The Grammar for Positive Integers.

Example of side actions

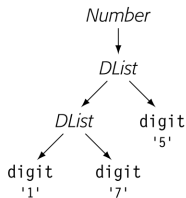
```
Number : DList          { return $1; } ;
DList  : DList digit    { $$ = $1 * 10 + CTToI($2); } ;
        | digit         { $$ = CTToI($1); } ;
```


Example: positive integer grammar

Iteration	State	Word	Stack	Action
0	—	<u>1</u>	\$ <i><Number 0></i>	<i>shift 2</i>
1	2	<u>7</u>	\$ <i><Number 0></i> <i><digit 2></i>	<i>reduce 3</i>
2	1	<u>7</u>	\$ <i><Number 0></i> <i><DList 1></i>	<i>shift 3</i>
3	3	<u>5</u>	\$ <i><Number 0></i> <i><DList 1></i> <i><digit 3></i>	<i>reduce 2</i>
4	1	<u>5</u>	\$ <i><Number 0></i> <i><DList 1></i>	<i>shift 3</i>
5	3	eof	\$ <i><Number 0></i> <i><DList 1></i> <i><digit 3></i>	<i>reduce 2</i>
6	1	eof	\$ <i><Number 0></i> <i><DList 1></i>	<i>accept</i>

■ **FIGURE 5.2** Parser Actions for the Number “175”.

Example: positive integer grammar



(a) Syntax Tree for "175" with the Left-recursive Grammar

Value(root)

if (root is type Number) then
 return Value(son(root))
else if (root is type DList) then
 return $10 \times \text{Value}(\text{left}(\text{root}))$
 + Value(right(root))
else if (root is type digit) then
 return integer(root's lexeme)

(b) Treewalk to Compute the Value from the Syntax Tree

■ **FIGURE 5.3** Treewalk Computations for the Positive Integer Grammar.

Which one is better?

A)

```
Number : DList          { return $1; } ;
DList  : DList digit    { $$ = $1 * 10 + CToI($2); } ;
        | digit         { $$ = CToI($1); } ;
```

B)

```
Number : DList          { return second($1); } ;
DList  : digit DList    { $$ = pair(10 * first($2),
                          first($2) * CToI($1) + second($2)); }
        | digit         { $$ = pair(10, CToI($1)); } ;
```

We may have to rewrite the grammar for simpler and faster computation

Example: construct an AST

Production	Syntax-Driven Actions
$Expr \rightarrow Expr + Term$	{ $$$ \leftarrow MakeNode2(\text{plus}, \$1, \$3);$ };
$Expr - Term$	{ $$$ \leftarrow MakeNode2(\text{minus}, \$1, \$3);$ };
$Term$	{ $$$ \leftarrow \$1;$ };
$Term \rightarrow Term \times Factor$	{ $$$ \leftarrow MakeNode2(\text{times}, \$1, \$3);$ };
$Term \div Factor$	{ $$$ \leftarrow MakeNode2(\text{divide}, \$1, \$3);$ };
$Factor$	{ $$$ \leftarrow \$1;$ };
$Factor \rightarrow (Expr)$	{ $$$ \leftarrow \$2;$ };
$number$	{ $$$ \leftarrow MakeLeaf(\text{number}, \text{lexeme});$ };
$name$	{ $$$ \leftarrow MakeLeaf(\text{name}, \text{lexeme});$ };

■ FIGURE 5.4 Building an Abstract Syntax Tree.

Example: $a - 2 \times b$

Example: translate to 3-address code

Example: $a - 2 \times b$

Production	Syntax-Driven Actions
$Expr \rightarrow Expr + Term$	{ $$$ \leftarrow NextRegister()$; $Emit(add, \$1, \$3, $$)$; };
$Expr - Term$	{ $$$ \leftarrow NextRegister()$; $Emit(sub, \$1, \$3, $$)$; };
$Term$	{ $$$ \leftarrow \1 ; };
$Term \rightarrow Term \times Factor$	{ $$$ \leftarrow NextRegister()$; $Emit(mult, \$1, \$3, $$)$; };
$Term \div Factor$	{ $$$ \leftarrow NextRegister()$; $Emit(div, \$1, \$3, $$)$; };
$Factor$	{ $$$ \leftarrow \1 ; };
$Factor \rightarrow (Expr)$	{ $$$ \leftarrow \2 ; };
number	{ $$$ \leftarrow NumberIntoReg(lexeme)$; };
name	{ $entry \leftarrow STLookup(lexeme)$; $$$ \leftarrow ValueIntoReg(entry)$; };

loadAI $r_{arp}, @a \Rightarrow r_1$

loadI 2 $\Rightarrow r_2$

loadAI $r_{arp}, @b \Rightarrow r_3$

mult $r_2, r_3 \Rightarrow r_4$

sub $r_1, r_4 \Rightarrow r_5$

ILOC Code for $a - 2 \times b$

■ **FIGURE 5.5** Emitting Three-Address Code for Expressions.

```
push  $\langle \text{INVALID}, \text{INVALID}, \text{INVALID} \rangle$  onto the stack
push  $\langle \text{start symbol}, s_0, \text{INVALID} \rangle$  onto the stack
word  $\leftarrow \text{NextWord}()$ 
while (true) do
    state  $\leftarrow$  state from triple at top of stack
    if Action[state,word] = "reduce  $A \rightarrow \beta$ " then
        value  $\leftarrow \text{PerformActions}(A \rightarrow \beta)$ 
        pop  $|\beta|$  triples from the stack
        state  $\leftarrow$  state from triple at top of stack
        push  $\langle A, \text{Goto}[\text{state}, A], \text{value} \rangle$  onto the stack
    else if Action[state,word] = "shift  $s_i$ " then
        push  $\langle \text{word}, s_i, \text{lexeme} \rangle$  onto the stack
        word  $\leftarrow \text{NextWord}()$ 
    else if Action[state,word] = "accept" and word = eof
        then break
    else throw a syntax error
report success /* executed the "accept" case */
```

■ **FIGURE 5.6** The Skeleton LR(1) Parser with Translation Support.

What if the information is not in the syntax tree?

- E.g. type, lifetime, visibility

Solution

- 1 Use global variable, **CurType**

```
Declaration → TypeSpec NameList { CurType ← invalid; };
TypeSpec    → int                { CurType ← int; };
              | float             { CurType ← float; };
NameList     → NameList , name    { err ← SetType($2, CurType); };
              | name              { err ← SetType($1, CurType); };
```

- 2 Change grammar, avoid overriding **CurType**

```
Declaration → int INameList
              | float FNameList
INameList   → NameList name { err ← SetType($2, int); };
              | name        { err ← SetType($1, int); };
FNameList   → NameList name { err ← SetType($2, float); };
              | name        { err ← SetType($1, float); };
```

a x 2

- 1 Load constant to a register, then perform multiplication

loadI 2 $\Rightarrow r_i$

mult $r_a, r_i \Rightarrow r_j$

ILOC Code for $a \times 2$

- 2 Use `multI`

- ▶ Current grammar:

★ $Term \rightarrow Term \times Factor$

★ $Factor \rightarrow \text{number}$

- ▶ Changed grammar:

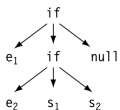
★ $Term \rightarrow Term \times \text{number}$

★ $Term \rightarrow Term \times \text{name}$

However, the new grammar does not support $2 \times a$. One better solution is to delay the optimization to later stages.

If structure

if e_1 then if e_2 then s_1 else s_2



AST for Nested If Statements

- 1 Evaluation the conditional expression (e_1)
- 2 Branch to s_1 section if the condition is true, s_2 otherwise. We need the labels of the first command in s_1 and s_2 parts.
- 3 At the end of each subpart (s_1 and s_2), jump to the next statement (*exit*). We also need the label for this statement.

From the if-else production rules

1	<i>Stmt</i>	→	if <i>Expr</i> then <i>Stmt</i>
2			if <i>Expr</i> then <i>WithElse</i> else <i>Stmt</i>
3			<i>Other</i>
4	<i>WithElse</i>	→	if <i>Expr</i> then <i>WithElse</i> else <i>WithElse</i>
5			<i>Other</i>

■ **FIGURE 5.7** The Unambiguous If-Then-Else Grammar.

We need to create the labels s_1 , s_2 , *exit* when we start working on the statement. We have to pass this label information to each subpart.

We may construct an object for each information and embed it to the variable.

- However, actions for each piece of information occur at different time.

Solution

- Then, we construct an empty string production for the actions.

WithElse \rightarrow if *Expr* *CreateBranch*
 then *WithElse ToExit1*
 else *WithElse ToExit2*

CreateBranch \rightarrow ϵ

ToExit1 \rightarrow ϵ

ToExit2 \rightarrow ϵ

■ **FIGURE 5.8** Creating Mid-Production Actions.

- 1 Introduction
- 2 Syntax-driven translation
- 3 Attribute grammar
- 4 Modeling the naming environment
- 5 Type systems
- 6 Storage layout

For a context-free grammar

- Each symbol associates with one or more attributes
- Each production is augmented with attribute computations
- There is no evaluation order between computations for the same production

Evaluation order? Dependency?

- We may evaluate the attribute during parsing if all required values are available when we construct the node

Synthesized attribute

- An attribute defined wholly in terms of the attributes of the node, **its children**, and constants
- We can compute this attribute during **bottom-up** parsing

Inherited attribute

- An attribute defined wholly in terms of the node's own attributes and those of its **siblings** or **its parent** in the parse tree, and constants
- We can compute this attribute during **top-down** parsing

Naming values

YACC notation

- $\$ \$$ refers to the result location for the current production, i.e. the left-hand-side symbol.
- $\$ 1, \$ 2, \dots, \$ n$ refer to the locations for the first, second, through the $\$ n$ symbol on the right-hand side.

Object-oriented notation

- Simply refer to a value as an attribute of the variable.
- E.g. *Number.value*

Exercise: Signed binary numbers

Set the value of *Number*

Production		
1	<i>Number</i>	→ <i>Sign List</i>
2	<i>Sign</i>	→ +
3	<i>Sign</i>	→ -
4	<i>List</i>	→ <i>Bit</i>
5	<i>List</i> ₀	→ <i>List</i> ₁ <i>Bit</i>
6	<i>Bit</i>	→ 0
7	<i>Bit</i>	→ 1

Code snippet

- 1 Introduction
- 2 Syntax-driven translation
- 3 Attribute grammar
- 4 Modeling the naming environment**
- 5 Type systems
- 6 Storage layout

Name resolution

- Static binding: determine the name-to-entity binding at the **compile time**
- Dynamic binding: defer the name-to-entity binding until **runtime**

Rule

At a point p in a program, an occurrence of name n refers to the entity named n that was created, explicitly or implicitly, in the scope that is lexically closest to p .

Static coordinate

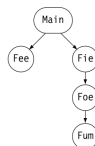
```

program Main0(input0, output0);
  var x1, y1, z1: integer;
  procedure Fee1;
    var x2: integer;
    begin { Fee1 }
      x2 := 1;
      y1 := x2 * 2 + 1
    end;
  procedure Fie1;
    var y2: real;
    procedure Foe2;
      var z3: real;
      procedure Fum3;
        var y4: real;
        begin { Fum3 }
          x1 := 1.25 * z3;
          Fee1;
          writeln('x = ', x1)
        end;
      begin { Foe2 }
        z3 := 1;
        Fee1;
        Fum3
      end;
    end;
  begin { Fie1 }
    Foe2;
    writeln('x = ', x1)
  end;
begin { Main0 }
  x1 := 0;
  Fie1
end.
    
```

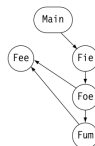
(a) PASCAL Program

Scope	x	y	z
Main	(1,0)	(1,4)	(1,8)
Fee	(2,0)	(1,4)	(1,8)
Fie	(1,0)	(2,0)	(1,8)
Foe	(1,0)	(2,0)	(3,0)
Fum	(1,0)	(4,0)	(3,0)

(b) Static Coordinates

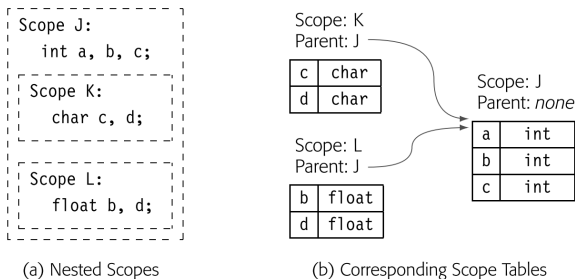


(c) Nesting Relationships



(d) Calling Relationships

■ FIGURE 5.9 Nested Lexical Scopes in PASCAL.



■ **FIGURE 5.10** Tables for the Lexical Hierarchy.

- Chains of symbol tables
- A search path

Via syntax-driven translation

- When the source language constructs enter and leave distinct scopes
 - ▶ Create tables, search paths

Examples

- Block demarcations
 - ▶ On entry: Create a new table and link it to the current scope
 - ▶ On exit: Mark the table as final (not in use)
- Variable declarations: Create entries in a local table
 - ▶ Existing variables: Populate attributes
 - ▶ New variable: Infer attributes
- References: search for information in the local table, and the entire search path

C

- Procedure names and global variables exist in the global scope
- Each procedure creates its own local scope for variables, parameters, and labels
- Blocks ({ and }), create their own local scopes
- **static** name has a global lifetime
 - ▶ **static** global name is visible inside the file, while **static** local name is visible locally

Python

- 3 kinds of scopes
 - ▶ A local function-specific scope
 - ▶ A global scope
 - ▶ A builtin scope e.g. **print**
 - ▶ The first use of **x** defines the scope of **x**
 - ★ If it is an assignment, **x** is in a local scope
 - ★ If it is a use, **x** is in a global scope

- A subclass will **inherit** all methods and members from its superclass
- Single inheritance vs multiple inheritance
- Each class definition has its own scope, with links to its superclass
- Overriding
- The compiler must map an <object, member> pair back to a specific definition

- Links to its superclass(es)
- Compile-time vs. runtime resolution
- Lookup with inheritance or search paths
- Construct scope tables using syntax-driven actions

Tables for the inheritance hierarchy

```
Class Point {  
    public int x, y;  
    private int z;  
    public void draw() {...};  
    public void move() {...};  
}  
  
Class ColorPoint extends Point {  
    private Color c;  
    public void draw() {...};  
    public void setc( Color x )  
        { this.c = x };  
}
```

(a) Class Definitions

Class: Point
Superclass: *none*

x	int	public
y	int	public
z	int	private
draw	void()	public
move	void()	public

Class: ColorPoint
Superclass: Point

c	Color	private
draw	void()	public
setc	void()	public

(b) Corresponding Scope Tables

■ **FIGURE 5.11** Tables for the Inheritance Hierarchy.

For example,

- `public`
- `private`
- `protected`
- `default`

A typical implementation will include a visibility tag in the symbol table record of each name

- 1 Introduction
- 2 Syntax-driven translation
- 3 Attribute grammar
- 4 Modeling the naming environment
- 5 Type systems**
- 6 Storage layout

Fundamental properties of a name

- Storage size
- Range
- Low-level representation
- Etc.

Behaviors of the program depends on data types (overloading)

- Conformable: the operators and arguments (operands) are well defined
- Efficient translation
- Casting
- Information for garbage collection

- Type signature
- Function prototype

```
if (tag(a) = tag(b)) then // take the short path
  switch (tag(a)) into {
    case SHORT: // use SHORT add
      value(c) ← value(a) + value(b)
      tag(c) ← SHORT
      break
    case INTEGER: // use INTEGER add
      value(c) ← value(a) + value(b)
      tag(c) ← INTEGER
      break
    case LONG INTEGER: // use LONG INTEGER add
      value(c) ← value(a) + value(b)
      tag(c) ← LONG INTEGER
      break
  }
else // take the long path
  (c, tag(c)) ← AddMixedTypes(a, tag(a), b, tag(b))
```

■ **FIGURE 5.12** Integer Addition with Runtime Type Checking.

Base types

- Integer, floating point, string
- Multiple size such as byte, word, double word, quadruple word

Compound and constructed types

- Graphs, trees, tables, lists, stacks, maps
- Arrays
- Strings
- Enumerated types
- Structures and variants
- Objects and classes

- Name equivalence
- Structural equivalence

```
struct Tree {  
    int value;  
    struct Tree *left;  
    struct Tree *right;  
}  
  
struct Bush {  
    int value;  
    struct Bush *left;  
    struct Bush *right;  
}
```

Production	Syntax-Driven Action
$Expr \rightarrow Expr + Term$	{ set_type(\$\$, \mathcal{F}_+ (type(\$1),type(\$3))); };
$Expr - Term$	{ set_type(\$\$, \mathcal{F}_- (type(\$1),type(\$3))); };
$Term$	{ set_type(\$\$,type(\$1)); };
$Term \rightarrow Term \times Factor$	{ set_type(\$\$, \mathcal{F}_\times (type(\$1),type(\$3))); };
$Term \div Factor$	{ set_type(\$\$, \mathcal{F}_\div (type(\$1),type(\$3))); };
$Factor$	{ set_type(\$\$,type(\$1)); };
$Factor \rightarrow (Expr)$	{ set_type(\$\$,type(\$2)); };
num	{ set_type(\$\$,type(num)); };
name	{ set_type(\$\$,type(name)); };

■ **FIGURE 5.13** Framework to Assign Types to Subexpressions.

- 1 Introduction
- 2 Syntax-driven translation
- 3 Attribute grammar
- 4 Modeling the naming environment
- 5 Type systems
- 6 Storage layout

Automatic variables

- Identical lifetime of its declaring scope
- Stored in the local data area

Activation record

- A region of memory for an invocation of a procedure
- Activation record pointer is normally a register pointing to the first address of AR

Static variables

- The lifetime is the same of the program
- Single storage choice with visibility tag or individual data area for each variable

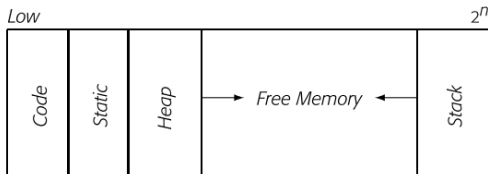
Irregular entities

- Dynamic allocation
- Stored in runtime heap

Temporary values

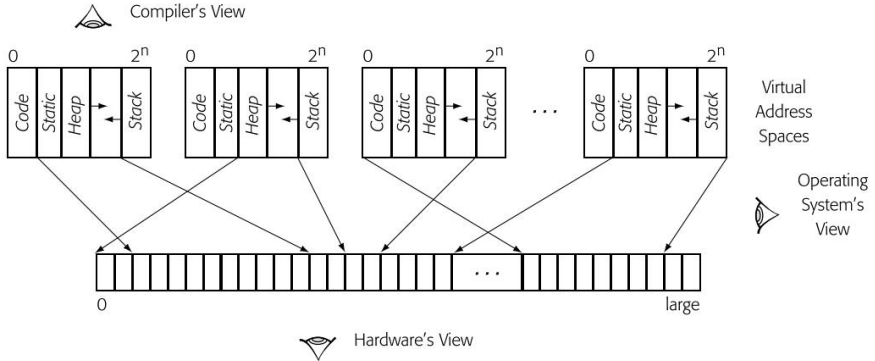
- Usually stored in registers
- Local data area if the size is known
- Heap if the size is unknown

Layout within a virtual address space



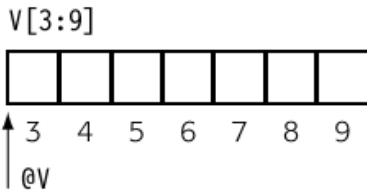
■ **FIGURE 5.14** Virtual Address-Space Layout.

Different views of the address space



■ **FIGURE 5.15** Different Views of the Address Space.

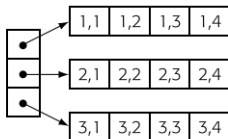
Internal layout for arrays



Vector Layout

1,1	1,2	1,3	1,4
2,1	2,2	2,3	2,4
3,1	3,2	3,3	3,4

(a) 3×4 Array



(c) Indirection Vectors

1,1	1,2	1,3	1,4	2,1	2,2	2,3	2,4	3,1	3,2	3,3	3,4
-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----

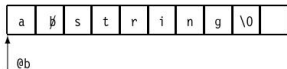
(b) Row-Major Order

1,1	2,1	3,1	1,2	2,2	3,2	1,3	2,3	3,3	1,4	2,4	3,4
-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----

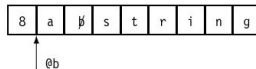
(d) Column-Major Order

■ **FIGURE 5.16** Two-Dimensional Array Layouts.

Internal layout for strings

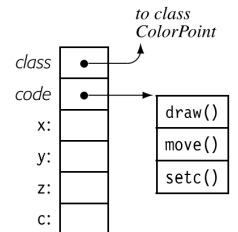


String with Null Termination

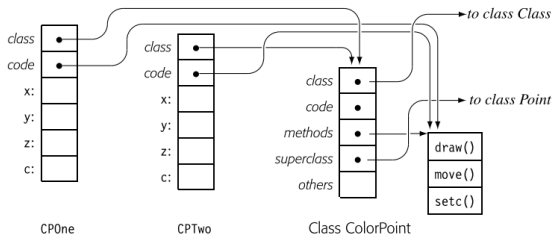


String with Explicit Length Field

Internal layout for structures

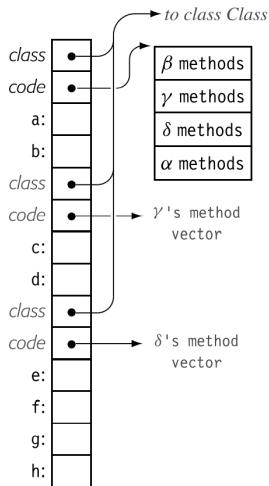


Object Record for Instance
Of Class ColorPoint



■ FIGURE 5.17 Multiple Instances of Class ColorPoint.

Object record for multiple inheritance



Object Record for α

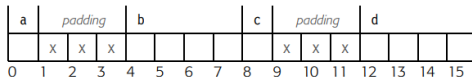
- Declaration before any executable statement appears
- Build up the symbol table during parsing and perform type inference on IR
- Build an IR with abstract references and perform type inference on the IR

Alignment restrictions and padding

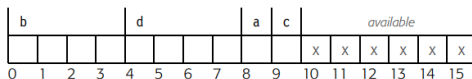
Name	Bytes	Constraint
------	-------	------------

a	1	$@a \bmod 1 = 0$
b	4	$@b \bmod 4 = 0$
c	1	$@c \bmod 1 = 0$
d	4	$@d \bmod 4 = 0$

(a) Variables and Their Alignments



(b) A Layout That Wastes Space



(c) A Better Layout

■ **FIGURE 5.18** Alignment Issues in Data-Area Offset Assignment.