# 26

# Code-Tuning Techniques

## Contents

## Related Topics

CODE TUNING HAS BEEN a popular topic during most of the history of computer programming. Consequently, once you've decided that you need to improve performance and that you want to do it at the code level, you have a rich set of techniques at your disposal.

This chapter focuses on improving speed and includes a few tips for making code smaller. Performance usually refers to both speed and size, but size reductions tend to come more from redesigning classes and data than from tuning code. Code tuning refers to small-scale changes rather than changes in larger-scale designs.

Few of the techniques in this chapter are so generally applicable that you'll be able to copy the example code directly into your programs. The main purpose of the discussion here is to illustrate a handful of code tunings that you can adapt to your situation.

The code-tuning changes described in this chapter might seem cosmetically similar to the refactorings described in Chapter 24. But refactorings are changes that improve a program's internal structure (Fowler 1999). The changes in this chapter might better be called "anti-refactorings." Far from "improving the

31
32
33
34

internal structure," these changes degrade the internal structure in exchange for gains in performance. This is true by definition. If they didn't degrade the internal structure, we wouldn't consider them to be optimizations; we would use them by default and consider them to be standard coding practice.

35 **CROSS-REFERENCE** Cod
36 e tunings are heuristics. For
37 more on heuristics, see
38 Section 5.3, "Design
   Building Blocks: Heuristics."
39
40
41

Some books present code tuning techniques as "rules of thumb" or cite research that suggests that a specific tuning will produce the desired effect. As you will soon see, the concept of "rules of thumb" applies poorly to code tuning. The only reliable rule of thumb is to measure the effect of each tuning in your environment. Thus this chapter presents a catalog of "things to try"—many of which won't work in your environment but some of which will work very well indeed.

42

# 26.1 Logic

43 **CROSS-REFERENCE** For
44 other details on using
   statement logic, see Chapters
   14 through 19.
45

Much of programming consists of manipulating logic. This section describes how to manipulate logical expressions to your advantage.

## Stop Testing When You Know the Answer

46

Suppose you have a statement like

47
48
49

```
if ( 5 < x ) and ( x < 10 ) then ...
```
Once you've determined that $x$ is less than 5, you don't need to perform the second half of the test.

50 **CROSS-REFERENCE** For
51 more on short-circuit
52 evaluation, see "Knowing
   How Boolean Expressions
53 Are Evaluated" in "Knowing
   How Boolean Expressions
54 Are Evaluated" in Section
55 19.1.
56

Some languages provide a form of expression evaluation known as "short-circuit evaluation," which means that the compiler generates code that automatically stops testing as soon as it knows the answer. Short-circuit evaluation is part of C++'s standard operators and Java's "conditional" operators.

If your language doesn't support short-circuit evaluation natively, you have to avoid using *and* and *or*, adding logic instead. With short-circuit evaluation, the code above changes to this:

57
58
59
60
61
62
63
64

```
if ( 5 < x ) then
   if ( x < 10 ) then ...
```
The principle of not testing after you know the answer is a good one for many other kinds of cases as well. A search loop is a common case. If you're scanning an array of input numbers for a negative value and you simply need to know whether a negative value is present, one approach is to check every value, setting a *negativeFound* variable when you find one. Here's how the search loop would look:

**C++ Example of Not Stopping After You Know the Answer**

```
negativeInputFound = False;
for ( i = 0; i < iCount; i++ ) {
    if ( input[ i ] < 0 ) {
        negativeInputFound = True;
    }
}
```

A better approach would be to stop scanning as soon as you find a negative value. Here are the approaches you could use to solve the problem:

- Add a *break* statement after the *negativeInputFound = True* line.

- If your language doesn't have *break*, emulate a *break* with a *goto* that goes to the first statement after the loop.

- Change the *for* loop to a *while* loop and check for *negativeInputFound* as well as for incrementing the loop counter past *iCount*.

- Change the *for* loop to a *while* loop, put a sentinel value in the first array element after the last value entry, and simply check for a negative value in the *while* test. After the loop terminates, see whether the position of the first found value is in the array or one past the end. Sentinels are discussed in more detail later in the chapter.

Here are the results of using the *break* keyword in C++ and Java:

| Language | Straight Time | Code-Tuned Time | Time Savings |
|---|---|---|---|
| **C++** | **4.27** | **3.68** | **14%** |
| Java | 4.85 | 3.46 | 29% |

*Note: (1) Times in these tables are given in seconds and are meaningful only for comparisons across rows of each table. Actual times will vary according to the compiler and compiler options used and the environment in which each test is run. (2) Benchmark results are typically made up of several thousand to many million executions of the code fragments to smooth out the sample-to-sample fluctuations in the results. (3) Specific brands and versions of compilers aren't indicated. Performance characteristics vary significantly from brand to brand and version to version. (4) Comparisons among results from different languages aren't always meaningful because compilers for different languages don't always offer comparable code-generation options. (5) The results shown for interpreted languages (PHP and Python) are typically based on less than 1% of the test runs used for the other languages. (6) Some of the "time savings" percentages  might not be exactly reproducible from the data in these tables due to rounding of the "straight time" and "code-tuned time" entries.*

99
100
101
102

The impact of this change varies a great deal depending on how many values you have and how often you expect to find a negative value. This test assumed an average of 100 values and assumed that a negative value would be found 50 percent of the time.

103

# Order Tests by Frequency

104
105
106
107

Arrange tests so that the one that's fastest and most likely to be true is performed first. It should be easy to drop through the normal case, and if there are inefficiencies, they should be in processing the uncommon cases. This principle applies to *case* statements and to chains of *if-then-else*s.

108
109

Here's a *Select-Case* statement that responds to keyboard input in a word processor:

110

**Visual Basic Example of a Poorly Ordered Logical Test**

111
112
113
114
115
116
117
118
119
120
121
122
123
124

```
Select inputCharacter
   Case "+", "="
      ProcessMathSymbol( inputCharacter )
   Case "0" To "9"
      ProcessDigit( inputCharacter )
   Case ",", ".", ":", ";", "!", "?"
      ProcessPunctuation( inputCharacter )
   Case " "
      ProcessSpace( inputCharacter )
   Case "A" To "Z", "a" To "z"
      ProcessAlpha( inputCharacter )
   Case Else
      ProcessError( inputCharacter )
End Select
```

125
126
127
128
129
130
131

The cases in this *case* statement are ordered in something close to the ASCII sort order. In a *case* statement, however, the effect is often the same as if you had written a big set of *if-then-elses*, so if you get an *<;$QS>a<;$QS>* as an input character, the program tests whether it's a math symbol, a punctuation mark, a digit, or a space before determining that it's an alphabetic character. If you know the likely frequency of your input characters, you can put the most common cases first. Here's the reordered *case* statement:

132

**Visual Basic Example of a Well-Ordered Logical Test**

133
134
135
136
137

```
Select inputCharacter
   Case "A" To "Z", "a" To "z"
      ProcessAlpha( inputCharacter )
   Case " "
      ProcessSpace( inputCharacter )
```

```
138        Case ",", ".", ":", ";", "!", "?"
139           ProcessPunctuation( inputCharacter )
140        Case "0" To "9"
141           ProcessDigit( inputCharacter )
142        Case "+", "="
143           ProcessMathSymbol( inputCharacter )
144        Case Else
145           ProcessError( inputCharacter )
146     End Select
```

147    Since the most common case is usually found sooner in the optimized code, the
148    net effect will be the performance of fewer tests. Here are the results of this
149    optimization with a typical mix of characters:

| Language | Straight Time | Code-Tuned Time | Time Savings |
|---|---|---|---|
| C# | 0.220 | 0.260 | -18% |
| Java | 2.56 | 2.56 | 0% |
| **Visual Basic** | **0.280** | **0.260** | **7%** |

150    *Note: Benchmarked with an input mix of 78 percent alphabetic characters, 17*
151    *percent spaces, and 5 percent punctuation symbols.*

152    The Visual Basic results are as expected, but the Java and C# results are not as
153    expected. Apparently that's because of the way *switch-case* statements are
154    structured in C++ and Java—since each value must be enumerated individually
155    rather than in ranges, the C++ and Java code doesn't benefit from the
156    optimization as the Visual Basic code does. This result underscores the
157    importance of not following any optimization advice blindly—specific compiler
158    implementations will significantly affect the results.

159    You might assume that the code generated by the Visual Basic compiler for a set
160    of *if-then-else*s that perform the same test as the *case* statement would be similar.
161    Here are those results:

| Language | Straight Time | Code-Tuned Time | Time Savings |
|---|---|---|---|
| C# | 0.630 | 0.330 | 48% |
| Java | 0.922 | 0.460 | 50% |
| **Visual Basic** | **1.36** | **1.00** | **26%** |

162    The results are quite different. For the same number of tests, the VB compiler
163    takes about 5 times as long in the unoptimized case, 4 times in the optimized
164    case. This suggests that the compiler is generating different code for the *case*
165    approach than for the *if-then-else* approach.

166
167
168
169

The improvement with *if-then-else*s is more consistent than it was with the *case* statements, but that's a mixed blessing. In C# and VB both versions of the *case* statement approach are faster than both versions of the *if-then-else* approach, whereas in Java both versions are slower.

170
171

This variation in results suggests a third possible optimization, described in the next section.

172

## Compare Performance of Similar Logic Structures

173
174
175
176

The test described above could be performed using either a *case* statement or *if-then-else*s. Depending on the environment, either approach might work better. Here is the data from the preceding two tables reformatted to present the "code-tuned" times comparing *if-then-else* and *case* performance:

| Language | *case* | *if-then-else* | Time Savings | Performance Ratio |
|---|---|---|---|---|
| C# | 0.260 | 0.330 | -27% | 1:1 |
| Java | 2.56 | 0.460 | 82% | 6:1 |
| Visual Basic | 0.260 | 1.00 | 258% | 1:4 |

177
178
179
180
181

These results defy any logical explanation. In one of the languages, *case* is dramatically superior to *if-then-else*, and in another, *if-then-else* is dramatically superior to *case*. In the third language, the difference is relatively small. You might think that because C# and Java share similar syntax for *case* statements, their results would be similar, but in fact their results are opposite each other.

182
183
184

This example clearly illustrates the difficulty of performing any sort of "rule of thumb" or "logic" to code tuning—there is simply no reliable substitute for *measuring* results.

185
186

## Substitute Table Lookups for Complicated Expressions

187
188
189
190
191

**CROSS-REFERENCE**  For details on using table lookups to replace complicated logic, see Chapter 18, "Table-Driven Methods."
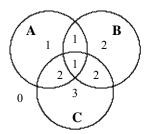
In some circumstances, a table lookup may be quicker than traversing a complicated chain of logic. The point of a complicated chain is usually to categorize something and then to take an action based on its category. As an abstract example, suppose you want to assign a category number to something based on which of Groups *A*, *B*, and *C* it falls into:

192

### G26xx01

194 Here's an example of the complicated logic chain that assigns the category
195 numbers:

196 **C++ Example of a Complicated Chain of Logic**

```
if ( ( a && !c ) || ( a && b && c ) ) {
   category = 1;
}
else if ( ( b && !a ) || ( a && c && !b ) ) {
   category = 2;
}
else if ( c && !a && !b ) {
   category = 3;
}
else {
   category = 0;
}
```

209 You can replace this test with a more modifiable and higher-performance lookup
210 table. Here's how:

211 **C++ Example of Using a Table Lookup to Replace Complicated Logic**

*This table definition is somewhat difficult to understand. Any commenting you can do to make table definitions readable helps.*

```
// define categoryTable
static int categoryTable[ 2 ][ 2 ][ 2 ] = {
   // !b!c   !bc   b!c   bc
      0,    3,    2,    2,   //   !a
      1,    2,    1,    1    //   a
};
...

category = categoryTable[ a ][ b ][ c ];
```

221 Although the definition of the table is hard to read, if it's well documented it
222 won't be any harder to read than the code for the complicated chain of logic was.
223 If the definition changes, the table will be much easier to maintain than the
224 earlier logic would have been. Here are the performance results:

| Language | Straight Time | Code-Tuned Time | Time Savings | Performance Ratio |
|----------|---------------|-----------------|--------------|-------------------|
| **C++** | **5.04** | **3.39** | **33%** | **1.5:1** |
| Visual Basic | 5.21 | 2.60 | 50% | 2:1 |

## Use Lazy Evaluation

225

226  One of my former roommates was a great procrastinator. He justified his laziness
227  by saying that many of the things people feel rushed to do simply don't need to
228  be done. If he waited long enough, he claimed, the things that weren't important
229  would be procrastinated into oblivion, and he wouldn't waste his time doing
230  them.

231  Lazy evaluation is based on the principle my roommate used. If a program uses
232  lazy evaluation, it avoids doing any work until the work is needed. Lazy
233  evaluation is similar to just-in-time strategies that do the work closest to when
234  it's needed.

235  Suppose, for example, that your program contains a table of 5000 values,
236  generates the whole table at startup time, and then uses it as the program
237  executes. If the program uses only a small percentage of the entries in the table,
238  it might make more sense to compute them as they're needed rather than all at
239  once. Once an entry is computed, it can still be stored for future reference
240  ("cached").

# 26.2 Loops

241

242 **CROSS-REFERENCE**  For
243 other details on loops, see
     Chapter 16, "Controlling
     Loops."

Because loops are executed many times, the hot spots in a program are often
inside loops. The techniques in this section make the loop itself faster.

## Unswitching

244

245  Switching refers to making a decision inside a loop every time it's executed. If
246  the decision doesn't change while the loop is executing, you can unswitch the
247  loop by making the decision outside the loop. Usually this requires turning the
248  loop inside out, putting loops inside the conditional rather than putting the
249  conditional inside the loop. Here's an example of a loop before unswitching:

**CODING HORROR**

**C++ Example of a Switched Loop**

```
for ( i = 0; i < count; i++ ) {
   if ( sumType == SUMTYPE_NET ) {
      netSum = netSum + amount[ i ];
   }
   else {
      grossSum = grossSum + amount[ i ];
   }
}
```

In this code, the test *if ( sumType == SUMTYPE_NET )* is repeated through each iteration even though it'll be the same each time through the loop. You can rewrite the code for a speed gain this way:

**C++ Example of an Unswitched Loop**

```
if ( sumType == SUMTYPE_NET ) {
   for ( i = 0; i < count; i++ ) {
      netSum = netSum + amount[ i ];
   }
}
else {
   for ( i = 0; i < count; i++ ) {
      grossSum = grossSum + amount[ i ];
   }
}
```

This is good for about a 20 percent time savings:

| Language | Straight Time | Code-Tuned Time | Time Savings |
|---|---|---|---|
| **C++** | **2.81** | **2.27** | **19%** |
| Java | 3.97 | 3.12 | 21% |
| Visual Basic | 2.78 | 2.77 | <1% |
| Python | 8.14 | 5.87 | 28% |

A hazard distinct to this case is that the two loops have to be maintained in parallel. If *count* changes to *clientCount*, you have to remember to change it in both places, which is an annoyance for you and a maintenance headache for anyone else who has to work with the code.

This example also illustrates a key challenge in code tuning—the effect of any specific code tuning is not predictable. The code tuning produced significant improvements in three of the four languages, but not in Visual Basic. To perform this specific optimization in this specific version of VB would produce less maintainable code without any offsetting gain in performance. The general

283 lesson is that you must measure the effect of each specific optimization to be
284 sure of its effect—no exceptions.

## Jamming

286 Jamming, or "fusion," is the result of combining two loops that operate on the
287 same set of elements. The gain lies in cutting the loop overhead from two loops
288 to one. Here's a candidate for loop jamming:

289 **Visual Basic Example of Separate Loops That Could Be Jammed**

```
For i = 0 to employeeCount - 1
   employeeName( i ) = ""
Next
...
For i = 0 to employeeCount - 1
   employeeEarnings( i ) = 0
Next
```

297 When you jam loops, you find code in two loops that you can combine into one.
298 Usually, that means the loop counters have to be the same. In this example, both
299 loops run from *0* to *employeeCount - 1*, so you can jam them. Here's how:

**CODING HORROR**

300 **Visual Basic Example of a Jammed Loop**

```
For i = 0 to employeeCount - 1
   employeeName( i ) = ""
   employeeEarnings( i ) = 0
Next
```

305 Here are the savings:

| Language | Straight Time | Code-Tuned Time | Time Savings |
|---|---|---|---|
| C++ | 3.68 | 2.65 | 28% |
| PHP | 3.97 | 2.42 | 32% |
| **Visual Basic** | **3.75** | **3.56** | **4%** |

306 *Note: Benchmarked for the case in which* employeeCount *equals* 100.

307 As before, the results vary significantly among languages.

308 Loop jamming has two main hazards. First, the indexes for the two parts that
309 have been jammed might change so that they're no longer compatible. Second,
310 you might not be able to combine the loops easily. Before you combine the
311 loops, make sure they'll still be in the right order with respect to the rest of the
312 code.

313 # Unrolling

314 The goal of loop unrolling is to reduce the amount of loop housekeeping. In
315 Chapter 25, a loop was completely unrolled, and 10 lines of code were shown to
316 be faster than 3. In that case, the loop that went from 3 to 10 lines was unrolled
317 so that all 10 array accesses were done individually.

318 Although completely unrolling a loop is a fast solution and works well when
319 you're dealing with a small number of elements, it's not practical when you have
320 a large number of elements or when you don't know in advance how many
321 elements you'll have. Here's an example of a general loop:

322 **Java Example of a Loop That Can Be Unrolled**

323 *Normally, you'd probably use*
324 *a* for *loop for a job like this,*
325 *but to optimize, you'd have to*
326 *convert to a* while *loop. For*
327 *clarity, a* while *loop is shown*
328 *here.*

```java
i = 0;
while ( i < count ) {
    a[ i ] = i;
    i = i + 1;
}
```

328 To unroll the loop partially, you handle two or more cases in each pass through
329 the loop instead of one. This unrolling hurts readability but doesn't hurt the
330 generality of the loop. Here's the loop unrolled once:

331 **CODING HORROR**

331 **Java Example of a Loop That's Been Unrolled Once**

```java
i = 0;
while ( i < count - 1 ) {
    a[ i ] = i;
    a[ i + 1 ] = i + 1;
    i = i + 2;
}
```

339 *These lines pick up the case*
340 *that might fall through the*
341 *cracks if the loop went by*
342 *twos instead of by ones.*

```java
if ( i == count ) {
    a[ count - 1 ] = count - 1;
}
```

342 The technique replaced the original *a[ i ] = i* line with two lines, and *i* is
343 incremented by *2* rather than by *1*. The extra code after the *while* loop is needed
344 when *count* is odd and the loop has one iteration left after the loop terminates.

345 When five lines of straightforward code expand to nine lines of tricky code, the
346 code becomes harder to read and maintain. Except for the gain in speed, its
347 quality is poor. Part of any design discipline, however, is making necessary
348 trade-offs. So, even though a particular technique generally represents poor
349 coding practice, specific circumstances may make it the best one to use.

350 Here are the results of unrolling the loop:

| Language | Straight Time | Code-Tuned Time | Time Savings |
|---|---|---|---|
| C++ | 1.75 | 1.15 | 34% |
| **Java** | **1.01** | **0.581** | **43%** |
| PHP | 5.33 | 4.49 | 16% |
| Python | 2.51 | 3.21 | -27% |

351    *Note: Benchmarked for the case in which* count *equals* 100.

352    A gain of 16 to 43 percent is respectable, although again you have to watch out
353    for hurting performance, as the Python benchmark shows. The main hazard of
354    loop unrolling is an off-by-one error in the code after the loop that picks up the
355    last case.

356    What if you unroll the loop even further, going for two or more unrollings? Do
357    you get more benefit? Here's the code for a loop unrolled twice:

358    **CODING HORROR**

**Java Example of a Loop That's Been Unrolled Twice**

```
i = 0;
while ( i < count - 2 ) {
   a[ i ] = i;
   a[ i + 1 ] = i+1;
   a[ i + 2 ] = i+2;
   i = i + 3;
}
if ( i <= count - 1 ) {
   a[ count - 1 ] = count - 1;
}
if ( i == count - 2 ) {
   a[ count -2 ] = count - 2;
}
```

372    Here are the results of unrolling the loop the second time:

| Language | Straight Time | Single Unrolled Time | Double Unrolled Time | Time Savings |
|---|---|---|---|---|
| C++ | 1.75 | 1.15 | 1.01 | 42% |
| **Java** | **1.01** | **0.581** | **0.581** | **43%** |
| PHP | 5.33 | 4.49 | 3.70 | 31% |
| Python | 2.51 | 3.21 | 2.79 | -12% |

373    *Note: Benchmarked for the case in which* count *equals* 100.

374    The results indicate that further loop unrolling can result in further time savings,
375    but not necessarily so, as the Java measurement shows. The main concern is how

1/13/2004 2:46 PM

376  Byzantine your code becomes. When you look at the code above, you might not
377  think it looks incredibly complicated, but when you realize that it started life a
378  couple of pages ago as a five-line loop, you can appreciate the trade-off between
379  performance and readability.

## Minimizing the Work Inside Loops

381  One key to writing effective loops is to minimize the work done inside a loop. If
382  you can evaluate a statement or part of a statement outside a loop so that only the
383  result is used inside the loop, do so. It's good programming practice, and, in
384  some cases, it improves readability.

385  Suppose you have a complicated pointer expression inside a hot loop that looks
386  like this:

387  **C++ Example of a Complicated Pointer Expression Inside a Loop**

```
for ( i = 0; i < rateCount; i++ ) {
   netRate[ i ] = baseRate[ i ] * rates->discounts->factors->net;
}
```

391  In this case, assigning the complicated pointer expression to a well-named
392  variable improves readability and often improves performance.

393  **C++ Example of Simplifying a Complicated Pointer Expression**

```
quantityDiscount = rates->discounts->factors->net;
for ( i = 0; i < rateCount; i++ ) {
   netRate[ i ] = baseRate[ i ] * quantityDiscount;
}
```

398  The extra variable, *quantityDiscount*, makes it clear that the *baseRate* array is
399  being multiplied by a quantity-discount factor to compute the net rate. That
400  wasn't at all clear from the original expression in the loop. Putting the
401  complicated pointer expression into a variable outside the loop also saves the
402  pointer from being dereferenced three times for each pass through the loop,
403  resulting in the following savings:

| Language | Straight Time | Code-Tuned Time | Time Savings |
|----------|--------------|-----------------|--------------|
| **C++** | **3.69** | **2.97** | **19%** |
| C# | 2.27 | 1.97 | 13% |
| Java | 4.13 | 2.35 | 43% |

404  *Note: Benchmarked for the case in which* rateCount *equals* 100*.*

405     Except for the Java compiler, the savings aren't anything to crow about,
406     implying that during initial coding you can use whichever technique is more
407     readable without worrying about the speed of the code until later.

408     ## Sentinel Values

409     When you have a loop with a compound test, you can often save time by
410     simplifying the test. If the loop is a search loop, one way to simplify the test is to
411     use a sentinel value, a value that you put just past the end of the search range and
412     that's guaranteed to terminate the search.

413     The classic example of a compound test that can be improved by use of a
414     sentinel is the search loop that checks both whether it has found the value it is
415     seeking and whether it has run out of values. Here's the code:

416     **C# Example of Compound Tests in a Search Loop**

417     
418     
419     *Here's the compound test.*
420     
421     
422     
423     
424     
425     
426     
427
428     
429     

```csharp
found = FALSE;
i = 0;
while ( ( !found ) && ( i < count ) ) {
   if ( item[ i ] == testValue ) {
      found = TRUE;
   }
   else {
      i++;
   }
}


if ( found ) {
   ...
```

430     In this code, each iteration of the loop tests for *!found* and for *i < count*. The
431     purpose of the *!found* test is to determine when the desired element has been
432     found. The purpose of the *i < count* test is to avoid running past the end of the
433     array. Inside the loop, each value of *item[]* is tested individually, so the loop
434     really has three tests for each iteration.

435     In this kind of search loop, you can combine the three tests so that you test only
436     once per iteration by putting a "sentinel" at the end of the search range to stop
437     the loop. In this case, you can simply assign the value you're looking for to the
438     element just beyond the end of the search range. (Remember to leave space for
439     that element when you declare the array.) You then check each element, and if
440     you don't find the element until you find the one you stuck at the end, you know
441     that the value you're looking for isn't really there. Here's the code:

**C# Example of Using a Sentinel Value to Speed Up a Loop**

442
443 `// set sentinel value, preserving the original value`
444 `initialValue = item[ count ];`
445 `item[ count ] = testValue;`
446
447 `i = 0;`
448 `while ( item[ i ] != testValue ) {`
449 `    i++;`
450 `}`
451
452 `// restore the value displaced by the sentinel`
453 `item[ count ] = initialValue;`
454
455 `// check if value was found`
456 `if ( i < count ) {`
457 `    ...`

*Remember to allow space for the sentinel value at the end of the array.*

458 When *item* is an array of integers, the savings can be dramatic:

| Language | Straight Time | Code-Tuned Time | Time Savings | Performance Ratio |
|---|---|---|---|---|
| **C#** | **0.771** | **0.590** | **23%** | **1.3:1** |
| Java | 1.63 | 0.912 | 44% | 2:1 |
| Visual Basic | 1.34 | 0.470 | 65% | 3:1 |

459 *Note: Search is of a 100-element array of integers.*

460 The Visual Basic results are particularly dramatic, but all the results are good.
461 When the kind of array changes, however, the results also change. Here are the
462 results when *item* is an array of single-precision floating-point numbers:

| Language | Straight Time | Code-Tuned Time | Time Savings |
|---|---|---|---|
| **C#** | **1.351** | **1.021** | **24%** |
| Java | 1.923 | 1.282 | 33% |
| Visual Basic | 1.752 | 1.011 | 42% |

463 *Note: Search is of a 100-element array of 4-byte floating-point numbers.*

464 As usual, the results vary significantly.

465 The sentinel technique can be applied to virtually any situation in which you use
466 a linear search—to linked lists as well as arrays. The only caveats are that you
467 must choose the sentinel value carefully and that you must be careful about how
468 you put the sentinel value into the array or linked list.

469 ## Putting the Busiest Loop on the Inside

470 When you have nested loops, think about which loop you want on the outside
471 and which you want on the inside. Following is an example of a nested loop that
472 can be improved.

---

473 **Java Example of a Nested Loop That Can Be Improved**

474
475
476
477
478
```
for ( column = 0; column < 100; column++ ) {
   for ( row = 0; row < 5; row++ ) {
      sum = sum + table[ row ][ column ];
   }
}
```

479 The key to improving the loop is that the outer loop executes much more often
480 than the inner loop. Each time the loop executes, it has to initialize the loop
481 index, increment it on each pass through the loop, and check it after each pass.
482 The total number of loop executions is 100 for the outer loop and 100 * 5 = 500
483 for the inner loop, for a total of 600 iterations. By merely switching the inner and
484 outer loops, you can change the total number of iterations to 5 for the outer loop
485 and 5 * 100 = 500 for the inner loop, for a total of 505 iterations. Analytically,
486 you'd expect to save about (600 – 505) / 600 = 16 percent by switching the
487 loops. Here's the measured difference in performance:

| Language | Straight Time | Code-Tuned Time | Time Savings |
|---|---|---|---|
| C++ | 4.75 | 3.19 | 33% |
| **Java** | **5.39** | **3.56** | **34%** |
| PHP | 4.16 | 3.65 | 12% |
| Python | 3.48 | 3.33 | 4% |

488 The results vary significantly, which shows once again that you have to measure
489 the effect in your particular environment before you can be sure your
490 optimization will help.

491 ## Strength Reduction

492 Reducing strength means replacing an expensive operation such as
493 multiplication with a cheaper operation such as addition. Sometimes you'll have
494 an expression inside a loop that depends on multiplying the loop index by a
495 factor. Addition is usually faster than multiplication, and if you can compute the
496 same number by adding the amount on each iteration of the loop rather than by
497 multiplying, the code will run faster. Here's an example of code that uses
498 multiplication:

**Visual Basic Example of Multiplying a Loop Index**

```
For i = 0 to saleCount - 1
    commission( i ) = (i + 1) * revenue * baseCommission * discount
Next
```

This code is straightforward but expensive. You can rewrite the loop so that you accumulate multiples rather than computing them each time. This reduces the strength of the operations from multiplication to addition. Here's the code:

**Visual Basic Example of Adding Rather Than Multiplying**

```
incrementalCommission = revenue * baseCommission * discount
cumulativeCommission = incrementalCommission
For i = 0 to saleCount - 1
    commission( i ) = cumulativeCommission
    cumulativeCommission = cumulativeCommission + incrementalCommission
Next
```

Multiplication is expensive, and this kind of change is like a manufacturer's coupon that gives you a discount on the cost of the loop. The original code incremented *i* each time and multiplied it by *revenue * baseCommission * discount*—first by 1, then by 2, then by 3, and so on. The optimized code sets *incrementalCommission* equal to *revenue * baseCommission * discount*. It then adds *incrementalCommission* to *cumulativeCommission* on each pass through the loop. On the first pass, it's been added once; on the second pass, it's been added twice; on the third pass, it's been added three times; and so on. The effect is the same as multiplying *incrementalCommission* by 1, then by 2, then by 3, and so on, but it's cheaper.

The key is that the original multiplication has to depend on the loop index. In this case, the loop index was the only part of the expression that varied, so the expression could be recoded more economically. Here's how much the rewrite helped in some test cases:

| Language | Straight Time | Code-Tuned Time | Time Savings |
|---|---|---|---|
| C++ | 4.33 | 3.80 | 12% |
| **Visual Basic** | **3.54** | **1.80** | **49%** |

*Note: Benchmark performed with* saleCount *equals* 20. *All computed variables are floating point.*

## 26.3 Data Transformations

Changes in data types can be a powerful aid in reducing program size and improving execution speed. Data-structure design is outside the scope of this

book, but modest changes in the implementation of a specific data type can also benefit performance. Here are a few ways to tune your data types.

## Use Integers Rather Than Floating-Point Numbers

Integer addition and multiplication tend to be faster than floating point. Changing a loop index from a floating point to an integer, for example, can save time. Here's an example:

**Visual Basic Example of a Loop That Uses a Time-Consuming Floating-Point Loop Index**

```
Dim i As Single
For i = 0 to 99
   x( i ) = 0
Next
```

Contrast this with a similar Visual Basic loop that explicitly uses the integer type:

**Visual Basic Example of a Loop That Uses a Timesaving Integer Loop Index**

```
Dim i As Integer
For i = 0 to 99
   x( i ) = 0
Next
```

How much difference does it make? Here are the results for this Visual Basic code and for similar code in C++ and PHP:

| Language | Straight Time | Code-Tuned Time | Time Savings | Performance Ratio |
|----------|---------------|-----------------|--------------|-------------------|
| C++ | 2.80 | 0.801 | 71% | 3.5:1 |
| PHP | 5.01 | 4.65 | 7% | 1:1 |
| **Visual Basic** | **6.84** | **0.280** | **96%** | **25:1** |

## Use the Fewest Array Dimensions Possible

Conventional wisdom maintains that multiple dimensions on arrays are expensive. If you can structure your data so that it's in a one-dimensional array rather than a two-dimensional or three-dimensional array, you might be able to save some time.

Suppose you have initialization code like this:

**Java Example of a Standard, Two-Dimensional Array Initialization**

560

561
562
563
564
565

```
for ( row = 0; row < numRows; row++ ) {
   for ( column = 0; column < numColumns; column++ ) {
      matrix[ row ][ column ] = 0;
   }
}
```

566
567
568

When this code is run with 50 rows and 20 columns, it takes twice as long with my current Java compiler as when the array is restructured so that it's one-dimensional. Here's how the revised code would look:

**Java Example of a One-Dimensional Representation of an Array**

569

570
571
572

```
for ( entry = 0; entry < numRows * numColumns; entry++ ) {
   matrix[ entry ] = 0;
}
```

573
574

Here's a summary of the results, with the addition of comparable results in several other languages:

| Language | Straight Time | Code-Tuned Time | Time Savings | Performance Ratio |
|---|---|---|---|---|
| C++ | 8.75 | 7.82 | 11% | 1:1 |
| C# | 3.28 | 2.99 | 9% | 1:1 |
| **Java** | **7.78** | **4.14** | **47%** | **2:1** |
| PHP | 6.24 | 4.10 | 34% | 1.5:1 |
| Python | 3.31 | 2.23 | 32% | 1.5:1 |
| Visual Basic | 9.43 | 3.22 | 66% | 3:1 |

575
576

*Note: Times for Python and PHP aren't directly comparable to times for the other languages because they were run <1% as many iterations as the other languages.*

577
578
579

The results of this optimization are excellent in Visual Basic and Java, good in PHP and Python, but mediocre in C++ and C#. Of course the C++ compiler's unoptimized time was easily the best of the group, so you can't be too hard on it.

580
581
582

This wide range of results also show the hazard of following any code-tuning advice blindly. You can never be sure until you try the advice in your specific circumstances.

583

# Minimize Array References

584
585
586
587

In addition to minimizing accesses to doubly or triply dimensioned arrays, it's often advantageous to minimize array accesses, period. A loop that repeatedly uses one element of an array is a good candidate for the application of this technique. Here's an example of an unnecessary array access:

H:\books\CodeC2Ed\Reviews\Web\26-TuningTechniques.doc

**C++ Example of Unnecessarily Referencing an Array Inside a Loop**

```
for ( discountType = 0; discountType < typeCount; discountType++ ) {
   for ( discountLevel = 0; discountLevel < levelCount; discountLevel++ ) {
      rate[ discountLevel ] = rate[ discountLevel ] * discount[ discountType ];
   }
}
```

The reference to *discount[ discountType ]* doesn't change when *discountLevel* changes in the inner loop. Consequently, you can move it out of the inner loop so that you'll have only one array access per execution of the outer loop rather than one for each execution of the inner loop. The next example shows the revised code.

**C++ Example of Moving an Array Reference Outside a Loop**

```
for ( discountType = 0; discountType < typeCount; discountType++ ) {
   thisDiscount = discount[ discountType ];
   for ( discountLevel = 0; discountLevel < levelCount; discountLevel++ ) {
      rate[ discountLevel ] = rate[ discountLevel ] * thisDiscount;
   }
}
```

Here are the results:

| Language | Straight Time | Code-Tuned Time | Time Savings |
|---|---|---|---|
| **C++** | **32.1** | **34.5** | **-7%** |
| C# | 18.3 | 17.0 | 7% |
| Visual Basic | 23.2 | 18.4 | 20% |

*Note: Benchmark times were computed for the case in which* typeCount *equals* 10 *and* levelCount *equals* 100.

As usual, the results vary significantly from compiler to compiler.

# Use Supplementary Indexes

Using a supplementary index means adding related data that makes accessing a data type more efficient. You can add the related data to the main data type, or you can store it in a parallel structure.

## String-Length Index

One example of using a supplementary index can be found in the different string-storage strategies. In C, strings are terminated by a byte that's set to 0. In Visual Basic string format, a length byte hidden at the beginning of each string indicates how long the string is. To determine the length of a string in C, a program has to start at the beginning of the string and count each byte until it

620  finds the byte that's set to 0. To determine the length of a Visual Basic string, the
621  program just looks at the length byte. Visual Basic length byte is an example of
622  augmenting a data type with an index to make certain operations—like
623  computing the length of a string—faster.

624  You can apply the idea of indexing for length to any variable-length data type.
625  It's often more efficient to keep track of the length of the structure rather than
626  computing the length each time you need it.

### Independent, Parallel Index Structure

627

628  Sometimes it's more efficient to manipulate an index to a data type than it is to
629  manipulate the data type itself. If the items in the data type are big or hard to
630  move (on disk, perhaps), sorting and searching index references is faster than
631  working with the data directly. If each data item is large, you can create an
632  auxiliary structure that consists of key values and pointers to the detailed
633  information. If the difference in size between the data-structure item and the
634  auxiliary-structure item is great enough, sometimes you can store the key item in
635  memory even when the data item has to be stored externally. All searching and
636  sorting is done in memory, and you have to access the disk only once, when you
637  know the exact location of the item you want.

## Use Caching

638

639  Caching means saving a few values in such a way that you can retrieve the most
640  commonly used values more easily than the less commonly used values. If a
641  program randomly reads records from a disk, for example, a routine might use a
642  cache to save the records read most frequently. When the routine receives a
643  request for a record, it checks the cache to see whether it has the record. If it
644  does, the record is returned directly from memory rather than from disk.

645  In addition to caching records on disk, you can apply caching in other areas. In a
646  Microsoft Windows font-proofing program, the performance bottleneck was in
647  retrieving the width of each character as it was displayed. Caching the most
648  recently used character width roughly doubled the display speed.

649  You can cache the results of time-consuming computations too—especially if the
650  parameters to the calculation are simple. Suppose, for example, that you need to
651  compute the length of the hypotenuse of a right triangle, given the lengths of the
652  other two sides. The straightforward implementation of the routine would look
653  like this:

654  **Java Example of a Routine That's Conducive to Caching**

```
655  double Hypotenuse(
656     double sideA,
```

```
657          double sideB
658          ) {
659          return Math.sqrt( ( sideA * sideA ) + ( sideB * sideB ) );
660     }
```

661     If you know that the same values tend to be requested repeatedly, you can cache
662     values this way:

---

663     **Java Example of Caching to Avoid an Expensive Computation**

```
664     private double cachedHypotenuse = 0;
665     private double cachedSideA = 0;
666     private double cachedSideB = 0;
667
668     public double Hypotenuse(
669          double sideA,
670          double sideB
671          ) {
672          // check to see if the triangle is already in the cache
673          if ( ( sideA == cachedSideA ) && ( sideB == cachedSideB ) ) {
674             return cachedHypotenuse;
675          }
676
677          // compute new hypotenuse and cache it
678          cachedHypotenuse = Math.sqrt( ( sideA * sideA ) + ( sideB * sideB ) );
679          cachedSideA = sideA;
680          cachedSideB = sideB;
681
682          return cachedHypotenuse;
683     }
```

684     The second version of the routine is more complicated than the first and takes up
685     more space, so speed has to be at a premium to justify it. Many caching schemes
686     cache more than one element, so they have even more overhead. Here's the
687     speed difference between these two versions:

| Language | Straight Time | Code-Tuned Time | Time Savings | Performance Ratio |
|---|---|---|---|---|
| C++ | 4.06 | 1.05 | 74% | 4:1 |
| **Java** | **2.54** | **1.40** | **45%** | **2:1** |
| Python | 8.16 | 4.17 | 49% | 2:1 |
| Visual Basic | 24.0 | 12.9 | 47% | 2:1 |

688     *Note: The results shown assume that the cache is hit twice for each time it's set.*

689     The success of the cache depends on the relative costs of accessing a cached
690     element, creating an uncached element, and saving a new element in the cache.

Success also depends on how often the cached information is requested. In some cases, success might also depend on caching done by the hardware. Generally, the more it costs to generate a new element and the more times the same information is requested, the more valuable a cache is. The cheaper it is to access a cached element and save new elements in the cache, the more valuable a cache is. As with other optimization techniques, caching adds complexity and tends to be error prone.

# 26.4 Expressions

Much of the work in a program is done inside mathematical or logical expressions. Complicated expressions tend to be expensive, so this section looks at ways to make them cheaper.

## Exploit Algebraic Identities

You can use algebraic identities to replace costly operations with cheaper ones. For example, the following expressions are logically equivalent:

```
not a and not B
not (a or B)
```

If you choose the second expression instead of the first, you can save a *not* operation.

Although the savings from avoiding a single *not* operation are probably inconsequential, the general principle is powerful. Jon Bentley describes a program that tested whether *sqrt(x) < sqrt(y)* (1982). Since *sqrt(x)* is less than *sqrt(y)* only when *x* is less than *y*, you can replace the first test with $x < y$. Given the cost of the *sqrt( )* routine, you'd expect the savings to be dramatic, and they are. Here are the results:

| Language | Straight Time | Code-Tuned Time | Time Savings | Performance Ratio |
|---|---|---|---|---|
| C++ | 7.43 | 0.010 | 99.9% | 750:1 |
| Visual Basic | 4.59 | 0.220 | 95% | 20:1 |
| Python | 4.21 | 0.401 | 90% | 10:1 |

## Use Strength Reduction

As mentioned earlier, strength reduction means replacing an expensive operation with a cheaper one. Here are some possible substitutions:

718          ●    Replace multiplication with addition.

719          ●    Replace exponentiation with multiplication.

720          ●    Replace trigonometric routines with their trigonometric identities.

721          ●    Replace *longlong* integers with *long*s or *int*s (but watch for performance
722               issues associated with using native-length vs. non-native–length integers)

723          ●    Replace floating-point numbers with fixed-point numbers or integers.

724          ●    Replace double-precision floating points with single-precision numbers.

725          ●    Replace integer multiplication-by-two and division-by-two with shift
726               operations.

727     Here is a detailed example. Suppose you have to evaluate a polynomial. If you're
728     rusty on polynomials, they're the things that look like

729     $Ax^2 + Bx + C$

730     The letters *A*, *B*, and *C* are coefficients, and *x* is a variable. General code to
731     evaluate an *n*th-order polynomial looks like this:

**Visual Basic Example of Evaluating a Polynomial**

```
value = coefficient( 0 )
For power = 1 To order
   value = value + coefficient( power ) * x^power
Next
```

737     If you're thinking about strength reduction, you'll look at the exponentiation
738     operator with a jaundiced eye. One solution would be to replace the
739     exponentiation with a multiplication on each pass through the loop, which is
740     analogous to the strength-reduction case a few sections ago in which a
741     multiplication was replaced with an addition. Here's how the reduced-strength
742     polynomial evaluation would look:

**Visual Basic Example of a Reduced-Strength Method of Evaluating a Polynomial**

```
value = coefficient( 0 )
powerOfX = x
For power = 1 to order
   value = value + coefficient( power ) * powerOfX
   powerOfX = powerOfX * x
Next
```

751     This produces a noticeable advantage if you're working with second-order
752     polynomials (polynomials in which the highest-power term is squared) or higher-
753     order polynomials.

| Language | Straight Time | Code-Tuned Time | Time Savings | Performance Ratio |
|---|---|---|---|---|
| Python | 3.24 | 2.60 | 20% | 1:1 |
| **Visual Basic** | **6.26** | **0.160** | **97%** | **40:1** |

754 If you're serious about strength reduction, you still won't care for those two
755 floating-point multiplications. The strength-reduction principle suggests that you
756 can further reduce the strength of the operations in the loop by accumulating
757 powers rather than multiplying them each time. Here's that code:

758 **Visual Basic Example of Further Reducing the Strength Required to**
759 **Evaluate a Polynomial**

```
760 value = 0
761 For power = order to 1 Step -1
762    value = ( value + coefficient( power ) ) * x
763 Next
764 value = value + coefficient( 0 )
```

765 This method eliminates the extra *powerOfX* variable and replaces the two
766 multiplications in each pass through the loop with one.

| Language | Straight Time | First Optimization | Second Optimization | Savings over First Optimization |
|---|---|---|---|---|
| Python | 3.24 | 2.60 | 2.53 | 3% |
| **Visual Basic** | **6.26** | **0.16** | **0.31** | **-94%** |

767 This is a good example of theory not holding up very well to practice. The code
768 with reduced strength seems like it should be faster, but it isn't. One possibility
769 is that decrementing a loop by *−1* instead of incrementing it by *+1* in Visual
770 Basic hurts performance, but you'd have to measure that hypothesis to be sure.

## Initialize at Compile Time

772 If you're using a named constant or a magic number in a routine call and it's the
773 only argument, that's a clue that you could precompute the number, put it into a
774 constant, and avoid the routine call. The same principle applies to
775 multiplications, divisions, additions, and other operations.

776 I once needed to compute the base-two logarithm of an integer, truncated to the
777 nearest integer. The system didn't have a log-base-two routine, so I wrote my
778 own. The quick and easy approach was to use the fact that

779 $$\log(x)_{base} = \log(x) / \log(base)$$

780 Given this identity, I could write a routine like this one:

781

782

783

784

**C++ Example of a Log-Base-Two Routine Based on System Routines**

```
unsigned int Log2( unsigned int x ) {
   return (unsigned int) ( log( x ) / log( 2 ) );
}
```

785  This routine was really slow, and since the value of *log(2)* never changed, I

786  replaced *log(2)* with its computed value, *0.69314718*. Then the code looked like

787  this:

788  **C++ Example of a Log-Base-Two Routine Based on a System Routine**

789  **and a Constant**

790
```
unsigned int Log2( unsigned int x ) {
   return (unsigned int) ( log( x ) / LOG2 );
}
```

791  LOG2 *is a named constant*

792  *equal to* 0.69314718.

793  Since *log()* tends to be an expensive routine, much more expensive than type

794  conversions or division, you'd expect that cutting the calls to the *log()* function

795  by half would cut the time required for the routine by about half. Here are the

796  measured results:

| Language | Straight Time | Code-Tuned Time | Time Savings |
|----------|---------------|-----------------|--------------|
| **C++**  | **9.66**      | **5.97**        | **38%**      |
| Java     | 17.0          | 12.3            | 28%          |
| PHP      | 2.45          | 1.50            | 39%          |

797  In this case, the educated guess about the relative importance of the division and

798  type conversions and the estimate of 50 percent were pretty close. Considering

799  the predictability of the results described in this chapter, the accuracy of my

800  prediction in this case proves only that even a blind squirrel finds a nut

801  occasionally.

802  ## Be Wary of System Routines

803  System routines are expensive and provide accuracy that's often wasted. Typical

804  system math routines, for example, are designed to put an astronaut on the moon

805  within ± 2 feet of the target. If you don't need that degree of accuracy, you don't

806  need to spend the time to compute it either.

807  In the previous example, the *Log2()* routine returned an integer value but used a

808  floating-point *log()* routine to compute it. That was overkill for an integer result,

809  so after my first attempt, I wrote a series of integer tests that were perfectly

810  accurate for calculating an integer $\log_2$. Here's the code:

811  **C++ Example of a Log-Base-Two Routine Based on Integers**

812
```
unsigned int Log2( unsigned int x ) {
```

```
813        if ( x < 2 ) return 0 ;
814        if ( x < 4 ) return 1 ;
815        if ( x < 8 ) return 2 ;
816        if ( x < 16 ) return 3 ;
817        if ( x < 32 ) return 4 ;
818        if ( x < 64 ) return 5 ;
819        if ( x < 128 ) return 6 ;
820        if ( x < 256 ) return 7 ;
821        if ( x < 512 ) return 8 ;
822        if ( x < 1024 ) return 9 ;
823        ...
824        if ( x < 2147483648 ) return 30;
825        return 31 ;
826     }
```

827     This routine uses integer operations, never converts to floating point, and blows
828     the doors off both floating-point versions. Here are the results:

| Language | Straight Time | Code-Tuned Time | Time Savings | Performance Ratio |
|----------|---------------|-----------------|--------------|-------------------|
| **C++**  | **9.66**      | **0.662**       | **93%**      | **15:1**          |
| Java     | 17.0          | 0.882           | 95%          | 20:1              |
| PHP      | 2.45          | 3.45            | -41%         | 2:3               |

829     Most of the so-called "transcendental" functions are designed for the worst
830     case—that is, they convert to double-precision floating point internally even if
831     you give them an integer argument. If you find one in a tight section of code and
832     don't need that much accuracy, give it your immediate attention.

833     Another option is to take advantage of the fact that a right-shift operation is the
834     same as dividing by two. The number of times you can divide a number by two
835     and still have a nonzero value is the same as the $\log_2$ of that number. Here's how
836     code based on that observation looks:

**CODING HORROR**

### C++ Example of an Alternative Log-Base-Two Routine Based on the Right-Shift Operator

```
839     unsigned int Log2( unsigned int x ) {
840        unsigned int i = 0;
841        while ( ( x = ( x >> 1 ) ) != 0 ) {
842           i++;
843        }
844        return i ;
845     }
```

846
847
848

To non-C++ programmers, this code is particularly hard to read. The complicated expression in the *while* condition is an example of a coding practice you should avoid unless you have a good reason to use it.

849
850
851

This routine takes about 350 percent longer than the longer version above, executing in 2.4 seconds rather than 0.66 seconds. But it's faster than the first approach, and adapts easily to 32-bit, 64-bit, and other environments.

852 **KEY POINT**
853
854
855

This example highlights the value of not stopping after one successful optimization. The first optimization earned a respectable 30-40 percent savings but had nowhere near the impact of the second optimization or third optimizations.

## Use the Correct Type of Constants

857
858
859
860
861

Use named constants and literals that are the same type as the variables they're assigned to. When a constant and its related variable are different types, the compiler has to do a type conversion to assign the constant to the variable. A good compiler does the type conversion at compile time so that it doesn't affect - run-time performance.

862
863
864
865

A less advanced compiler or an interpreter generates code for a runtime conversion, so you might be stuck. Here are some differences in performance between the initializations of a floating-point variable $x$ and an integer variable $i$ in two cases. In the first case, the initializations look like this:

866
867

```
x = 5
i = 3.14
```

868
869

and require type conversions, assuming $x$ is a floating point variable and $i$ is an integer In the second case, they look like this:

870
871

```
x = 3.14
i = 5
```

872

and don't require type conversions. Here are the results:

| Language | Straight Time | Code-Tuned Time | Time Savings | Performance Ratio |
|---|---|---|---|---|
| C++ | 1.11 | 0.000 | 100% | not measurable |
| C# | 1.49 | 1.48 | <1% | 1:1 |
| Java | 1.66 | 1.11 | 33% | 1.5:1 |
| Visual Basic | 0.721 | 0.000 | 100% | not measurable |
| PHP | 0.872 | 0.847 | 3% | 1:1 |

873

The variation among compilers is once again notable.

# Precompute Results

874

A common low-level design decision is the choice of whether to compute results
on the fly or compute them once, save them, and look them up as needed. If the
results are used many times, it's often cheaper to compute them once and look
them up the rest of the time.

This choice manifests itself in several ways. At the simplest level, you might
compute part of an expression outside a loop rather than inside. An example of
this appeared earlier in the chapter. At a more complicated level, you might
compute a lookup table once when program execution begins, using it every time
thereafter, or you might store results in a data file or embed them in a program.

In a space-wars video game, for example, the programmers initially computed
gravity coefficients for different distances from the sun. The computation for the
gravity coefficients was expensive and affected performance. The program
recognized relatively few distinct distances from the sun, however, so the
programmers were able to precompute the gravity coefficients and store them in
a 10-element array. The array lookup was much faster than the expensive
computation.

Suppose you have a routine that computes payment amounts on automobile
loans. The code for such a routine would look like this:

**Java Example of a Complex Computation That Could Be Precomputed**

```java
double ComputePayment(
   long loanAmount,
   int months,
   double interestRate
   ) {
   return loanAmount /
     (
     ( 1.0 - Math.pow( ( 1.0 + ( interestRate / 12.0 ) ), -months ) ) /
     ( interestRate / 12.0 )
     );
}
```

The formula for computing loan payments is complicated and fairly expensive.
Putting the information into a table instead of computing it each time would
probably be cheaper.

How big would the table be? The widest-ranging variable is *loanAmount*. The
variable *interestRate* might range from 5 percent through 20 percent by quarter
points, but that's only 61 distinct rates. *months* might range from 12 through 72,
but that's only 61 distinct periods. *loanAmount* could conceivably range from

912
913

$1000 through $100,000, which is more entries than you'd generally want to handle in a lookup table.

914
915
916
917

Most of the computation doesn't depend on *loanAmount*, however, so you can put the really ugly part of the computation (the denominator of the larger expression) into a table that's indexed by *interestRate* and *months*. You recompute the *loanAmount* part each time. Here's the revised code:

918

**Java Example of Precomputing a Complex Computation**

919
920
921
922
923
924
925
926
927

*The new variable* interestIndex *is created to provide a subscript into the* loanDivisor *array.*

```java
double ComputePayment(
   long loanAmount,
   int months,
   double interestRate
   ) {
   int interestIndex =
      Math.round( ( interestRate - LOWEST_RATE ) * GRANULARITY * 100.00 );
   return loanAmount / loanDivisor[ interestIndex ][ months ];
}
```

928
929

In this code, the hairy calculation has been replaced with the computation of an array index and a single array access. Here are the results of the change:

| Language | Straight Time | Code-Tuned Time | Time Savings | Performance Ratio |
|---|---|---|---|---|
| **Java** | **2.97** | **0.251** | **92%** | **10:1** |
| Python | 3.86 | 4.63 | -20% | 1:1 |

930
931
932
933
934

Depending on your circumstances, you would need to precompute the *loanDivisor* array at program initialization time or read it from a disk file. Alternatively, you could initialize it to *0*, compute each element the first time it's requested, store it, and look it up each time it's requested subsequently. That would be a form of caching, discussed earlier.

935
936
937
938
939

You don't have to create a table to take advantage of the performance gains you can achieve by precomputing an expression. Code similar to the code in the previous examples raises the possibility of a different kind of precomputation. Suppose you have code that computes payments for many loan amounts, as shown here.

940
941

**Java Example of a Second Complex Computation That Could Be Precomputed**

942
943
944
945

```java
double ComputePayments(
   int months,
   double interestRate
   ) {
```

```
946          for ( long loanAmount = MIN_LOAN_AMOUNT; loanAmount < MAX_LOAN_AMOUNT;
947             loanAmount++ ) {
948             payment = loanAmount / (
949                ( 1.0 – Math.pow( 1.0+(interestRate/12.0), - months ) ) /
950                ( interestRate/12.0 )
951                );
952             ...
953          }
954       }
```

*952  The following code would do*
*953  something with payment here;*
*954  for this example's point, it*
*955  doesn't matter what.*

Even without precomputing a table, you can precompute the complicated part of
the expression outside the loop and use it inside the loop. Here's how it would
look:

**Java Example of Precomputing the Second Complex Computation**

```
959   double ComputePayments(
960      int months,
961      double interestRate
962      ) {
963      long loanAmount;
964      double divisor = ( 1.0 – Math.pow( 1.0+(interestRate/12.0). - months ) ) /
965         ( interestRate/12.0 );
966      for ( long loanAmount = MIN_LOAN_AMOUNT; loanAmount <= MAX_LOAN_AMOUNT;
967         loanAmount++ ) {
968         payment = loanAmount / divisor;
969         ...
970      }
971   }
```

*964  Here's the part that's*
*965  precomputed.*

This is similar to the techniques suggested earlier of putting array references and
pointer dereferences outside a loop. The results for Java in this case are
comparable to the results of using the precomputed table in the first
optimization:

| Language | Straight Time | Code-Tuned Time | Time Savings | Performance Ratio |
|----------|---------------|-----------------|--------------|-------------------|
| **Java** | **7.43** | **0.24** | **97%** | **30:1** |
| Python | 5.00 | 1.69 | 66% | 3:1 |

Python improved here, but not in the first optimization attempt. Many times
when one optimization does not produce the desired results, a seemingly similar
optimization will work as expected.

Optimizing a program by precomputation can take several forms:

● Computing results before the program executes and wiring them into
  constants that are assigned at compile time

982
983

- Computing results before the program executes and hard-coding them into variables used at run time

984
985

- Computing results before the program executes and putting them into a file that's loaded at run time

986
987

- Computing results once, at program startup, and then referencing them each time they're needed

988
989

- Computing as much as possible before a loop begins, minimizing the work done inside the loop

990
991

- Computing results the first time they're needed and storing them so that you can retrieve them when they're needed again

992

## Eliminate Common Subexpressions

993
994
995
996

If you find an expression that's repeated several times, assign it to a variable and refer to the variable rather than recomputing the expression in several places. The loan-calculation example has a common subexpression that you could eliminate. Here's the original code:

997

**Java Example of a Common Subexpression**

998
999
1000
1001

```
payment = loanAmount / (
     ( 1.0 – Math.pow( 1.0 + ( interestRate / 12.0 ), -months ) ) /
     ( interestRate / 12.0 )
   );
```

1002
1003
1004
1005
1006

In this sample, you can assign *interestRate/12.0* to a variable that is then referenced twice rather than computing the expression twice. If you have chosen the variable name well, this optimization can improve the code's readability at the same time that it improves performance. The next example shows the revised code.

1007

**Java Example of Eliminating a Common Subexpression**

1008
1009
1010
1011
1012

```
monthlyInterest = interestRate / 12.0;
payment = loanAmount / (
     ( 1.0 – Math.pow( 1.0 + monthlyInterest, -months ) ) /
     monthlyInterest
   );
```

1013

The savings in this case don't seem impressive:

| Language | Straight Time | Code-Tuned Time | Time Savings |
| --- | --- | --- | --- |
| **Java** | **2.94** | **2.83** | **4%** |
| Python | 3.91 | 3.94 | -1% |

1/13/2004 2:46 PM

1014    It appears that the *Math.pow()* routine is so costly that it overshadows the
1015    savings from subexpression elimination. Or possibly the subexpression is already
1016    being eliminated by the compiler. If the subexpression were a bigger part of the
1017    cost of the whole expression or if the compiler optimizer were less effective, the
1018    optimization might have more impact.

## 26.5 Routines

1020    One of the most powerful tools in code tuning is a good routine decomposition.
Small, well-defined routines save space because they take the place of doing jobs
separately in multiple places. They make a program easy to optimize because
you can refactor code in one routine and thus improve every routine that calls it.
Small routines are relatively easy to rewrite in assembler. Long, tortuous
routines are hard enough to understand on their own; in assembler they're
impossible.

**CROSS-REFERENCE** For details on working with routines, see Chapter 7, "High-Quality Routines."

### Rewrite Routines In Line

1028    In the early days of computer programming, some machines imposed prohibitive
1029    performance penalties for calling a routine. A call to a routine meant that the
1030    operating system had to swap out the program, swap in a directory of routines,
1031    swap in the particular routine, execute the routine, swap out the routine, and
1032    swap the calling routine back in. All this swapping chewed up resources and
1033    made the program slow.

1034    Modern computers collect a far smaller toll for calling a routine. Here are the
1035    results of putting a string-copy routine in line:

| Language | Routine Time | Inline-Code Time | Time Savings |
| --- | --- | --- | --- |
| C++ | 0.471 | 0.431 | 8% |
| Java | 13.1 | 14.4 | -10% |

1036    In some cases, you might be able to save a few nanoseconds by putting the code
1037    from a routine into the program directly where it's needed using a language
1038    feature like C++'s *inline* keyword. If you're working in a language that doesn't
1039    support *inline* directly but that does have a macro preprocessor, you can use a
1040    macro to put the code in, switching it in and out as needed. But modern
1041    machines—and "modern" means any machine you're ever likely to work on—
1042    impose virtually no penalty for calling a routine. As the example shows, you're
1043    as likely to degrade performance by keeping code inline as to optimize it.

1044 # 26.6 Recoding in Assembler

1045 One longstanding piece of conventional wisdom that shouldn't be left
1046 unmentioned is the advice that when you run into a performance bottleneck, you
1047 should recode in assembler. Recoding in assembler tends to improve both speed
1048 and code size. Here is a typical approach to optimizing with assembler:

1049 1.  Write 100 percent of an application in a high-level language.

1050 2.  Fully test the application, and verify that it's correct.

1051 **CROSS-REFERENCE**  For
1052 details on the phenomenon of
1053 a small percentage of a
1054 program accounting for most
of its run time, see "The
Pareto Principle" in Section
1055 25.2.

3.  If performance improvements are needed after that, profile the application to
    identify hot spots. Since about 5 percent of a program usually accounts for
    about 50 percent of the running time, you can usually identify small pieces
    of the program as hot spots.

1055 4.  Recode a few small pieces in assembler to improve overall performance.

1056 Whether you follow this well-beaten path depends on how comfortable you are
1057 with assembler, how well-suited the problem is to assembler, and on your level
1058 of desperation.

1059 I got my first exposure to assembler on the DES encryption program I mentioned
1060 in the previous chapter. I had tried every optimization I'd ever heard of, and the
1061 program was still twice as slow as the speed goal. Recoding part of the program
1062 in assembler was the only remaining option. As an assembler novice, about all I
1063 could do was make a straight translation from a high-level language to
1064 assembler, but I got a 50 percent improvement even at that rudimentary level.

1065 Suppose you have a routine that converts binary data to uppercase ASCII
1066 characters. The next example shows the Delphi code to do it.

1067 **Delphi Example of Code That's Better Suited to Assembler**

```
1068  procedure HexExpand(
1069     var source: ByteArray;
1070     var target: WordArray;
1071     byteCount: word
1072  );
1073  var
1074     index: integer;
1075     lowerByte: byte;
1076     upperByte: byte;
1077     targetIndex: integer;
1078  begin
```

```
1079                      targetIndex := 1;
1080                      for index := 1 to byteCount do begin
1081                         target[ targetIndex ] := ( (source[ index ] and $F0) shr 4 ) + $41;
1082                         target[ targetIndex+1 ] := (source[ index ] and $0f) + $41;
1083                         targetIndex := targetIndex + 2;
1084                      end;
1085                   end;
```

1086    Although it's hard to see where the fat is in this code, it contains a lot of bit
1087    manipulation, which isn't exactly Delphi's forte. Bit manipulation is assembler's
1088    forte, however, so this code is a good candidate for recoding. Here's the
1089    assembler code:

1090    **Example of a Routine Recoded in Assembler**

```
1091    procedure HexExpand(
1092       var source;
1093       var target;
1094       byteCount : Integer
1095    );
1096        label
1097        EXPAND;
1098
1099        asm
1100              MOV    ECX,byteCount      // load number of bytes to expand
1101              MOV    ESI,source         // source offset
1102              MOV    EDI,target         // target offset
1103              XOR    EAX,EAX            // zero out array offset
1104
1105        EXPAND:
1106              MOV    EBX,EAX            // array offset
1107              MOV    DL,[ESI+EBX]       // get source byte
1108              MOV    DH,DL              // copy source byte
1109
1110              AND    DH,$F              // get msbs
1111              ADD    DH,$41             // add 65 to make upper case
1112
1113              SHR    DL,4               // move lsbs into position
1114              AND    DL,$F              // get lsbs
1115              ADD    DL,$41             // add 65 to make upper case
1116
1117              SHL    BX,1               // double offset for target array offset
1118              MOV    [EDI+EBX],DX       // put target word
1119
1120              INC    EAX                // increment array offset
1121              LOOP   EXPAND             // repeat until finished
1122           end;
```

1123
1124
1125
1126

Rewriting in assembler in this case was profitable, resulting in a time savings of 41 percent. It's logical to assume that code in a language that's more suited to bit manipulation—C++, for instance—would have less to gain than Delphi code would. Here are the results:

| Language | High-Level Time | Assembler Time | Time Savings |
|----------|-----------------|----------------|--------------|
| C++      | 4.25            | 3.02           | 29%          |
| **Delphi** | **5.18**      | **3.04**       | **41%**      |

1127
1128
1129
1130

The "before" picture in this measurements reflects the two languages' strengths at bit manipulation. The "after" picture looks virtually identical, and it appears that the assembler code has minimized the initial performance differences between Delphi and C++.

1131
1132
1133
1134

The assembler routine shows that rewriting in assembler doesn't have to produce a huge, ugly routine. Such routines are often quite modest, as this one is. Sometimes assembler code is almost as compact as its high-level-language equivalent.

1135
1136
1137
1138
1139
1140
1141
1142

A relatively easy and effective strategy for recoding in assembler is to start with a compiler that generates assembler listings as a by-product of compilation. Extract the assembler code for the routine you need to tune, and save it in a separate source file. Using the compiler's assembler code as a base, hand-optimize the code, checking for correctness and measuring improvements at each step. Some compilers intersperse the high-level-language statements as comments in the assembler code. If yours does, you might keep them in the assembler code as documentation.

CC2E.COM/2672

1143

## CHECKLIST: Code-Tuning Techniques

1144

**Improve Both Speed and Size**

1145    ❑ Substitute table lookups for complicated logic

1146    ❑ Jam loops

1147    ❑ Use integer instead of floating-point variables

1148    ❑ Initialize data at compile time

1149    ❑ Use constants of the correct type

1150    ❑ Precompute results

1151    ❑ Eliminate common subexpressions

1152    ❑ Translate key routines to assembler

**Improve Speed Only**

❑ Stop testing when you know the answer

❑ Order tests in *case* statements and *if-then-else* chains by frequency

❑ Compare performance of similar logic structures

❑ Use lazy evaluation

❑ Unswitch loops that contain *if* tests

❑ Unroll loops

❑ Minimize work performed inside loops

❑ Use sentinels in search loops

❑ Put the busiest loop on the inside of nested loops

❑ Reduce the strength of operations performed inside loops

❑ Change multiple-dimension arrays to a single dimension

❑ Minimize array references

❑ Augment data types with indexes

❑ Cache frequently used values

❑ Exploit algebraic identities

❑ Reduce strength in logical and mathematical expressions

❑ Be wary of system routines

❑ Rewrite routines in line

# 26.7 The More Things Change, the More They Stay the Same

You might expect that performance attributes of systems would have changed somewhat in the 10 years since I wrote the first edition of *Code Complete*, and in some ways they have. Computers are dramatically faster and memory is more plentiful. In the first edition, I ran most of the tests in this chapter 10,000 to 50,000 times to get meaningful, measurable results. For this edition I had to run most tests 1 million to 100 million times. When you have to run a test 100 million times to get measurable results, you have to ask whether anyone will ever notice the impact in a real program. Computers have become so powerful that for many common kinds of programs the level of performance optimization discussed in this chapter has become irrelevant.

In other ways, performance issues have hardly changed at all. People writing desktop applications may not need this information, but people writing software

1187    for embedded systems, real-time systems, and other systems with strict speed or
1188    space restrictions can still benefit from this information.

1189    The need to measure the impact of each and every attempt at code tuning has
1190    been a constant since Donald Knuth published his study of Fortran programs in
1191    1971. According to the measurements in this chapter, the effect of any specific
1192    optimization is actually *less predictable* than it was 10 years ago. The effect of
1193    each code tuning is affected by the programming language, compiler, compiler
1194    version, code libraries, library versions, and compiler settings, among other
1195    things.

1196    Code tuning invariably involves tradeoffs among complexity, readability,
1197    simplicity, and maintainability on the one hand and a desire to improve
1198    performance on the other. It introduces a high degree of maintenance overhead
1199    because of all the reprofiling that's required.

1200    I have found that insisting on *measurable improvement* is a good way to resist
1201    the temptation to optimize prematurely and to enforce a bias toward clear,
1202    straightforward code. If an optimization is important enough to haul out the
1203    profiler and measure the optimization's effect, then it's probably important
1204    enough to allow—as long as it works. But if an optimization isn't important
1205    enough to haul out the profiling machinery, then it isn't important enough to
1206    degrade readability, maintainability, and other code characteristics. The impact
1207    of unmeasured code tuning on performance is speculative at best, whereas the
1208    impact on readability is as certain as it is detrimental.

CC2E.COM/2679

# Additional Resources

1209

1210    My favorite reference on code tuning is *Writing Efficient Programs* (Bentley,
1211    Englewood Cliffs, N.J.: Prentice Hall, 1982). The book is out of print, but worth
1212    reading if you can find it. It's an expert treatment of code tuning, broadly
1213    considered. Bentley describes techniques that trade time for space and space for
1214    time. He provides several examples of redesigning data types to reduce both
1215    space and time. His approach is a little more anecdotal than the one taken here,
1216    and his anecdotes are interesting. He takes a few routines through several
1217    optimization steps so that you can see the effects of first, second, and third
1218    attempts on a single problem. Bentley strolls through the primary contents of the
1219    book in 135 pages. The book has an unusually high signal-to-noise ratio—it's
1220    one of the rare gems that every practicing programmer should own.

1221    Appendix 4 of Bentley's *Programming Pearls*, 2d Ed. (2000), contains a
1222    summary of the code tuning rules from his earlier book.

1223    You can also find a full array of technology-specific optimization books. Several
1224    are listed below, and the web link to the left contains an up-to-date list.

1225  CC2E.COM/2686    Booth, Rick. *Inner Loops : A Sourcebook for Fast 32-bit Software Development*,
1226    Boston, Mass.: Addison Wesley, 1997.

1227    Gerber, Richard. *Software Optimization Cookbook: High-Performance Recipes*
1228    *for the Intel Architecture*, Intel Press, 2002.

1229    Hasan, Jeffrey and Kenneth Tu. *Performance Tuning and Optimizing ASP.NET*
1230    *Applications*, Apress, 2003.

1231    Killelea, Patrick. *Web Performance Tuning, 2d Ed*, O'Reilly & Associates, 2002.

1232    Larman, Craig and Rhett Guthrie. *Java 2 Performance and Idiom Guide*,
1233    Englewood Cliffs, N.J.: Prentice Hall, 2000.

1234    Shirazi, Jack. *Java Performance Tuning*, O'Reilly & Associates, 2000.

1235    Wilson, Steve and Jeff Kesselman. *Java Platform Performance: Strategies and*
1236    *Tactics*, Boston, Mass.: Addison Wesley, 2000.

1237    # Key Points

1238    ● Results of optimizations vary widely with different languages, compilers,
1239    and environments. Without measuring each specific optimization, you'll
1240    have no idea whether it will help or hurt your program.

1241    ● The first optimization is often not the best. Even after you find a good one,
1242    keep looking for one that's better.

1243    ● Code tuning is a little like nuclear energy. It's a controversial, emotional
1244    topic. Some people think it's so detrimental to reliability and maintainability
1245    that they won't do it at all. Others think that with proper safeguards, it's
1246    beneficial. If you decide to use the techniques in this chapter, apply them
1247    with care.