```
1   #include <x86intrin.h>
2   #define UNROLL (4)
3   #define BLOCKSIZE 32
4   void do_block (int n, int si, int sj, int sk,
5                     double *A, double *B, double *C)
6   {
7     for ( int i = si; i < si+BLOCKSIZE; i+=UNROLL*8 )
8       for ( int j = sj; j < sj+BLOCKSIZE; j++ ) {
9           __m512d c[UNROLL];
10          for (int r=0;r<UNROLL;r++)
11            c[r] =  _mm512_load_pd(C+i+r*8+j*n); //[ UNROLL];
12
13          for( int k = sk; k < sk+BLOCKSIZE; k++ )
14          {
15              __m512d bb = _mm512_broadcastsd_pd(_mm_load_sd(B+j*n+k));
16              for (int r=0;r<UNROLL;r++)
17                c[r] = _mm512_fmadd_pd(_mm512_load_pd(A+n*k+r*8+i), bb, c[r]);
18            }
19
20          for (int r=0;r<UNROLL;r++)
21            _mm512_store_pd(C+i+r*8+j*n, c[r]);
22        }
23    }
24
25  void dgemm (int n, double* A, double* B, double* C)
26  {
27    for ( int sj = 0; sj < n; sj += BLOCKSIZE )
28      for ( int si = 0; si < n; si += BLOCKSIZE )
29        for ( int sk = 0; sk < n; sk += BLOCKSIZE )
30          do_block(n, si, sj, sk, A, B, C);
31  }
```

**FIGURE 5.48  Optimized C version of DGEMM from Figure 4.80 using cache blocking.** These changes are the same ones found in Figure 5.21. The assembly language produced by the compiler for the do_block function is nearly identical to Figure 4.81. Once again, there is no overhead to call the do_block because the compiler inlines the function call.