

**Problem 1.1.1.** Write code for multiplying a pair of  $n \times n$  matrices ( $C \leftarrow AB$ ):

```
for  $i = 1, 2, \dots, n$  do  
  for  $j = 1, 2, \dots, n$  do  
     $c_{ij} \leftarrow 0$   
    for  $k = 1, 2, \dots, n$  do  
       $c_{ij} \leftarrow c_{ij} + a_{ik}b_{kj}$   
    end for  
  end for  
end for
```

Do this in your favorite interpreted language (MATLAB, Python, Ruby,  $\dots$ ). Use your code for multiplying a pair of  $100 \times 100$  matrices. If your language has a built-in operation for performing matrix multiplication, use this to compute the matrix product of your test matrices. How much faster is the built-in operation?

**Problem 1.1.2.** The matrix operation  $C \leftarrow C + AB$  on  $n \times n$  matrices can be written as three nested loops: The order of the loops can be changed; in fact, any order of “for i”, “for j”, and “for k” can be used

```
for  $i = 1, 2, \dots, n$  do
  for  $j = 1, 2, \dots, n$  do
    for  $k = 1, 2, \dots, n$  do
       $c_{ij} \leftarrow c_{ij} + a_{ik}b_{kj}$ 
    end for
  end for
end for
```

giving a total of  $3! = 6$  possible orderings. In a *compiled language* (such as Fortran, C/C++, Java, or Julia), time the six possible orderings. Which one is faster. Can you explain why?

**Problem 1.1.3.** A rule of thumb for high-performance computing is that when a data item is read in, we should use it as much as possible, rather than re-reading that data item many times and doing a little computing on it each time. In matrix multiplication, this can be achieved by subdividing each matrix into  $b \times b$  blocks. For  $n \times n$  matrices  $A$  and  $B$ , we can let  $m = n/b$  and write:

$$AB = \begin{pmatrix} A_{11} & A_{12} & \cdots & A_{1m} \\ A_{21} & A_{22} & \cdots & A_{2m} \\ \vdots & \vdots & \ddots & \vdots \\ A_{m1} & A_{m2} & \cdots & A_{mm} \end{pmatrix} \begin{pmatrix} B_{11} & B_{12} & \cdots & B_{1m} \\ B_{21} & B_{22} & \cdots & B_{2m} \\ \vdots & \vdots & \ddots & \vdots \\ B_{m1} & B_{m2} & \cdots & B_{mm} \end{pmatrix} = \begin{pmatrix} C_{11} & C_{12} & \cdots & C_{1m} \\ C_{21} & C_{22} & \cdots & C_{2m} \\ \vdots & \vdots & \ddots & \vdots \\ C_{m1} & C_{m2} & \cdots & C_{mm} \end{pmatrix} = C$$

Then for  $i, j = 1, 2, \dots, n$ ,  $C_{ij} \leftarrow \sum_{k=1}^m A_{ik}B_{kj}$ . As long as we can keep all of  $A_{ik}$ ,  $B_{kj}$ , and  $C_{ij}$  in cache memory at once, the update  $C_{ij} \leftarrow C_{ij} + A_{ik}B_{kj}$  can be done without any memory transfers once  $A_{ik}$  and  $B_{kj}$  have been loaded into memory. Write out a blocked version of the matrix multiplication algorithm and count the number of memory transfers with the blocked and original matrix multiplication algorithms.

**Problem 1.1.4.** Implement the original and unrolled inner product algorithms in Algorithm 1.1.3 as a function in your favorite programming language. Time the function for taking the product of two vectors of  $10^6$  entries. Use the simpler version to check the correctness of the unrolled version. Note that there may be differences due to roundoff error. Put the function call inside a loop to perform the inner product  $10^3$  times. Is there any difference in the times of the two algorithms? Note that interpreted languages, such as Python and MATLAB, may see little or no difference in timings. This is probably due to the fact that the time savings for the unrolled version is negligible compared to the overhead of interpretation. Also, interpreted languages will typically “box” and “unbox” the raw floating point value, converting it to and from a data structure that contains information about the type of the object as well as the object itself.

**Problem 1.1.5.** The following pseudo-code is designed to provoke a “stack overflow” error:

```
function OVERFLOW( $n$ )  
  if  $n = 2^k$  for some  $k$  then  
    print( $n$ )  
  end if  
  OVERFLOW( $n + 1$ )  
end function
```

Implement in your favorite programming language. How does it behave on your computer? How large a value of  $n$  is printed out before overflow occurs?

**Problem 1.1.6.** Dynamic memory allocation is memory allocation that occurs at run-time. It is an essential in all modern computational systems. Pick a programming language. How does this language allocate memory or objects? Do you need to explicitly de-allocate memory or objects in your programming language? How is this done?

**Problem 1.1.7.** There are different ways of automatically de-allocating unusable objects (*garbage collection*) in programming systems. In this exercise, we look at two of them. Describe reference counted garbage collection and “mark and sweep” garbage collection. What are their strengths and weaknesses?

**Problem 1.1.8.** *Memory leaks* are a kind of bug that can be hard to find and remove. These arise when memory is allocated for objects that are never de-allocated, and so eventually take up all available memory. Even if a system has garbage collection, this can still occur. Explain how this might happen.



**Problem 1.1.9.** Memory allocation and de-allocation can lead to fragmentation of the memory allocation system over time, so that there may be a great deal of memory available, but no large object can be allocated because the available memory is fragmented into small pieces. Read about and describe the SmallTalk double indirection scheme that allows SmallTalk to de-fragment the memory allocation system.