# Molecular modelling: scientific and technological boundaries

David A. Pensak

E. I. du Pont de Nemours & Company, Inc.

Abstract – Classical molecular modelling techniques are quite demanding
of computer time. Even the simplest molecular mechanics technique makes
demands which are proportional to the square of the number of atoms.
Molecular orbital calculations can grow at, depending on the type of
approximations made, up to the fifth power of the size of the molecule.
There are a number of ways in which this demand can be satisfied:
increased central processor speed, new computer architectures (including
parallel processing), and fundamentally new algorithms. Each will be
considered in the context of how accessible they are (and will be) to
practising chemists, who, after all, are the very reason for modelling
research to be undertaken.

Molecular Modelling is a generic term that is used to refer to virtually anthing which is
done to depict, describe, or evaluate any aspect of the properties or structure of a molecule
that requires the use of a computer. For the purposes of this discussion, we will restrict
ourselves to a computation of molecular geometries and properties. Computer graphics will
be incorporated only to the extent that it can be used to illustrate the results of the
computations.

While molecular modelling is promoted by some as having no limits on its ability to impact
experimental chemistry, there are some quite definite boundaries that will be with us for
the next several years. They fall into three separate and distinct categories: inadequate
scientific understanding, non-optimal algorithmic implementation, and insufficient computer
time.

Since you cannot compute that which you do not understand (except in a very few machine
learning environments), those chemical problems where either the mechanism or the process is
known cannot be computed with any significant accuracy. As but one example of this, consider
the problem of the class of photosystem type II herbicides, which function by interception
of an electron during the photosynthesis process. There are a wide range of molecules which
have the appropriate electrode reduction potential to do this. Attempts to correlate
reduction potential with different computed parameters (such as HOMO-LUMO gap, LUMO energy,
etc.) have been moderately successful in closely related families and totally unsuccessful
between structurally different systems. Further, there are virtually no chemists who can
look at a structure and predict its reduction potential.

There are multiple ways of programming essentially any calculation. Some are more efficient
than others. As a flagrant example, there is a text string searching problem where the best
known algorithm runs faster on a TRS-80 personal computer than some other algorithms run on
a Cray X-MP supercomputer. In more practical cases, simple replacement of some of the more
computationally intensive modules from some scientific programs with ones from professionally
developed and maintained libraries have resulted in execution time improvements of factors
of ten or more.

A classical problem that all molecular orbital theorists have to confront is diagonalization
of real symmetric matrices to find the eigenvalues and eigenvectors (energy levels and
orbitals, respectively). Over the last twenty years techniques for this have improved by
more than a factor of ten and are more stable to degenerate vectors as well. Within the
last year however, a fundamental advance has been made by Dongarra at Argonne National
Laboratories which should permit significant additional speed ups on both serial and parallel
machines.

Optimization consists of finding the point or points in n-dimensional space which result in
the maximum (or minimum) value of a defined function such as total strain energy. If the
evaluation of energy is extremely fast, efficiency makes very little overall difference. If,
however, the energy is expensive to compute (such as an ab-initio total energy) it becomes

very important to not have to evaluate the energy more times than absolutely necessary. Recent new algorithms, such as those of Karmarkar at Bell Telephone Laboratories, hold hope for quite considerable improvements.

There are, however, some problems which, even though efficiently coded, require more CPU time than is now available (or is likely to be available in the coming few years). There are multiple possible approaches (described below) to deal with this.

Classical molecular modelling techniques can be extremely demanding of computer time. Even the simplest molecular mechanics technique makes demands which are proportional to the square of the number of atoms. Molecular orbital calculations can grow at, depending on the type of approximations made, up to the fifth power of the size of the molecule. There are a number of ways in which this demand can be satisfied: increased central processor speed, new computer architectures (including parallel processing), and fundamentally new algorithms.

Each will be considered in the context of how accessible they are (and will be) to practising chemists, who, after all, are the very reason for modelling research to be undertaken.

The personal computer of today is roughly comparable in power to the mainframe of 5-7 years ago and the supercomputer of 10-12 years ago. One can assume roughly a 30% improvement in price performance ratio every year. Mini-supercomputers with a speed of almost 100 million instructions per second are becoming available for less than $100,000. Real supercomputers of greater than one billion instructions per second are commercially available today (though they cost more than $10,000,000). The fundamental nature of how we perform scientific discovery will radically change as speeds in this range (and greater) become readily available to chemists.

Reduced instruction set computers (RISC) are just beginning to come onto the market. They reverse the trend of the last decade to make ever more complex sets of instructions that the central processor can utilize. It has now been shown that significantly greater price performance can be achieved by using only a very limited set of instructions and constructing the more complex ones from them. In 1988, RISC machines can run at 10 million instructions per second (MIPS) but by 1990 desk top workstations capable of speeds of at least 50 MIPS will be available that will cost less than $10,000.

Traditional computers have just one central processor, which can execute one instruction at a time. A growing trend is the generation of computers with multiple processors which execute in parallel with each other. Very few problems will be able to take advantage of multiple processors to the fullest degree, but significant improvements are attainable. For the next few years, the rate limiting steps will be the recognition of problems which can practically benefit from parallelization and then the development of useful algorithms to exploit the parallelism.

There are several different ways that parallel computers can be constructed. They can all share the same physical memory, in which case they are called tightly coupled. They can each have their own memories and communicate only on specific program command, in which they are called loosely coupled. They can also be separate processors but the communication between them is tightly controlled by the hardware only. These are called systolic processors and may be thought of as similar to the human circulatory system where on each beat of the heart, all the blood cells move, regardless of where they are or what they are doing.

Even when specialized computer hardware is not available, the very exercise of thinking about how a problem should be structured to fit effectively on a parallel machine can lead to great benefits. For example, one heavily used molecular orbital plotting program was analysed to see how to restructure it to fit effectively on a parallel processor. During the course of this analysis, a new algorithm was discovered which runs approximately 100 times faster than the old one (on the same hardware). The system studied was the plotting program of Jorgensen and Salem. Since it had been originally designed almost 15 years ago when memory was quite expensive, design decisions had been made that were far from optimal in today's computational climate. Instead of cutting the molecule into a number of horizontal and vertical slices and then contouring the electron density in each slice, individual processors are assigned to work on individual atoms. This avoids any computation on spatial points which are too far removed from an atom to be a significant contributor. An additional benefit is that the symmetry properties of each atomic orbital can be explicitly considered without regard to the coordinate system of the slices being cut from the molecule.

The very way in which we use computed data can change as a function of how rapidly we can retrieve it. For example, a large database of single crystal X-ray structures was originally distributed in a relatively computer independent manner. When it was re-written to take maximal advantage of the hardware, it sped up by a factor of over 2000. What makes this result unusual was the way that the scientists began to use the data (often in ways we had not anticipated). We will consider how user expectations of computational speed effects both the type and nature of the calculations attempted.

The goal of this research is to develop computational techniques, both hardware and software, that will permit interaction with models (both physical and computational) in a manner which will greatly increase the insight that can be developed in support of experimentation.

Closely associated with the computer hardware research has been an effort to develop graphical representations which will convey greater chemical information in a single image than the static pictures that we are used to seeing. For example, it is possible to show in a single frame all the conformational states of a flipping cyclohexane, complete with energetics. The difficulty is that chemists are not comfortable looking at these pictures so their information transmittal value is limited. We have two choices then, either to 'reprogram' the experimentalists so that they can appreciate the increased information content of these images, or radically change the nature of the presentation. The ramifications of both are being explored.

When we are actually able to increase our abilities to compute by two orders of magnitude (and more) the number and type of questions that we will be able to ask will fundamentally change. Model studies will be considerably more realistic and perhaps more audacious systems will be pursued. This will not, however, cause a revolution in experimental chemistry. It will stimulate more rigorous, and at the same time, more imaginative explorations. It may also increase efficiency, as it will become increasingly possible to gain meaningful data from experiments that fail (as well as from those that succeed).

One of the fundamental problems that we have no idea how to overcome is how to teach individuals to effectively use the wealth of new information that the computer will make available to them. The brain has never had to learn how to think and communicate in dimensions higher than three. Consider data which can be described as a spiral in 3 dimensions. It is easy to make a mental picture of that. Similarly, it is simple to visualize data which lies in a plane in 3 dimensions. If, however, the data is 6 dimensional and 3 dimensions are planar and 3 are in spiral, we are totally unable to visualize it. Even if we could make a mental picture, we are completely incapable to communicating it to other scientists. To what degree we will be able to retrain ourselves is not clear.

### BIBLIOGRAPHY

1. Color and the Computer. Edited by H. John Durrett

2. Computer Graphics: A Programming Approach, by Steven Harrington

3. Parallel Computing: Theory and Comparisons, by G. Jack Lipovski and Marioslaw Malek

4. Quantum Pharmacology, by W.G. Richards