Parallelization of the QC-lib Quantum Computer Simulator Library

Ian Glendinning and Bernhard Ömer

VCPC

European Centre for Parallel Computing at Vienna Liechtensteinstraße 22, A-1090 Vienna, Austria http://www.vcpc.univie.ac.at/qc/

Abstract. We report on work to parallelize QC-lib, a C++ library for the simulation of quantum computers at an abstract functional level. After a brief introduction to quantum computing, we give an outline of QC-lib, then describe its parallelization using MPI, and present performance measurements made on a Beowulf cluster. Using more processors allowed larger problems to be solved, and reasonable speedups were obtained for the Hadamard transform and Grover's quantum search algorithm.

1 Introduction

Quantum computers are devices that process information using physical phenomena unique to quantum mechanics, and which have the potential to be able to solve certain problems such as prime factorization spectacularly faster than any conventional computer [1]. In a classical computer the basic unit of information is the bit, a two-state device that can represent the values 0 and 1. The quantum analogue of the bit is a two-state quantum system, such as an electron's spin or a photon's polarization, which has come to be known as a qubit. The difference between a qubit and a bit is that a qubit can exist not only in the states 0 and 1, but also in a mixture of both of them, called a superposition state. Furthermore, whereas a register of n bits can be in any one of 2^n states, storing one of the numbers 0 to $2^n - 1$, a register of n qubits can be in a superposition of all 2^n states, and a function applied to a quantum register in a superposition state acts on all 2^n values at the same time! This is known as quantum parallelism, and it is one of the key ingredients in the power of quantum computers.

Unfortunately, when a quantum register in a superposition state is measured, the result obtained is only one of the 2^n possible values, at random. However all is not lost, as the probabilities of measuring the different values can be manipulated by operating on a quantum register with quantum gates, which are the quantum analogue of logic gates. Quantum algorithms consist of sequences of quantum gate operations and optionally measurements, and it turns out that algorithms exist that are able to exploit quantum parallelism, and to leave an output register in a state where the probability of obtaining the value that is the answer to the problem is very close to one, giving an advantage over classical algorithms.

However, building quantum computers is a huge technological challenge, and quantum computing hardware is not currently available outside physics research labs, so simulators present an attractive alternative for experimenting with quantum algorithms. Furthermore, they offer the only way to run programs on more than seven qubits, which is the current state of the art in experimental hardware. Simulators also help debugging of quantum programs, allowing direct examination of the quantum state, which is not possible in physical quantum computers. Simulators suffer from a problem, which is that their execution time and memory requirements increase exponentially with the number of qubits. Parallelization alleviates this problem, allowing more qubits to be simulated in the same time or the same number to be simulated in less time. Many simulators exist, but few for parallel systems. Niwa et al. [2] describe one and review related work.

2 Qubits, Registers and Gates

The state of a qubit can be represented by a two-dimensional complex vector of length 1. The states that are the quantum analogues of 0 and 1 are called the computational basis vectors, and they are written $|0\rangle$ and $|1\rangle$, in a notation due to Dirac. In terms of vectors, they are conventionally defined to be

$$|0\rangle = \begin{pmatrix} 1\\0 \end{pmatrix}, \ |1\rangle = \begin{pmatrix} 0\\1 \end{pmatrix} \ , \tag{1}$$

and a general qubit state is

$$\alpha|0\rangle + \beta|1\rangle = \begin{pmatrix} \alpha \\ \beta \end{pmatrix} , \qquad (2)$$

where α and β are complex numbers called *amplitudes*. Measurement of the state always gives either $|0\rangle$, with probability $|\alpha|^2$, or $|1\rangle$, with probability $|\beta|^2$, which is consistent with the *normalization* condition that the vector's length is 1, which is $|\alpha|^2 + |\beta|^2 = 1$.

The state of an n-qubit register can be represented as a 2^n -dimensional complex vector of length 1. If we call the ith basis state $|i\rangle$, where $0 \le i \le 2^n - 1$, then a general n-qubit state has the form

$$\sum_{i=0}^{2^n-1} \alpha_i |i\rangle , \qquad (3)$$

where α_i is the *i*th complex component of the vector representing the state, $|\alpha_i|^2$ is the probability that measurement will give the value *i*, and the normalization condition is $\sum_{i=0}^{2^n-1} |\alpha_i|^2 = 1$. The numbers labelling the basis states are often written in binary, to show the value of each component qubit in the register. For example, the computational basis vectors for a two qubit register are

$$|00\rangle = \begin{pmatrix} 1\\0\\0\\0 \end{pmatrix}, \ |01\rangle = \begin{pmatrix} 0\\1\\0\\0 \end{pmatrix}, \ |10\rangle = \begin{pmatrix} 0\\0\\1\\0 \end{pmatrix}, \ |11\rangle = \begin{pmatrix} 0\\0\\0\\1 \end{pmatrix} \ . \tag{4}$$

Any *n*-qubit gate (operator) can be represented as a $2^n \times 2^n$ unitary matrix, i.e. a complex matrix U with the property that $U^{\dagger}U = I$. The operation of a gate on a quantum register is implemented by matrix multiplication. The only non-trivial classical single-bit gate is the NOT gate, but there are many non-trivial single-qubit gates, for example the Hadamard gate:

$$H = \frac{1}{\sqrt{2}} \begin{pmatrix} 1 & 1\\ 1 & -1 \end{pmatrix} . \tag{5}$$

This gate is useful because applying it to either of the basis states produces an equal mixture of both of them: $H|0\rangle=\frac{1}{\sqrt{2}}(|0\rangle+|1\rangle)$ and $H|1\rangle=\frac{1}{\sqrt{2}}(|0\rangle-|1\rangle)$.

The prototypical multi-qubit gate is the *controlled*-NOT or CNOT gate. It has two inputs, known as the *control* and *target* qubits, and two outputs. If the control qubit is set to 0, the target qubit is unchanged, and if the control qubit is set to 1, the target qubit is flipped $(|c,t\rangle)$:

$$|00\rangle \to |00\rangle; \ |01\rangle \to |01\rangle; \ |10\rangle \to |11\rangle; \ |11\rangle \to |10\rangle \ .$$
 (6)

3 QC-Lib

QC-lib is a C++ library for the simulation of quantum computers at an abstract functional level [3], and it is used as the back end of an interpreter for the QCL Quantum Computation Language [4]. Its main features are basis vectors of arbitrary length (not limited to word length), efficient representation of quantum states using hash tables, nesting of substates and arbitrary combinations of qubits, composition and tensor product of operators (gates), and easy addition of new operators using class inheritance. The top-level classes in QC-lib are:

bitvec - arbitrary length bit vectors which represent basis states
term a basis vector with a complex amplitude
termlist - a list of terms: the internal representation of a quantum state
quState - user class for quantum states
opOperator - user class for quantum operators

The data structure of class termlist is a linear array in combination with a hash table. Only terms with non-zero amplitudes are stored, and the array and hash table are dynamically doubled in size if the array fills up.

The class quState contains two sub-classes, quBaseState which contains actual state information, and quSubState which represents substates. An object of class quBaseState represents the state of the whole quantum memory, and the quSubState class can be used to allocate subregisters. A quBaseState object contains two termlist objects. One contains the terms in the current state, and the other is a term buffer to accumulate the result of an operation on the state.

A sub-class of opOperator called opMatrix implements the representation of an n-qubit operator as a $2^n \times 2^n$ complex matrix, storing the non-zero elements of each row in an array of lists. However, most operators are simpler, only working

on a few qubits, or substituting one basis vector for another, with or without a phase factor, and op@perator has sub-classes for a number of such special cases, such as permutation of qubits, quantum functions, general single-qubit operations, and the CNOT gate.

The following example program uses QC-lib to implement the Hadamard transform, which is the Hadamard gate applied to every qubit in a register:

4 Parallelization

The parallelization is being carried out in the message-passing style of programming, using the MPI message passing interface. The parallelization strategy is to distribute the representation of the quantum memory. Each processor stores a subset of the terms, and operators are applied only to local terms. The result of a local operation in general includes terms that are not owned by the processor, which must be communicated to their owning processors. The program is SPMD, with each processor running a copy of the same code.

The data distribution scheme for terms is that on 2^n processors, the n least significant qubits of their basis states are interpreted as the number of the processor that owns them. For example, the processor allocation of basis states for a four-qubit register on the four processors P_0-P_3 is:

P_0	$ 0000\rangle$, $ 0100\rangle$, $ 1000\rangle$, $ 1100\rangle$
P_1	$ 0001\rangle, 0101\rangle, 1001\rangle, 1101\rangle$
P_2	$ 0010\rangle, 0110\rangle, 1010\rangle, 1110\rangle$
P_3	$ 0011\rangle$, $ 0111\rangle$, $ 1011\rangle$, $ 1111\rangle$

4.1 Communication Pattern for Single-Qubit Operators

Consider a general single-qubit operator $U=\begin{pmatrix}u_{11} & u_{12} \\ u_{21} & u_{22}\end{pmatrix}$ operating on a single qubit with state $\alpha|0\rangle+\beta|1\rangle$. For simplicity, assume that there are just two

processors, and the qubit in question is the least significant one, and so determines the data distribution. After the operation of U locally on each processor, terms are created that are not owned by the processor, and so communication is necessary. Specifically, $\beta u_{12}|0\rangle$ has to be sent from $P_1 \to P_0$ and $\alpha u_{21}|1\rangle$ from $P_0 \to P_1$:

$$\begin{array}{c|c}
P_0 & \overline{\alpha|0\rangle} \\
P_1 & \overline{\beta|1\rangle} & \xrightarrow{U} & \overline{\alpha U|0\rangle} \\
P_1 & \overline{\beta U|1\rangle} & = \overline{\begin{array}{c}
\alpha u_{11}|0\rangle + \alpha u_{21}|1\rangle} \\
\overline{\beta u_{12}|0\rangle + \beta u_{22}|1\rangle} & \xrightarrow{\text{Comm}} \overline{\begin{array}{c}
(\alpha u_{11} + \beta u_{12})|0\rangle} \\
\overline{(\alpha u_{21} + \beta u_{22})|1\rangle}
\end{array}}$$

When a single-qubit operator is applied to the ith qubit of a register, it is applied to the ith qubit of every term in the superposition state of the register, leaving the other qubits in each term unchanged. For each term in the initial state, at most two terms are therefore created (if one or more of the u_{ij} are zero, less terms will be created). If the ith qubit is not one of the qubits that determines the data distribution, then no communication is necessary, as both new terms are owned by the original processor. Otherwise, one of the new terms is locally owned and the other one is remotely owned and must be communicated. For each processor, all the remotely owned terms are owned by the same other processor, as a single bit has been flipped in the distribution key. Symmetrically, the remote processor potentially creates terms that are owned by the local one. In general therefore, bi-directional communication is needed between disjoint pairs of processors.

If the basis vectors of an n-qubit register are thought of as the coordinates of the corners of an n-dimensional hypercube, such that the ith qubit represents a coordinate of 0 or 1 in the ith dimension, then the communication pattern generated by an operation on the ith qubit can be visualized as being along parallel edges of the hypercube with data movement in its ith dimension.

4.2 Parallelization of Single-Qubit Operators

The representation of the distributed quantum memory has been encapsulated in class quBaseState, without altering the operator invocation mechanism in the sequential version of QC-lib. In the parallel version, each quBaseState object has a second term buffer (a termlist object), to accumulate terms destined for remote processors. The local and remote terms that result from the action of single-qubit operators are accumulated separately, and when the accumulation is complete, the remote term buffer is sent to the remote process, and reciprocally, the buffer that the remote process sent is received. This exchange of data between pairs of processes is implemented using the MPI function MPI_Sendrecv(). The received terms are then merged into the local term buffer. Finally, the term buffer is swapped with the termlist for the current state, so that the new terms become the current state, and the old one is ready to be used to accumulate the result of the next operation.

Currently only a few of the operator sub-classes of op0perator have been parallelized, but they include the general single-qubit operator and the CNOT gate, which together are universal for quantum computation, though they don't necessarily offer the most efficient way of achieving a particular operation.

5 Performance Measurements

Performance measurements were made on a Beowulf cluster with 16 compute nodes, each having a 3,06 GHz Pentium 4 processor and 2 GByte of 266 MHz dual channel DDR-RAM, running Linux 2.4.19. The nodes were connected by two independent networks, Gigabit Ethernet for message-passing communication, and Fast Ethernet for job submission, file access, and system administration.

Figure 1 shows how the run time for the Hadamard transform varies with the problem size (number of qubits) for different numbers of processors. The run time for the sequential code is also shown, and it can be clearly seen from the semi-logarithmic plot that it increases by a constant factor for every extra qubit. Each extra qubit corresponds to a doubling of the problem size, which is consistent with the slope of the line, which corresponds to an increase in run time by a factor of 2.05 for each extra qubit. The parallel version of the program run on one processor has very little overhead compared with the sequential version, which is not surprising, as it does not perform any communication. When larger numbers of processors are used for smaller problem sizes there is more of an overhead, but as the problem size increases, the scaling behaviour of the parallel code becomes similar to that of the sequential version.

Figure 2 shows how the speedup varies with the number of processors, for various problem sizes. The speedup is relative to the run time of the sequential code for the same problem size, but since the sequential code ran out of memory for the larger problem sizes, the sequential run time for those cases had to be estimated, which was done by by linearly extrapolating a least squares fit to the the logarithm of the run time versus the number of qubits. For small problem sizes there is actually a decrease in speedup for larger numbers of processors, which is to be expected for any problem size if enough processors are used. For larger problem sizes reasonable speedups were obtained, up to 9.4 on 16 processors for the largest problem size.

The Hadamard transform is a sufficiently simple operation, that it was possible to make make runs that used all the system's available memory within a reasonable run time. It was found that both the sequential version of the program and the parallel version on one processor could simulate a maximum of 25 qubits, and that each doubling of the number of processors increased the maximum number of qubits by one, up to a maximum of 29 qubits for all 16 processors. This is consistent with the amount of memory needed by the program to store quantum states, which in the sequential version of the code, for quantum states with no more than 2^{32} terms, is 64 bytes per term. The parallel version of the code contains an extra termlist to buffer terms that need to be communicated, which means that it potentially needs up to 50% more memory, but

in the Hadamard transform very few terms need to be communicated, as only changes to the least significant qubits cause communication, and these are processed first, but large numbers of terms are not created until the most significant qubits are processed. Operating on n qubits, the Hadamard transform ultimately produces a state with 2^n terms, so the amount of memory needed to represent it is approximately $64(2^n) = 2^{n+6}$ bytes, and 25 qubits require 2 GByte, which was the amount of memory available. Not all algorithms use such a large number of terms with non-zero amplitudes, and as only these are explicitly represented by QC-lib, some algorithms can be run for still larger numbers of qubits.

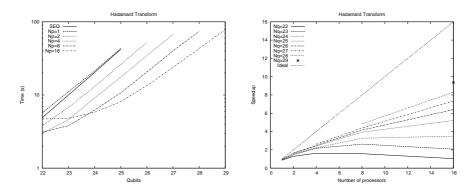


Fig. 1. Run time for Hadamard transform Fig. 2. Speedup for Hadamard transform

Grover's quantum search algorithm [5] was also implemented using QC-lib, and the run time was measured for the portion of the algorithm up to the point when the first measurement is made. This was done to allow meaningful comparison of run times, as in the full version of the algorithm, there is a finite probability that the search will fail, and that it has to be repeated, so the run time is non-deterministic. Figure 3 shows how the run time for Grover's algorithm varies with the problem size, for different numbers of processors. The behaviour is similar to that in the case of the Hadamard transform, but in this case the slope of the line for the sequential run time corresponds to an increase in run time by a factor of 3.10 for each extra qubit. This is partly accounted for by the fact that the number of iterations in the main loop of the algorithm increases by a factor of $\sqrt{2}$ for each extra qubit, and taken together with the doubling of the problem size, that would imply an increase in run time by a factor of 2.82.

Figure 4 shows how the speedup varies with the number of processors, for various problem sizes. Significantly better speedups are obtained than for the Hadamard transformation. For just 17 qubits, which is less than the minimum number of qubits considered for the Hadamard transform, a speedup of 11.2 was obtained on 16 processors, and it is to be expected that with more qubits the performance would be even better.

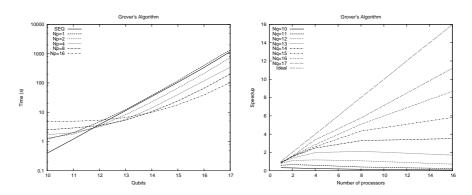


Fig. 3. Run time for Grover's algorithm

Fig. 4. Speedup for Grover's algorithm

6 Conclusion

Sufficient functionality of QC-lib has been implemented to simulate universal quantum computation. We have implemented the Hadamard transform and Grover's algorithm using our library, and have made performance measurements for these codes. Promising speedups were obtained. Future work will include the implementation of static control of the distribution of qubits by the programmer, more operators, implementation of Shor's prime factorization algorithm [6], and the investigation of its performance. In the longer term, dynamic redistribution of qubits will be implemented and a load-balancing strategy will be developed.

Acknowledgements

This work was partially supported by the Special Research Program SFB F011 "AURORA" of the Austrian Science Fund FWF.

References

- 1. Nielsen, M.A., Chuang, I.L.: Quantum Computation and Quantum Information. Cambridge University Press (2000)
- 2. Niwa, J., Matsumoto, K., Imai, H.: General-purpose parallel simulator for quantum computing (2002) http://arXiv.org/abs/quant-ph/0201042.
- 3. Ömer, B.: Simulation of quantum computers (1996) http://tph.tuwien.ac.at/~oemer/doc/qcsim.ps.
- 4. Ömer, B.: Quantum Programming in QCL. Master's thesis, Vienna University of Technology (2000) http://tph.tuwien.ac.at/~oemer/doc/quprog/index.html.
- Grover, L.K.: A fast quantum mechanical algorithm for database search. In: Proc. of the 28th annual ACM Symposium on the Theory of Computation (Philadelphia, Pennsylvania), New York, ACM Press (1996) 212–219
- 6. Shor, P.W.: Polynomial-time algorithms for prime factorization and discrete logarithms on a quantum computer. SIAM J. Comp. **26** (1997) 1484–1509