



SGI BIOINFORMATICS PERFORMANCE REPORT

Summer 1999

INTRODUCTION.....	2
AMDAHL'S LAW: EVALUATING PARALLEL PERFORMANCE	3
FASTA 3.2.....	5
FASTA 3.2: Single Processor Performance.....	6
FASTA 3.2: Parallelism and Throughput.....	7
FASTA 3.2: Throughput	8
CLUSTAL W 1.74	9
CLUSTAL W: Single Processor Performance.....	10
CLUSTAL W: Parallel Speedups	11
BLAST.....	12
High-Throughput BLAST (HT-BLAST).....	12
HT-BLAST: Uniprocessor Performance	13
HT-BLAST: Parallel Speedup	14
SMITH-WATERMAN	15
Smith-Waterman on SGI Origin 2000.....	16
D ² _CLUSTER v1.21	17
SRS PARALLEL PERFORMANCE.....	18
PROTEOMICS: HIGH-THROUGHPUT SEQUEST	19
HT-SEQUEST: Relative Uniprocessor Performance.....	20
HT-SEQUEST: Parallel Speedups.....	21
APPENDIX: SGI 1400	22
FASTA on SGI 1400L	23
BLAST on SGI 1400L	24
BLAST on SGI 1400L	25
ACKNOWLEDGMENTS	26

Introduction

The SGI™ Origin™ 2000 and SGI™ Origin™ 200 servers exhibit excellent performance and scalability on genomic sequence searching algorithms like BLAST, FASTA, Smith-Waterman, CLUSTAL W, d2_cluster and SRS. The high performance and scalability of these codes stem from the MIPS® R12000™ and MIPS R10000® processors and the supporting memory and I/O subsystems of the SGI™ Origin™ servers, combined with highly tuned codes that maintain good locality of data.

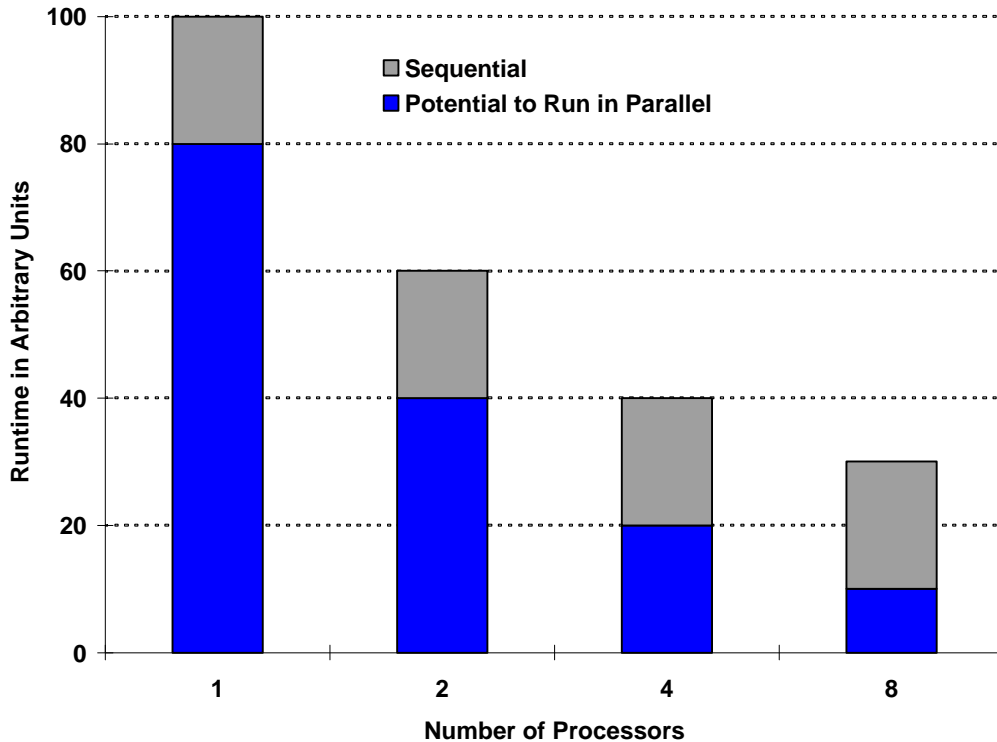
Computational performance can be evaluated from two perspectives: throughput performance addresses the need to process a large batch of jobs in the shortest amount of time, and turnaround performance, which measures how quickly one individual job is completed. Delivering high-throughput and fast-turnaround performance for the full range of bioinformatics applications is a key goal of the SGI™ team of bioinformatics professionals. This report gives examples of how SGI Origin systems meet both types of computational demands.

For detailed technical support information on BLAST, FASTA and CLUSTALW running on SGI systems please visit the SGI public web site: <http://www.sgi.com/chembio/resources> You will find general information, technical information, porting notes, and known problems and fixes.

For more information on how SGI can help you solve your computational problems, please contact your local SGI sales representative or distributor, visit our Web site at <http://www.sgi.com/chembio> or contact us at chembio@sgi.com

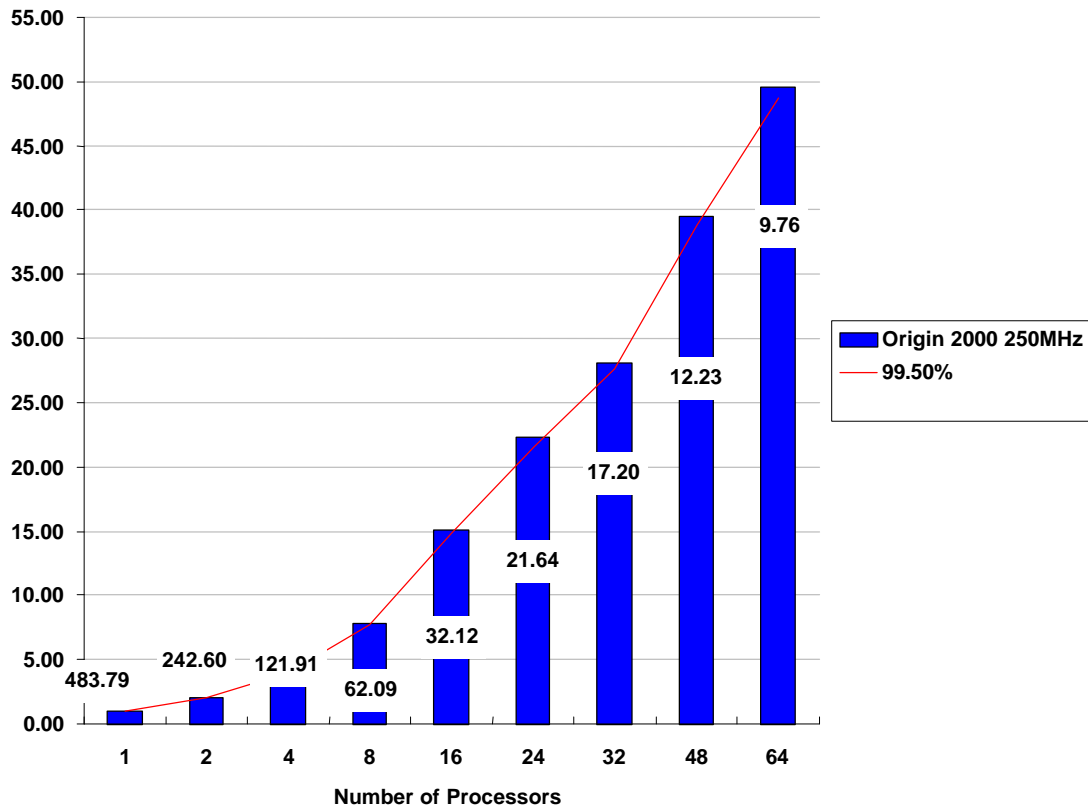
Amdahl's Law: Evaluating Parallel Performance

Amdahl's Law defines the increase in speed that can be gained by parallelizing an application. It states that performance improvements due to parallelization are limited by the fraction of the code that is not running in parallel. For example, if 80% of the run time of a hypothetical application can be run in parallel (blue or dark bars), then as more processors are added, the non-parallel portion of the code (red striped or gray bars) dominates the performance. When running this application with eight processors, the non-parallel or serial run time dominates the overall run time and very little, if any, speed increase can be achieved by using additional processors.



Amdahl's Law is of utmost importance in understanding and gauging parallel performance; thus, included on each of the following parallel speedup graphs is the theoretical performance curve predicted by Amdahl's Law. Large deviations from this curve could indicate performance anomalies.

Here, the example from a BLAST run shows how the measured speedup (bars) closely matches the theoretical speedup (curve, corresponding to 99.50% parallelism), indicating that no extraneous hardware or system software events negatively affect the parallel performance of BLAST on the SGI Origin 2000 system.



FASTA 3.2

FASTA 3.2 is a collection of sequence comparison programs for biological sequence databases. These programs are widely used in bioinformatics, and can be used to search sequence databases, evaluate similarity scores, and identify periodic structures based on local sequence similarity. One of the programs in FASTA 3.2 is FASTA, which uses the FASTA algorithm (see references below) to compare a protein sequence query to a protein sequence library or a DNA sequence query to a DNA sequence library.

Another program in FASTA 3.2 is FASTX, which uses the FASTX algorithm to compare a DNA sequence query to a protein sequence library, translating the DNA sequence in three frames and following frame shifts in the alignment.

The following performance graphs demonstrate the excellent throughput performance of FASTA and the parallel scalability of FASTX on the SGI Origin server.

References

W. R. Pearson and D. J. Lipman (1988) Proc. National Academy of Science, USA, Improved tools for biological sequence comparison. 85:2444-2448

W. R. Pearson (1996), Methods Enzymol, Effective protein sequence comparison. 266:227-258

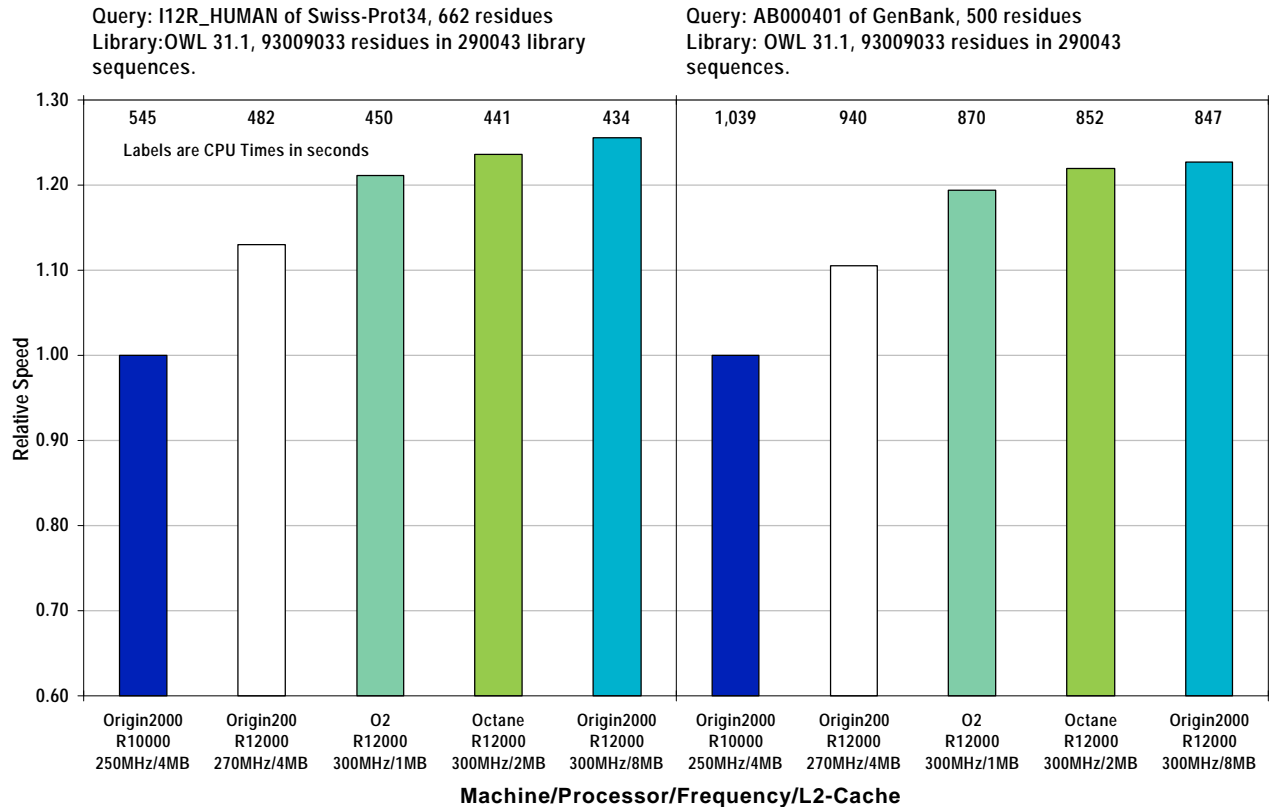
FASTA 3.2: Single Processor Performance

In the following graphs the relative performance of several computers is shown using two different executables from FASTA (version 3.2, revision 01), with both cases using a k-value of 1. Also in both cases the library searched is as follows: OWL 31.1, 93009033 residues in 290043 sequences.

The graph on the left uses FASTA3 with the following query sequence:
Sequence |I2R_HUMAN of Swiss-Prot34, 662 residues.

The graph on the right shows the performance of fastx3 with the query sequence:
Sequence AB000401 of GenBank, 500 residues.

The SGI Origin 2000 with MIPS R12000 300 MHz processors runs these cases more than 1.2X faster than the same system equipped with MIPS R10000 250 MHz processors.



FASTA 3.2: Parallelism and Throughput

Two different ways of exploiting the Origin system's multiprocessing flexibility and outstanding system properties are shown in the next two pages.

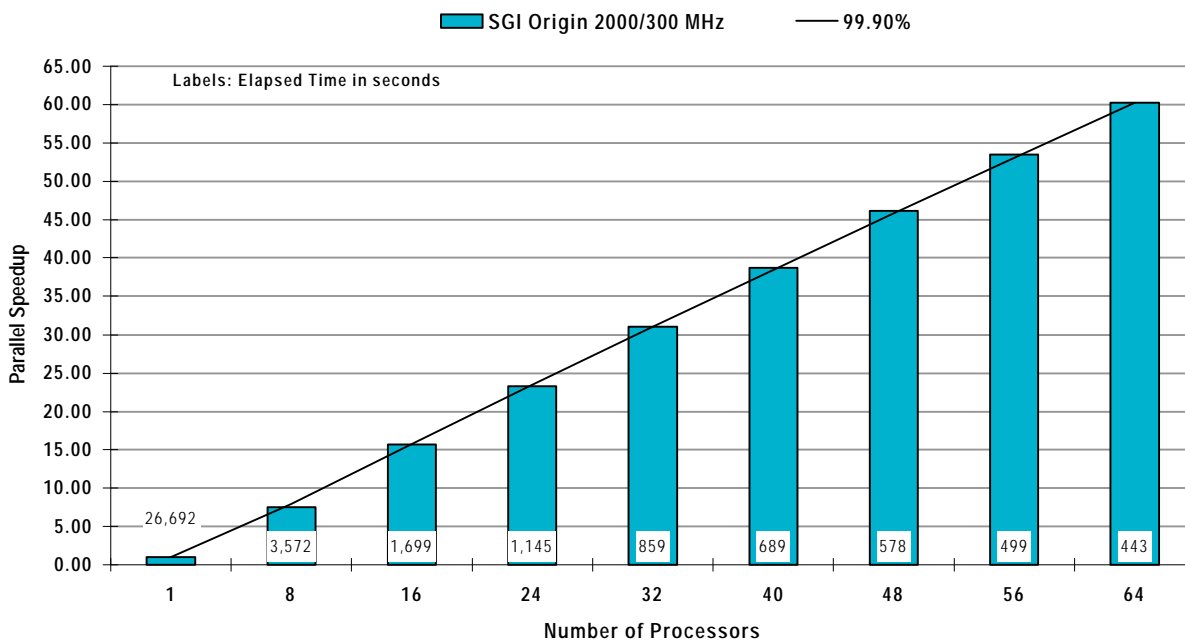
FASTA 3.2 Revision 01 (k=1), fastx3_t

QUERY SEQUENCE: Sequence AB002150 of GenBank, 15,684 residues

SEARCH LIBRARY: OWL 31.1, 93,009,033 residues in 290,043 library sequences

Parallel Speedup

Good parallel performance on a multiprocessor system allows the user to achieve fast turnaround time of a single job, enabling runs to be made in a much shorter amount of time than on a single-processor system. The benchmark in this graph shows the excellent parallel scalability that can be achieved by a single FASTX job on a 64-CPU SGI Origin 2000 system. The results show how the solution time of a single processor FASTX job can be reduced from almost seven and a half hours to just under seven and a half minutes by using 64 R12000 processors running at 300 MHz. For the scientist, this decreased turnaround time offers better interactivity for quickly investigating specific problems and making faster research decisions.



FASTA 3.2: Throughput

Revision 01 (k=6), fasta3

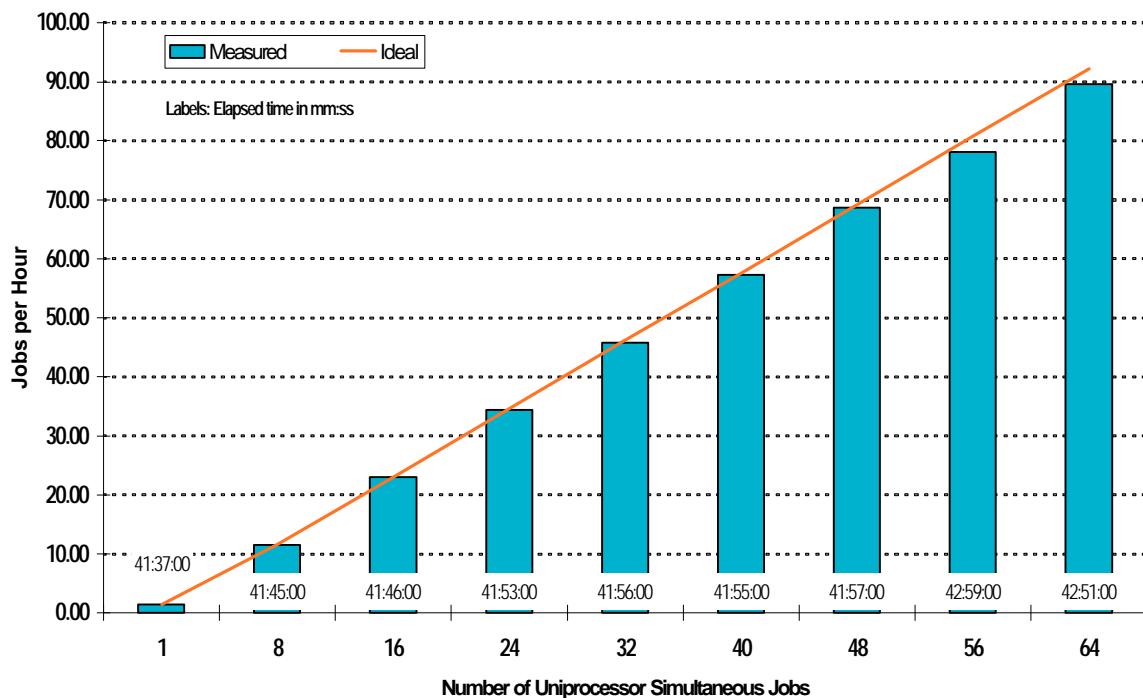
Source modified to support memory mapped databases

QUERY SEQUENCE: Sequence AB000401 of GenBank, 500 residues.

SEARCH LIBRARY: GenBank 111.0, 2,562,214,506 residues in 3,504,701 sequences.

Throughput

Good throughput efficiency on a multiprocessor system allows the user to achieve fast turnaround times of multiple jobs, enabling many more runs to be made in almost the same time as a single run. The benchmark below shows the excellent throughput performance that can be achieved by launching multiple FASTA jobs on a 64-CPU SGI Origin 2000 system equipped with MIPS R12000 300 MHz processors. The results show that by simply using more processors, the number of queries that can be processed per hour increases in near direct proportion to the number of processors applied. For the scientist, this increased throughput rate offers higher productivity for processing and analyzing many sequences on a daily basis.



CLUSTAL W 1.74

CLUSTAL W is a widely used, multiple sequence alignment program for biological sequences. The program uses the CLUSTAL W algorithms (see reference below) to progressively align multiple protein or DNA sequences.

SGI has developed a parallel version of CLUSTALW, which performs the pairwise alignments and guide tree formation using multiple processors. For large numbers of sequences (greater than 500 sequences), computations are the most time-consuming in CLUSTAL W.

Reference

Thompson, J.D., Higgins, D.G., and Gibson, T.J. (1994) CLUSTAL W: Improving the sensitivity of progressive multiple sequence alignment through sequence weighting, positions-specific gap penalties and weight matrix choice. *Nucleic Acids Research*. 22:4673-4680.

CLUSTAL W: Single Processor Performance

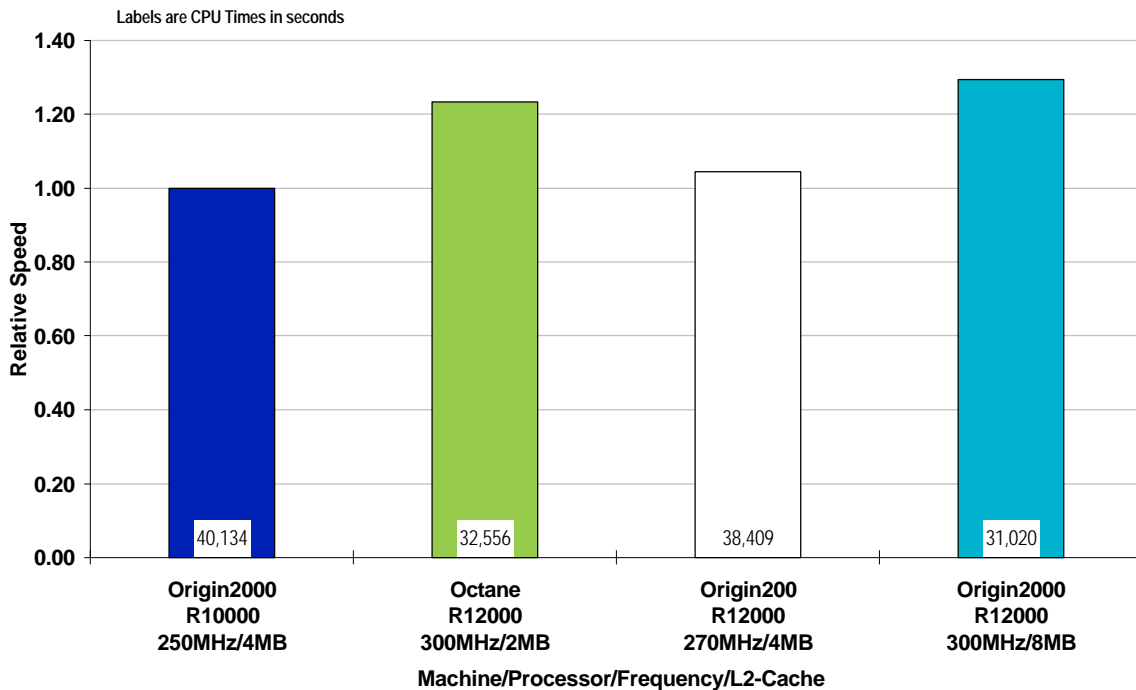
CLUSTAL W 1.74

1000 sequences (390,255 amino acids) are from the G Protein Coupled Receptor family.

Relative Single-Processor Speed

In this graph the relative performance of the MIPS R12000 processor is shown in different computer systems compared to the MIPS R10000 250 MHz processor in an SGI Origin 2000 system.

Researchers using the MIPS R12000 300 MHz processors in an SGI Origin 2000 computer system can expect about a 30% performance improvement compared to the MIPS R10000 processor.

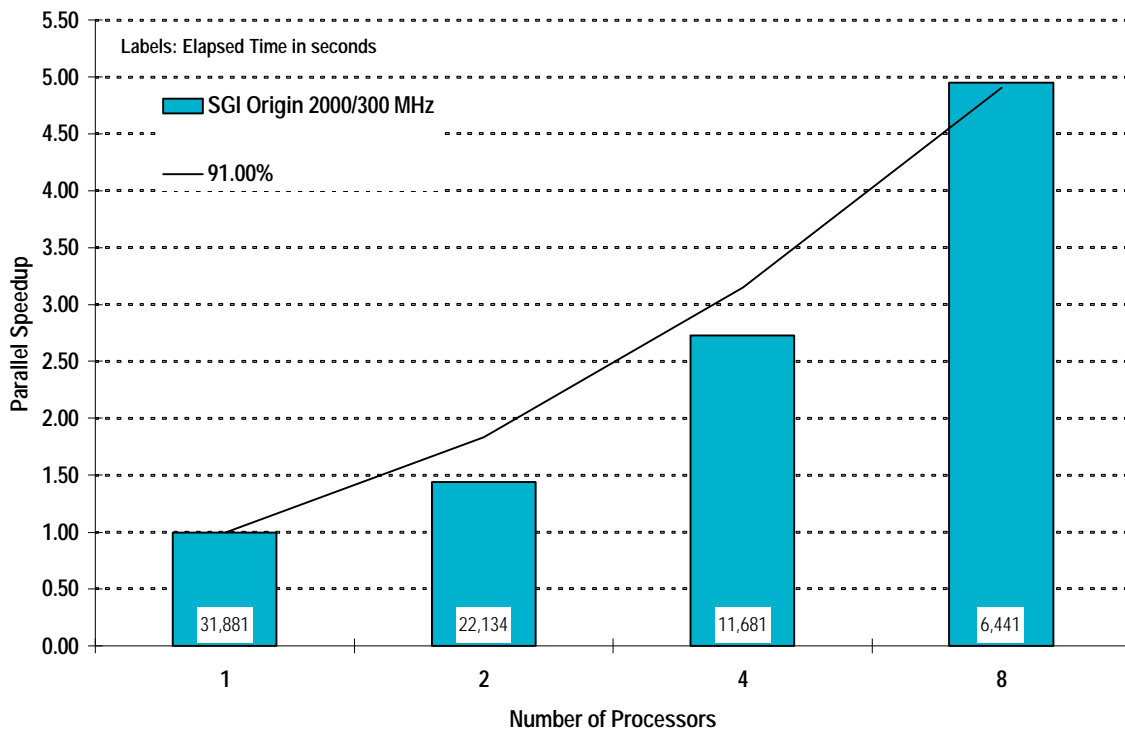


CLUSTAL W: Parallel Speedups

CLUSTAL W 1.74

There are 1000 sequences (390,255 amino acids) from the G Protein Coupled Receptor family.

The turnaround for this particular job can be improved by a factor of five when eight processors of an SGI Origin 2000 system with MIPS R12000 300 MHz processors are used.



BLAST

BLAST (Basic Local Alignment Search Tool) performs fast database searching combined with rigorous statistics for judging the significance of matches. Queries and database sequences can be in one of two formats—nucleotide or protein (amino acids).

Five BLAST programs search all combinations of query and database sequences. The BLAST algorithm is described in S.F. Altschul, W. Gish, W. Miller, E.W. Myers, and D.J. Lipman, *J. Mol. Biol.* 215, 403 (1990).

The BLAST benchmarks include queries against DNA (GenBank) databases. When additional processors are added, BLAST shows near-linear scalability. BLAST is ideally suited to MIPS R10000 and MIPS R12000 based Origin servers. This is due to the excellent integer processing capability and the high bandwidth to the memory system. SGI Origin 2000, with its large memory capacity, can memory-map even the largest sequence database files and thus eliminate I/O as the rate determining step in BLAST execution speed.

High-Throughput BLAST (HT-BLAST)

As the next few pages will show, BLAST running on SGI computers can effectively scale to large numbers of processors. BLAST was developed to handle “ad-hoc” searches, that is single-query comparisons against standard databases. It has also been used to do batch processing of large volumes of sequences. Annotation of newly sequenced EST’s using multiple databases is an example of this type of application of BLAST. The software systems which perform screening are generally built from scripts which manage the input set of sequences, invoke “out-of-the-box” BLAST executables (blastn, or blastx), and collect and sort the summary results. This design suffers from inefficiencies resulting from the repeated restarting of BLAST. This includes, in addition to generic startup overhead, the time wasted in continually remapping the subject databases. Also, scalability is limited when the query sequences are short, as in the case of EST data.

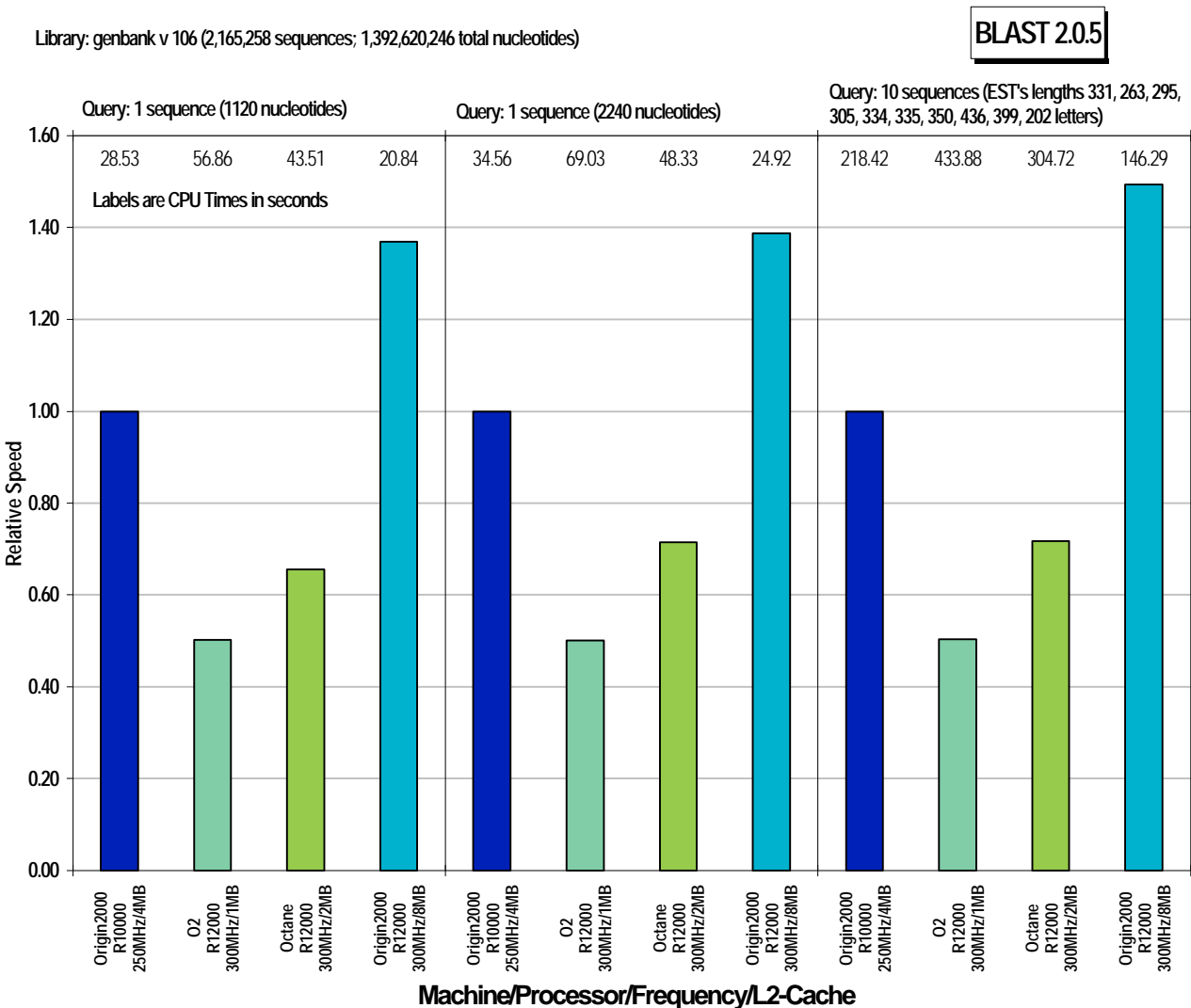
SGI has developed an alternative system that consists of a modified BLAST and a driver (High Throughput BLAST, or HT-BLAST). This version of BLAST allows multiple sequences to be compared against multiple databases via a single invocation of the code. The output of HT-BLAST is a summary of the High Scoring Pair information generated during the searches. The code saves on startup overhead through the reuse of data structures and elimination of the need to remap the databases. Also, all parallel constructs are removed from BLAST, resulting in increased single-processor speed. Parallelism has been relocated to the driver that distributes blocks of sequences to multiple processors running HT-BLAST. It uses a dynamically scheduled loop to maintain load balance. Since the independent tasks are blocks of sequences compared to multiple databases, the parallel grain-size can be much greater than it is for unmodified BLAST. Thus, scaling to large numbers of processors is accomplished even for short sequences and small databases.

HT-BLAST executables are available via <http://www.sgi.com/chembio/resources/blast/>.

HT-BLAST: Uniprocessor Performance

The relative performance of HT-BLAST on several computer systems based on MIPS R12000 processors is shown in the next graph. A comparison is made with SGI Origin 2000 equipped with MIPS R10000 processors running at 250 MHz.

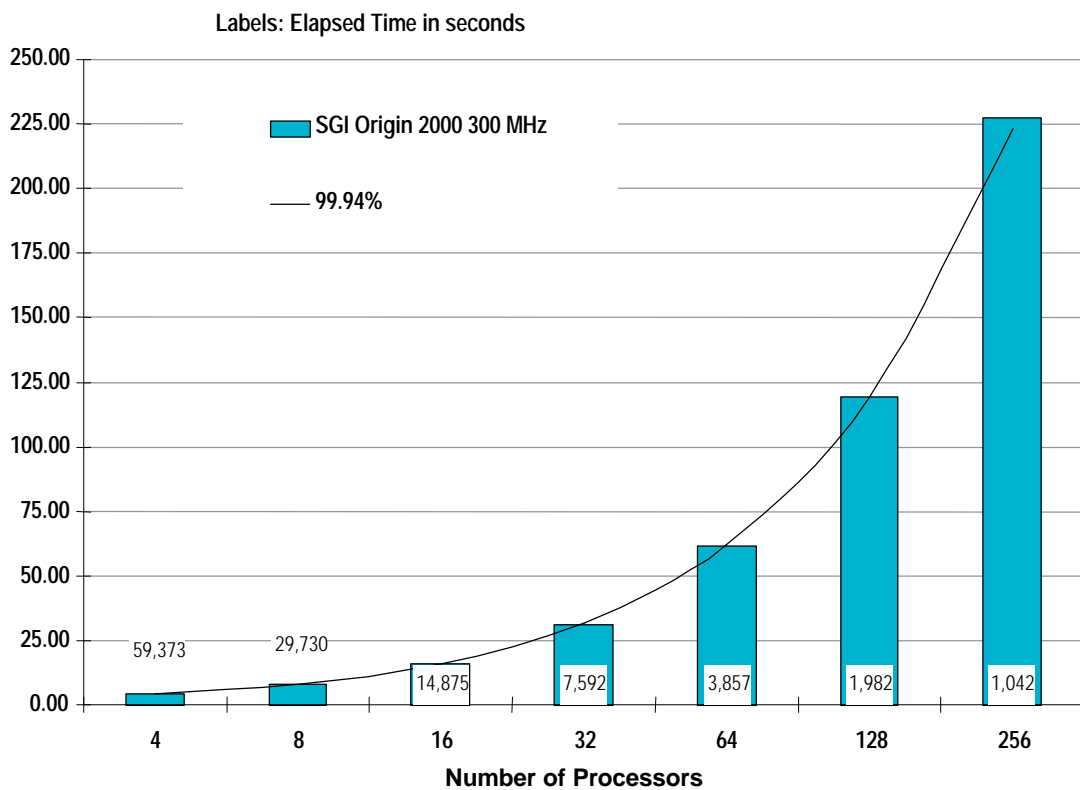
Executable: HT-BLAST performance based on v. 2.0.5 from NCBI .



HT-BLAST: Parallel Speedup

The great benefits of using HT-BLAST become very clear when searching a large number of sequences in very large databases. The benchmark consisted of 2,500 EST ranging in size from 253 to 509 bp and two databases were genpept: 423,994 sequences; 130,988,238 total aa's gb111: 4,207,758 sequences; 3,916,984,546 total nt's.

In the case shown in the graph, the turnaround time for this search is reduced to around 17 minutes when 256 processors of an SGI Origin 2000 system with MIPS R12000 300 MHz and 64GB of main memory, are harnessed for this task. On four processors this job takes about 16.5 hours. A single processor run is estimated to take almost three days.



Smith-Waterman

University of Virginia's professor Bill Pearson's Smith-Waterman algorithm was used to measure two different comparisons:

E. coli vs. *Synechocystis* and
500-nucleotide DNA query vs. the GenBank 111 database.

The comparisons were made with Pearson's PVM parallelized, distributed memory version of the Smith-Waterman algorithm and with Pearson's pthreads-parallelized, shared memory version of the application. An advantage of SGI Origin 2000 technology is the ability to run either distributed memory or shared memory parallelized applications.

The *E.coli* comparison showed a reduction in time from over 14 hours on one processor of an SGI Origin 2000 system to less than 20 minutes on 64 processors. The "moderate" parallel speedup is due to the relatively small size of the database—the *Synechocystis* database is only 1.6MB in size.

However, the 500-nucleotide search showed impressive scalability of 118X on 128 processors, reducing the time from more than 46 hours on one processor to less than 25 minutes on 128 processors. The increase in scalability is attributed to the larger size of the GenBank 111 database of 4.5GB, and to the ability of SGI Origin 2000 to process large amounts of data effectively.

Both comparisons were made on SGI Origin 2000 servers. The *E. coli* search was run on a server configured with 128 MIPS R12000 300 MHz-8MB L2 processors and 32GB of memory and running IRIX 6.5. The 500-nucleotide search was run on a server configured with 128 MIPS R12000 300 MHz-8 MB L2 cache processors and 66GB of memory running IRIX 6.5.

Smith-Waterman on SGI Origin 2000

E.coli

EXECUTABLE: FASTA 3.2 Release 5, pvcompsw/c.workgsw executables, top five hits and alignments displayed.

QUERY SEQUENCES: Synechocystis, 1,033,205 residues in 3,169 sequences.

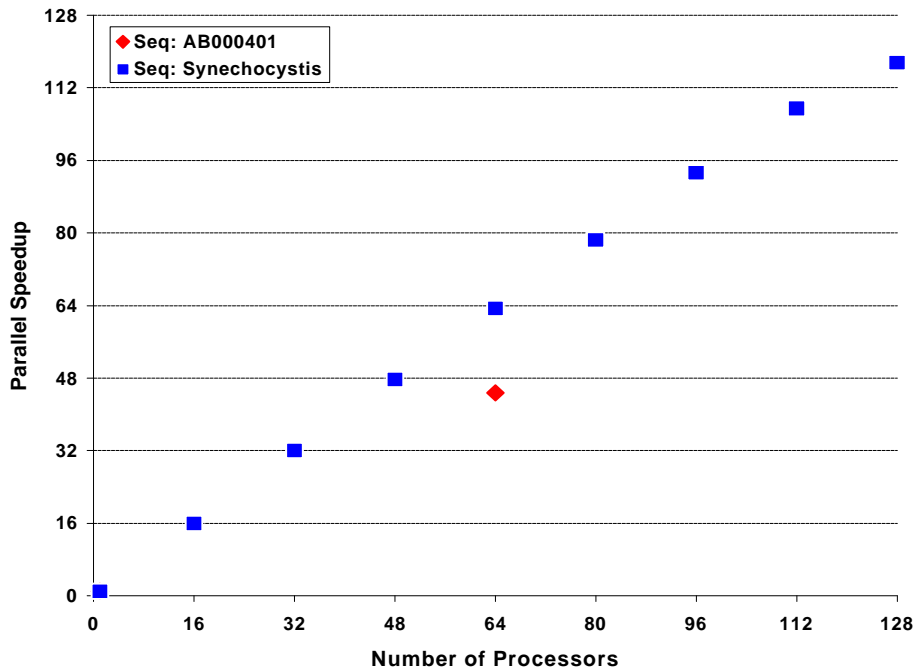
SEARCH LIBRARY: E.coli, 1358990 residues in 4289 sequences.

500-Nucleotide

EXECUTABLE: FASTA 3.2 Release 5, ssearch3_t executable, top 10 hits and alignments displayed.

QUERY SEQUENCE: Sequence AB000401 of GenBank, 500 residues.

SEARCH LIBRARY: GenBank 111.0, 3,916,981,906 residues in 4,207,758 sequences.



d2_cluster v1.21

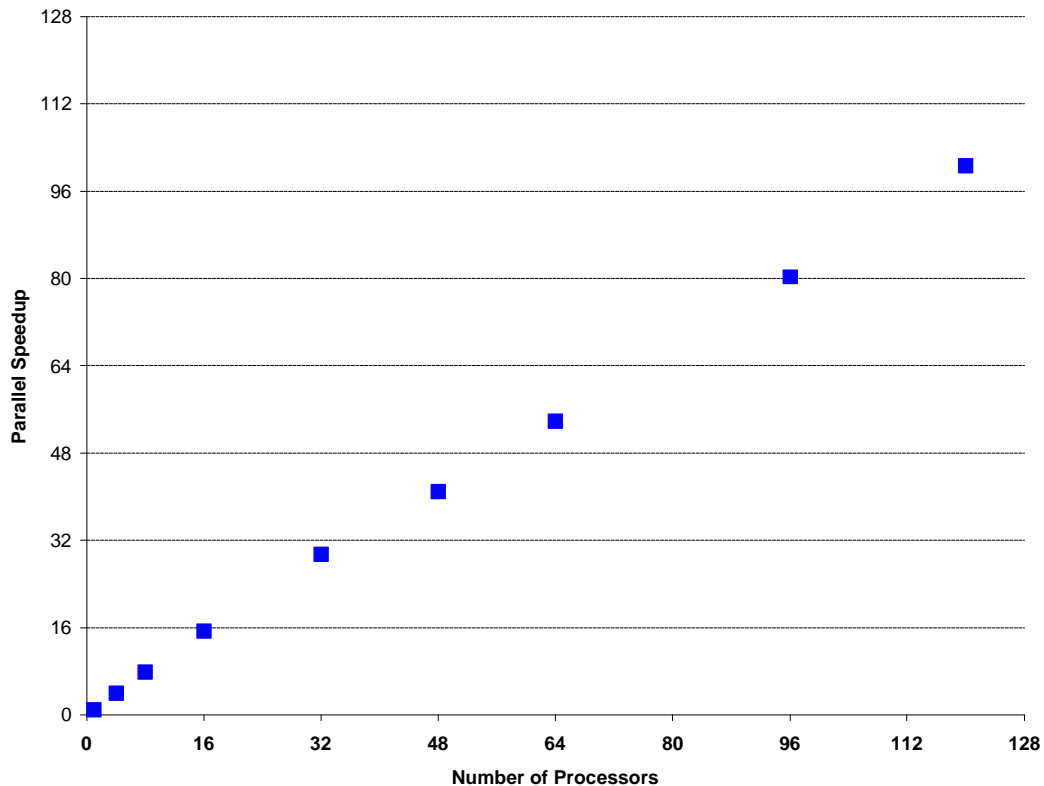
The d2_cluster, developed at the University of Houston, clusters a set of gene sequences by using the d^2 method. The program complexity scales as N^2 where N is the number of sequences. The program is used by the South African Bioinformatics Institute (SANBI).

Data set: 15,876 ESTs from lung tissue

Each word is 6 bases long, ignore any sequences less than 50 bases, compare sequences using a floating window of 150 bases, 0.04 d^2 score for clustering.

SGI Origin 2000 with 128 MIPS R12000 300 Mhz processors.

The time for this clustering was reduced from almost 3.5 hours on one processor to only two minutes on 120 processors. SGI's Bioinformatics team helped to parallelize this application efficiently.



SRS Parallel Performance

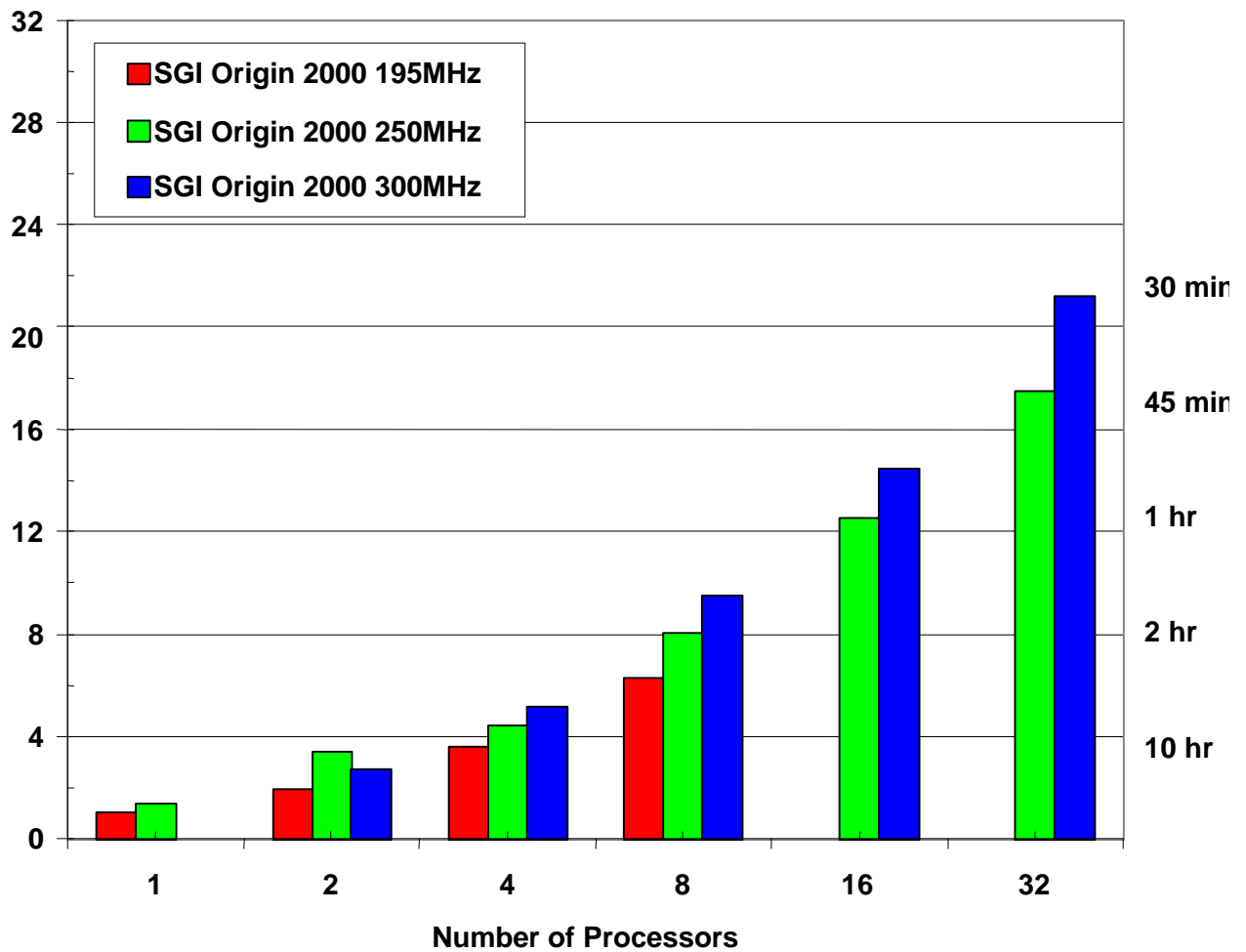
SRS (Sequence Retrieval System) is a widely used database navigation system for any biological sequence data.

Lion bioScience AG licensed SRS from the European Bioinformatics Institute (EBI) and founded a United Kingdom subsidiary directed by the author of SRS, Dr. Thure Etzold.

Timings have been measured for a complex indexing of the EMBL database—10.6GB in size.

The time for this indexing was reduced from an average of 10 hours to just over 30 minutes on an SGI Origin 2000 system with 32 MIPS R12000 300 MHz processors (24GB of memory).

Relative Speedups



Proteomics: High-Throughput SEQUEST

The "High-Throughput" techniques used in HT-BLAST can also be applied to other applications that require large-scale data processing operations, such as those in the field of proteomics. Proteomics is the study of PROTEins expressed by genOMEs, or tissues, and as a field has developed rapidly over the last several years. It involves analyzing the complex of intercellular proteins and comparing results across time, disease, or treatment state, or among tissues.

The newest and most productive techniques for identifying and analyzing intracellular proteins use mass spectrometers. When coupled in tandem and employing the latest ionization technology, these instruments can analyze thousands of peptides per day using femto- to pico-molar samples. The volume of data output makes rapid data reduction an essential part of proteomics via mass spectrometer.

Although the mass profiles generated can be used to "fingerprint" a set of proteins, they are more frequently used as the subjects of homology searches against either protein databases or translations of nucleotide databases. Finnigan Corporation, a division of Thermoquest Corporation, distributes one such program called SEQUEST, which compares experimental mass profiles with database entries in order to determine the amino acid sequence(s) and thus the protein(s) and organism(s) that best correspond to the mass spectrum being analyzed. The mass spectra data that is input to SEQUEST has been preprocessed by another program called EXTRACT_MS (University of Washington), which 'extracts' the raw data from a mass spectrometer, averages or combines multiple scans for the same daughter ion, and then generates a series of output files for homology searching by SEQUEST.

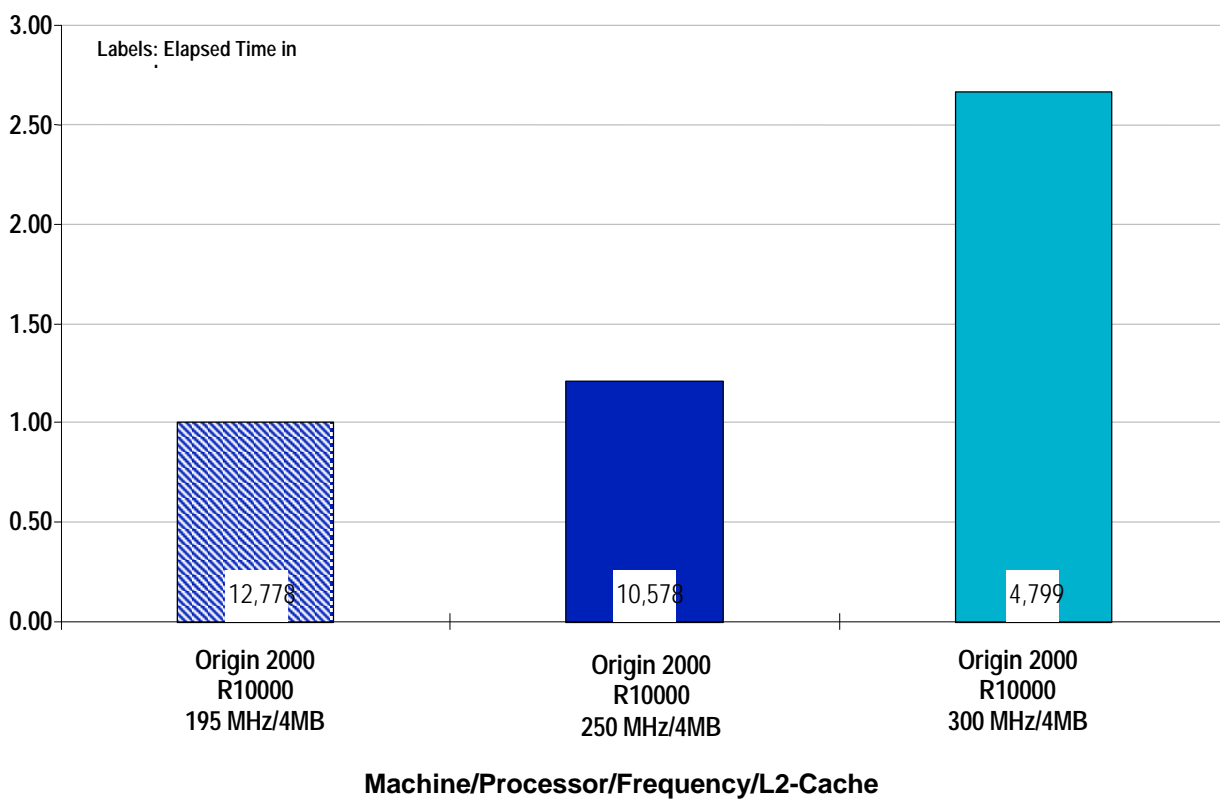
In collaboration with Thermoquest, SGI has developed an alternative system that combines and streamlines processing using EXTRACT_MS and SEQUEST called High-Throughput SEQUEST (HT-SEQUEST). HT-SEQUEST takes advantage of multiple processors by implementing a high-throughput scheme based on multiple instantiations of SEQUEST, each processing output generated as a queue by a modified version of EXTRACT_MS.

In HT-SEQUEST, no modifications were made to SEQUEST itself, but EXTRACT_MS was modified to store the information later passed to SEQUEST in a data structure rather than a file, immediately upon completion of the processing of a set of scans corresponding to a particular daughter ion. The EXTRACT_MS portion of HT-SEQUEST completes the averaging and format translation of all scans before continuing. Then, EXTRACT_MS enters the parallel regime, where each parallel thread chooses the next available element of the above data structure, outputs its contents to a temporary file, and initiates a SEQUEST execution to process that file. The scheduling of the parallel loop is dynamic, so the data structure serves as a queue whose exhaustion signals the completion of processing.

HT-SEQUEST: Relative Uniprocessor Performance

The performance on MIPS R10000 and MIPS R12000 processors is shown in the next graph. A comparison is made with an SGI Origin 2000 system equipped with MIPS R10000 processors running at 195 MHz.

HT-SEQUEST based on SEQUEST v.C2 and EXTRACT_MS v.2 (rel. 0.5).
The single data set processed contained meaningful scans for 606 daughter ions.

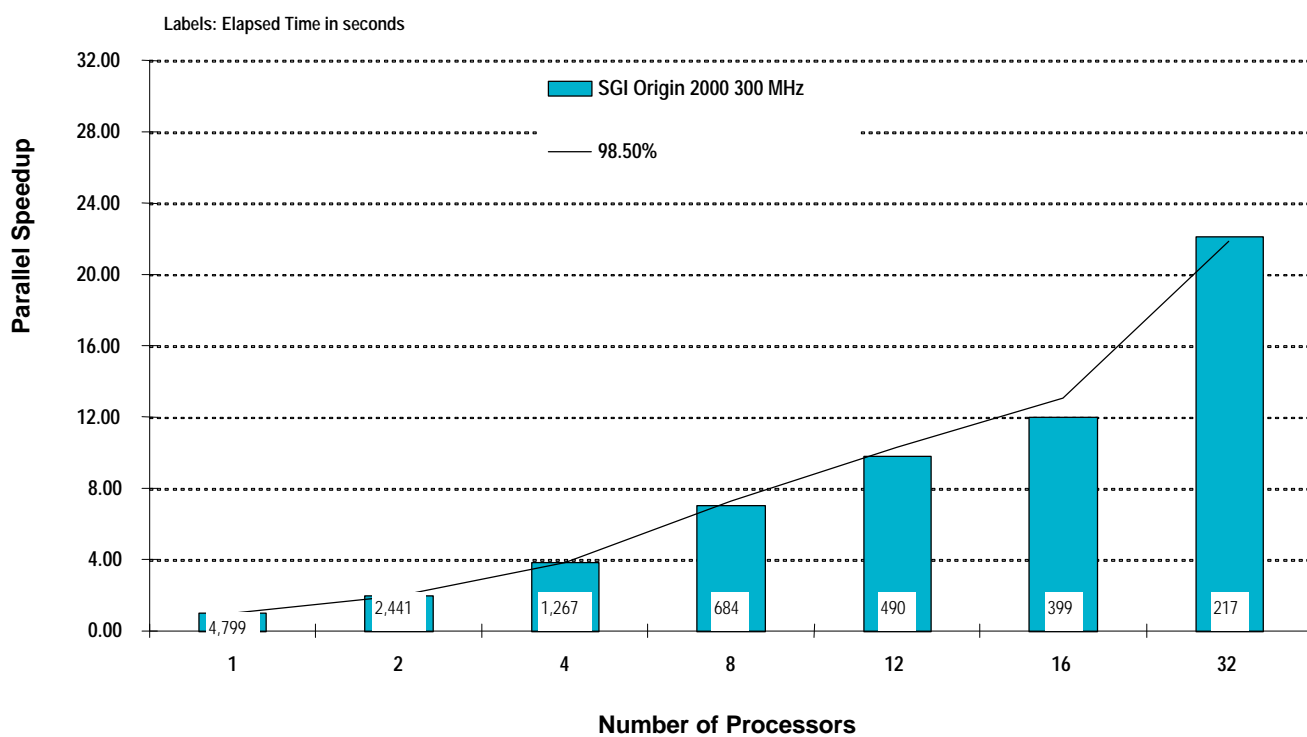


HT-SEQUEST: Parallel Speedups

HT-SEQUEST based on SEQUEST v.C2 and EXTRACT_MS v.2 (rel. 0.5).

The single data set processed contained meaningful scans for 606 daughter ions.

In the example below, HT-SEQUEST reduces the turnaround time for analyzing the data set from approximately 80 minutes on one processor to only four minutes when 32 processors of an SGI Origin 2000 system with MIPS R12000 300 MHz are used.



Appendix: SGI 1400

SGI recently introduced the SGI™ 1400 multiprocessor server, based on Intel® IA32-based microprocessors and capable of running both Microsoft® Windows NT® and the Linux® operating systems.

While the SGI Origin family of servers offers superb scalability and, where needed, 64-bit capability, the SGI 1400's uniprocessor performance on bioinformatics applications such as FASTA and BLAST is around 30% slower than that of the SGI Origin 2000 system with MIPS R12000 300 MHz processors

The following benchmarks show both uniprocessor and multiprocessor performance of FASTA, and BLAST on the following system:

SGI 1400, 4x500 MHz/2MB Pentium® III Xeon™, 1GB memory, SGI Linux 1.0 environment with Red Hat 6.0.

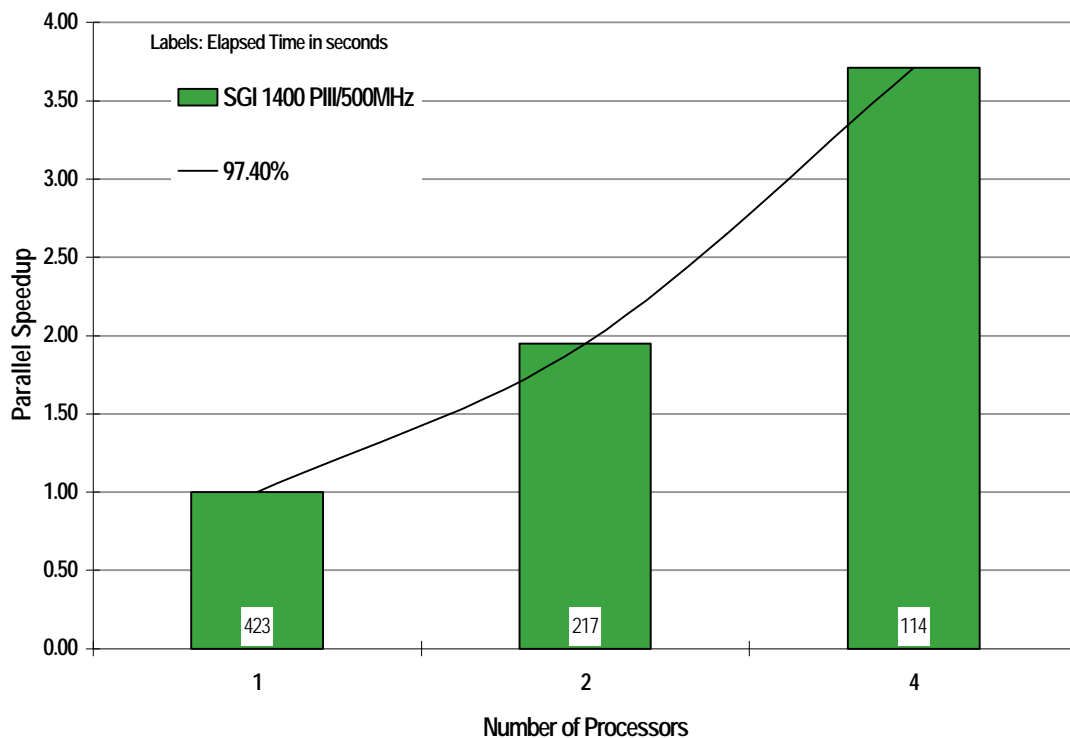
FASTA on SGI 1400L

FASTA version 3.2t01

QUERY SEQUENCE: I12R_HUMAN of Swiss-Prot34 (662 aa) with the
SEARCH LIBRARY: OWL 29.5 database (66,618,840 residues in 210,232 sequences).

Parallel Speedups

The compiler used was the PGI 3.0-4 Compiler



BLAST on SGI 1400L

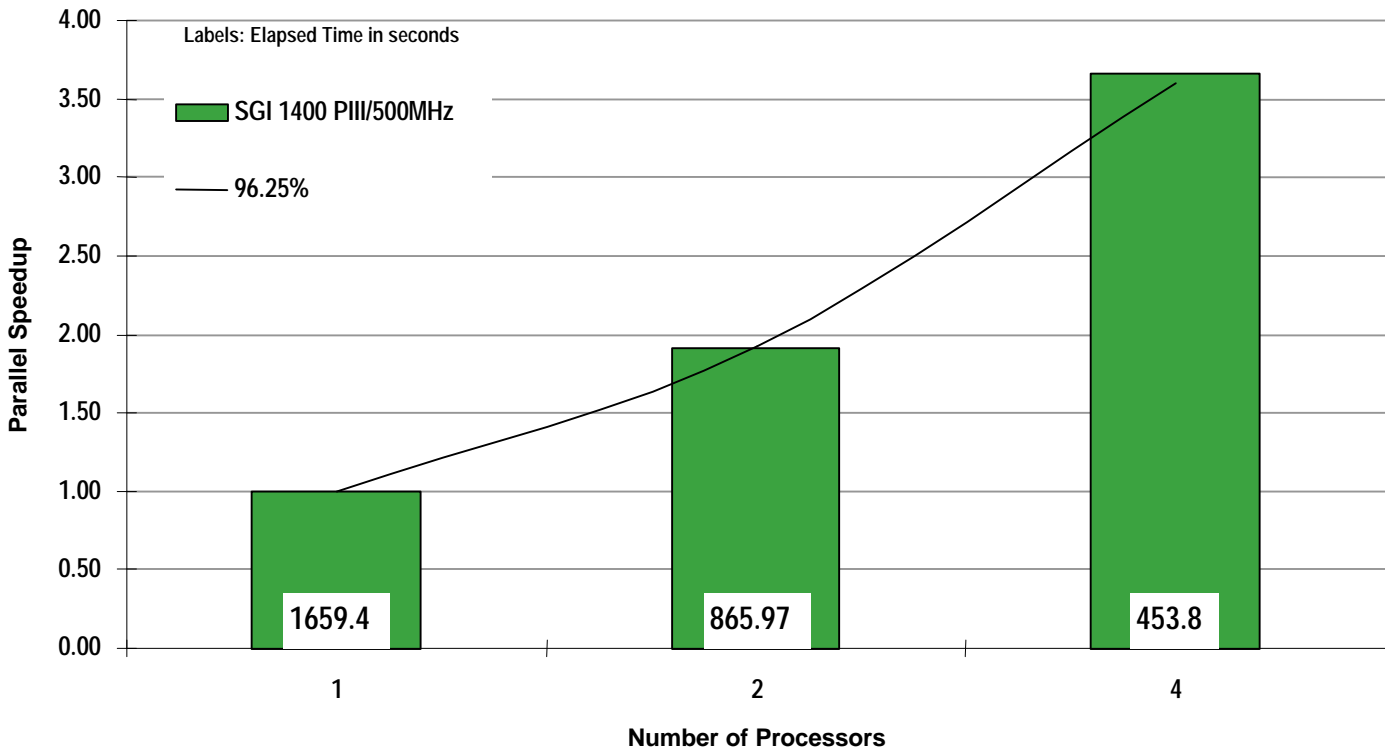
Based on v. 2.0.9 from NCBI

100 EST sequences ranging in length from 143 to 489 basepairs.

Library Accumulated daily updates of GenBank for April, May and the first half of June, 1999.
(616,612 sequences; 1,330,738,448 total nt's).

Parallel Speedups

For this and the next example, the BLAST executable was created using the gcc compiler.

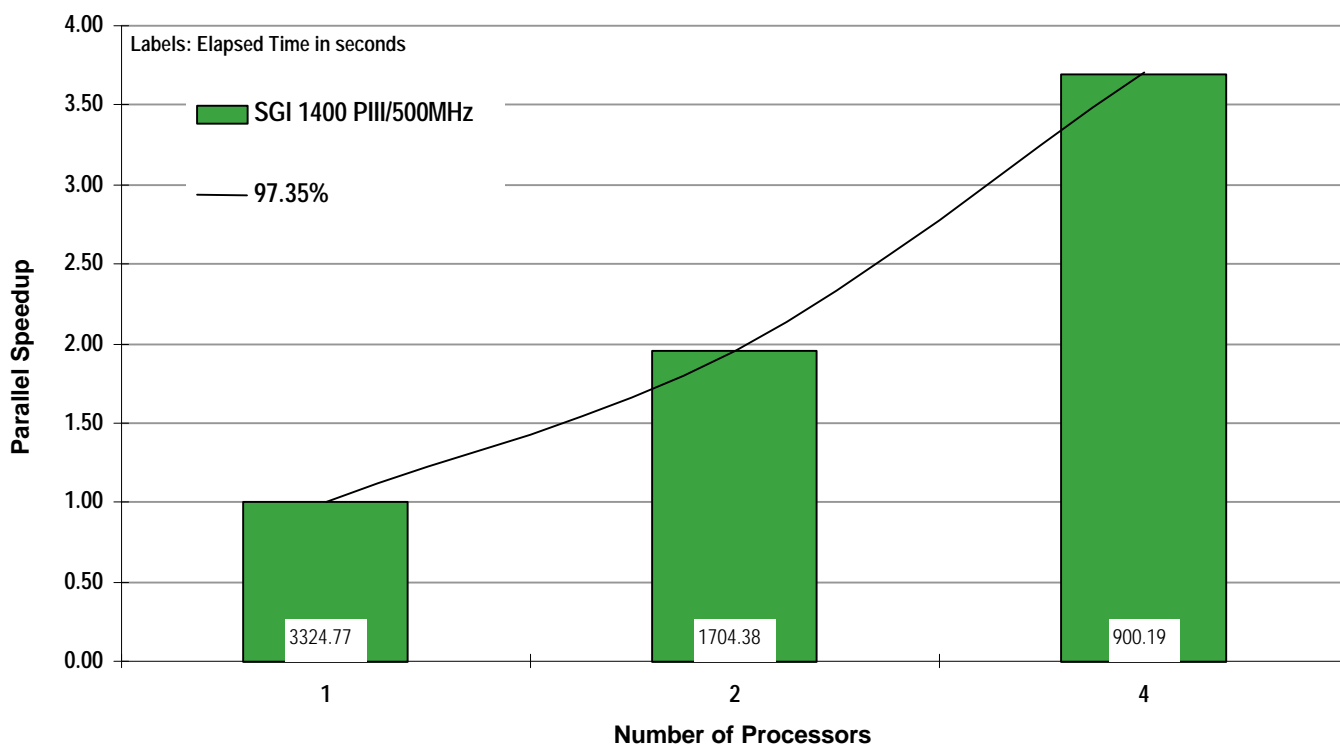


BLAST on SGI 1400L

Based on v. 2.0.9 from NCBI.

100 EST sequences ranging in length from 143 to 489 basepairs.

Library: accumulated daily updates of GenBank versus dailies and nr (322,347 sequences; 97,848,900 aa's).



Acknowledgments

BLAST

Dr. Warren Gish
Washington University
<http://blast.wustl.edu/blast/README.html>

Dr. Thomas Madden,
National Center for Biotechnology Information
<http://www.ncbi.nlm.nih.gov/BLAST>

CLUSTAL W

Dr. Toby Gibson and Dr. Julie Thompson
European Molecular Biology Laboratory

Dr. Des Higgins
University of County Cork

d2-CLUSTER

Dr. Winston Hide
South African National Bioinformatics Institute, University of Western Cape
<http://www.sanbi.ac.za>

FASTA

Dr. William R. Pearson
University of Virginia
<http://www.med.virginia.edu/~wrp/pearson.html>

SEQUEST and EXTRACT_MS

Dr. Jim Shofstahl
Finnigan Corporation
<http://www.finnigan.com/>

Dr. Jimmy Eng and Dr. John Yates
University of Washington
Dept. of Molecular Biotechnology
<http://thompson.mbt.washington.edu/sequest/>

SRS

Lion bioScience
<http://www.lion-ag.de>

SEQUEST

Thermoquest, Inc.
408 965-6515
<http://www.thermoquest.com>

© 1999 Silicon Graphics, Inc. All rights reserved. Specifications subject to change without notice. Silicon Graphics, O2 and Octane are registered trademarks, and SGI, Origin and the SGI logo are trademarks of Silicon Graphics, Inc. MIPS and R10000 are registered trademarks, and R12000 is a trademark of MIPS Technologies, Inc. R10000 and R12000 are trademarks or registered trademarks used under license by Silicon Graphics, Inc. Intel and Pentium are registered trademarks, and Xeon is a trademark, of Intel Corporation. Linux is a registered trademark of Linus Torvalds. Microsoft, Windows, and Windows NT are registered trademarks of Microsoft Corporation.

