

# The Tera Computer System

Robert Alverson

David Callahan  
Allan Porterfield

Daniel Gummings  
Burton Smith

Brian Koblenz

Tera Computer Company  
Seattle, Washington USA

## 1 Introduction

The Tera architecture was designed with several major goals in mind. First, it needed to be suitable for very high speed implementations, *i.e.*, admit to fast clock periods and be scalable to many processors. The goal will be achieved; a maximum configuration of 4096 processors, 512 memory units, 256 I/O cache units, and 4096 I/O processors, and a clock period less than 3 nanoseconds. The abstract architecture is scalable essentially without limit (although a particular implementation is not). The only requirement is that the number of streams increase more rapidly than the number of physical processors. Although this means that the number of streams is sublinear in the number of instruction streams, it still increases linearly with the number of processors. The price/performance ratio of Tera is unmatched, and puts Tera's high performance within economic reach.

Second, it was important that the architecture be applicable to a wide spectrum of problems. Programs that do not vectorize well, perhaps because of a preponderance of scalar operations or too many conditional branches, will execute efficiently as long as there is sufficient parallelism to keep the processor busy. Virtually any parallelism available in the workload can be turned into speedup at the program level parallelism within program multiuser time- and space-sharing.

Third, it was important that the architecture be applicable to a wide spectrum of problems. Programs that do not vectorize well, perhaps because of a preponderance of scalar operations or too many conditional branches, will execute efficiently as long as there is sufficient parallelism to keep the processor busy. Virtually any parallelism available in the workload can be turned into speedup at the program level parallelism within program multiuser time- and space-sharing.

## 2 Interconnection Network

The interconnection network is a three-dimensional mesh of computer packet-switching nodes, each connected to some of its neighbors. Each link contains a basic block containing source and destination addresses, operation, and 64 data bits in bits simultaneously on every clock tick. Some nodes are also linked to *resources*, *i.e.*, processor memory units, I/O processors, and I/O cache units instead of locating the processors on one side of work and memories on the other (in what Robert has called a "dancehall" configuration [5]), they are distributed more-or-less uniformly throughout the network. This permits data to be placed in units near the appropriate processor when the time is available and otherwise generally maximizes the distance between possibly interfering resources.

\*This research was supported by the United States Defense Advanced Research Projects Agency under Contract MD972-89-C-0002. The views and conclusions contained in this document are those of Tera Computer Company and should not be interpreted as representing the official policies, either expressed or implied, of DARPA or the U.S. Government.





ple: each instruction contains a three bit branch prediction field, and allow several streams that explicitly specifies how many instructions are fetched simultaneously. The stream will issue before encountering an instruction that would exceed the system limit on the number depends on the current one. Since seven is the maximum protection domain can reserve. possible lookahead value, at most eight instructions can be concurrently executing a CREATE operation to create twenty-four operations can be concurrently executing increments *scur*, generates the stream from each stream. A stream is ready to issue a new stream using one of its own target instruction when all instructions with lookahead level *l* have target *TO* from its own *TO* registers referring to the new instruction have completed. These registers in the new stream from if each stream maintains a lookahead of seven target registers. The newly created streams are needed to hide 72 ticks of latency quickly begin executing useful work in context.

Lookahead across one or more branch operations is created as long as significant store handled by specifying the minimum of all instructions necessary. The QUIT operation terminates the variant branch operations JUMP\_OFTEN that executes it, and decrements both JUMP\_SELDOM, for high- and low-probability branches the QUIT\_PRESERVE operation only decrements respectively, facilitate optimization by security by giving up a stream without surrendering lookahead along the less likely path. There are realisation.

SKIP\_OFTEN and SKIP\_SELDOM operations. The overall Each protection domain has a *retry limit* that determines approach is philosophically similar to an exposed window. If a memory reference can find lookahead except that the quanta are instructions as a protection's full/empty bit (see section ticks. it will trap. If a synchronization is not successful long time, then possibly a heavier weight mechanism that avoids busy waiting should be used to wait for synchronization. The retry limit should be the product of trap processing overhead, which defines the protection domain.

#### 4.4 Protection Domains

Each processor supports as many as 16 active protection domains that define the program memory, data memory, and number of streams allocated to the processor and invoke the heavier weight mechanisms using that processor. Each execution stream is assigned to a protection domain, but which domain (or which processor, for that matter) is not known to the user program. In this sense, a protection domain is a virtual processor and may be moved from one physical processor to another.

#### 4.5 Privilege Levels

The privilege levels apply to each stream independently. There are four levels of privilege: user, supervisor, kernel, and IPL. IPL level operates in absolute protection domain has two pairs of map bases and must be the highest privilege level. User registers that describe the region of each map and kernel levels use the program and data to it. The upper 2048 data segments and 1024 address translation, and represent increments pages are not relocated by the map bases, of privilege level. The data map entries define the by the operating system. Any active protection domain to read and write each segment can use all of either or both maps. The map program map entries define the *exact* level needed to read, write, or execute the segment or page, which is stored as part of the privileged stream the segment or page was read, written, and is not available to a user-level stream. appropriate; and the distribution (for the hardware) operations are provided to all streams.

The number of streams available to a protection domain is regulated by three quantities *slim*, *scur*, and *level*. The *LEVEL\_ENTER* operation sets the current privilege level associated with each protection domain. The *level* is incremented if the current number of streams executing in the protection domain is greater than the *level* recorded by *scur*; it is incremented when a stream reaches a point that can accept a call from a higher level. A trap occurs if the current *level* is greater than the *level* recorded by *scur*. The *LEVEL\_RETURN* operation is used to decrease the *level* if it has exceeded *sres*, the number of streams reserved to the original privilege level. A protection domain. The operations for reserving streams are the same as the current privilege level.



does not combine fetch-and-add operation  
memory location.

## 6 Arithmetic

The numeric data types directly supported by the architecture include:

- 64 bit two complement integers
- 64 bit unsigned integers
- 64 bit floating point numbers
- 64 bit complex numbers

Operations on these types include addition, subtraction, multiplication, conversion, and comparison. Reciprocal reduction and dataflow is provided for using Newton's method.

Other types are supported indirectly, including:

- 8, 16, and 32 bit two complement integers
- 8, 16, and 32 bit unsigned integers
- arbitrary length unsigned integers
- 32 bit floating point numbers
- 128 bit "doubled precision" numbers

The shorter integers are sign- or zero-extended to the appropriate length as they are loaded from memory, and truncated to the appropriate length as they are stored. The fundamental support for arbitrary length arithmetic is provided by the operations `INTEGER_ADD_MUL`, `UPPER_ADD_MUL`, and `CARRY_ADD_TEST` that together implement  $64 \times n$  bit unsigned multiply-add in approximately  $2 \frac{1}{2}$  instructions.

The 32 bit floating point numbers are simply the real parts of the 64 bit complex type with imaginary parts set to zero. The 128 bit "doubled precision" type was pointed out to us by Kahan [1, 7, 4]; it represents a real number  $R$  as the unevaluated sum of two 64 bit floating point numbers  $r$  and  $\rho$ , where  $\rho$  is insignificant with respect to  $r$  and as near as possible to  $R-r$ . Support for this type is provided by `FLOAT_ADD_LOWER` which (with `FLOAT_ADD`) implements "doubled precision" addition in six instructions, and by `FLOAT_MUL_ADD` which rounds only once and is used to implement "doubled precision" multiplication in five instructions.

## References

- [1] T. J. Dekker. A floating-point technique for extending the available precision. *Numerische Mathematik*, 18:224-242, 1971.
- [2] M. Flynn. Some computer organizations and their effects. *IEEE Transactions on Computers*, C-21(9):948-960, September 1972.
- [3] A. Gottlieb, R. Grishman, C. P. Kruskal, McAuliffe, L. Rudolph, and M. Snir. The NYU ultracomputer - designing an MIMD shared memory parallel computer. *IEEE Transactions on Computers*, C-32(2):175-189, 1984.
- [4] W. Kahan. Doubled-precision IEEE standard floating-point arithmetic. Unpublished manuscript, February 1987.
- [5] R. M. Keller. Rediflow: A proposed architecture combining reduction and dataflow. In *PAW83: Proceedings of the 1983 Parallel Architecture Workshop*, University of Colorado, Boulder, 1983.
- [6] J. T. Kuehn and B. J. Smith. The Horizon supercomputer system: Architecture and software. In *Proceedings of Supercomputing '88*, Orlando, Florida, November 1988.
- [7] S. Linnainmaa. Software for doubled-precision floating-point computations. *ACM Transactions on Mathematical Software*, 7:272-283, 1981.
- [8] A. Norton and E. Melton. A class of boolean transformations for conflict-free pointer stride access. In *Proceedings of the 1987 International Conference on Parallel Processing*, pages 247-254, August 1987.
- [9] J. Erankka, P. Pittelli and David Smitley. Analytical network for a shared memory architecture. In *Proceedings of Supercomputing '88*, Orlando, Florida, November 1988.
- [10] M. R. Thistle and B. J. Smith. A process architecture for Horizon. In *Proceedings of Supercomputing '88*, pages 35-41, Orlando, Florida, November 1988.