

BlackWidow Hardware System Overview

Brick Stephenson

Cray, Inc.

(bas@cray.com)

ABSTRACT: *This paper describes the hardware features of the BlackWidow scalable vector multiprocessor. Advances in ASIC technology and design techniques have enabled the BlackWidow project to improve processor and network performance, and dramatically reduce the cost per MFLOP. We provide an overview of the BlackWidow system packaging and network topology, processor improvements, maintenance system, and reliability features.*

KEYWORDS: Vector Processors, BlackWidow, Fat Tree, Clos Network

1. Introduction

The BlackWidow is the project code name for the next generation scalable vector multiprocessor from Cray Inc. The system is currently in development with several components in initial checkout. Processor blades are targeted for power up in summer 2006 with production shipment planned for 2007.

As with the Cray X1 [1] system, the BlackWidow (BW) system uses a Cray proprietary custom vector processor design. The system is a distributed shared memory multiprocessor, with each *node* consisting of four BW processors. The system is designed to run demanding scientific applications requiring significant memory and network bandwidth. The BW system uses a highly parallel memory system that is built from standard DDR2 memory devices and is driven by a custom controller chip (Weaver). The Weaver memory controller contains the L3 cache, memory directory for cache coherence, and DDR2 memory manager.

In a large-scale multiprocessor the interconnection network is the principal determinant of remote memory latency, and point-to-point bandwidth between two processors. Each BW processor can have several thousands of outstanding memory requests pipelined in the network, so an efficient network is critical to application scalability. The BW network implements a high-radix folded-Clos (fat tree) topology built from a radix-64 router chip. This network scales to thousands of processors while delivering high network bandwidth and low latency between processors. For example, a 4K processor system has a network diameter of only five hops. The network can be configured to achieve a balance between network bandwidth and cost per customer

requirements including adding specially designed router cabinets for large systems.

The BlackWidow system will be loosely coupled to a Cray XT [2] system (like XT3 or its follow-on), which will serve as the IO and service access into the BlackWidow system. Although the network topology of the Cray XT is a 3-D torus and the BlackWidow is a fat tree [3], the systems can still interoperate through a *protocol bridge* chip (StarGate) housed on a Bridge blade. The protocol Bridge blade plugs directly into the BlackWidow mechanical packaging and connects to the XT system using a network cable. This combination of BlackWidow and Cray XT systems provides customers with a heterogeneous computing environment, where the two systems share login nodes and a common file system.

A flexible and high-density packaging solution is critical to keeping cost down. The BW system provides two modular packaging options: air cooled for small system configurations, and liquid cooled for large systems that require more efficient heat removal from the computer room.

2. Processor Improvements

An aggressive custom vector and scalar pipeline is at the heart of the BW processor. Through extensive use of custom CMOS design practices, we achieve significant improvement in scalar performance, including lower cache hit latency. In addition, several architectural improvements were added to support additional operations in the NV-2 instruction set, reliability and availability features, and interconnection network.

Instruction Set

The BlackWidow processor implements the Cray NV-2 vector ISA (instruction set architecture), which is a

variant of the Cray NV-1 vector ISA used by its predecessor, the Cray X1. Some of the ISA enhancements include:

- Inclusive-OR version of the bit matrix multiply (Bit matrix compare)
- Vector AMO {Add,And,Or,Xor}
- Vector AMO w/Fetch {Add,And,Or,Xor}
- Versions of Gather and Scatter with Sword (32-bit) indices
- Immediate logicals, integer multiply, and conditional move

Scalar Improvements

A major goal of the BlackWidow project was to improve the overall scalar performance over the Cray X1 system. This was accomplished by: faster clock speeds, architecture improvements, more custom implementation within the scalar core, and reduced latencies. Some of these improvements include:

- 4-way instruction dispatch
- Active instruction window enlarged
- Speculative scalar loads
- Number of outstanding branches increased
- D-Cache hit-time reduced
- D-Cache protected from vector traffic
- Level-2 cache hit-time reduced
- Local Memory latency reduced

Vector Improvements

The BW processor is a more aggressive vector processor than prior Cray systems. Vector improvements include:

- Maximum vector length increased
- Number of vector masks increased
- Removed the mod-32 register usage restriction
- Full speed bit-matrix multiply

The vector core (or pipes) functional unit groupings (FUGs) are carried over from the Cray X1 processor. FUG1 contains the following integer and floating-point functional units:

- Integer Add/Subtract
- Floating-point Add/Subtract
- Logical {AND, OR, XOR, ANDNOT, MASK}
- Integer Compare
- Floating-point Compare

FUG2 contains the following integer and floating-point functional units:

- Integer Multiply
- Floating-point Multiply
- Shift

FUG3 contains the following integer, floating-point functional, and special units:

- Floating-point Divide
- Floating-point Square Root
- Leading Zero
- Pop Count
- Copy Sign
- Absolute Value
- Logical
- Integer Convert
- Floating-point Converts
- Vector Merge

3. Node Organization

BlackWidow node is organized as a 4-way SMP, where the four processors within the node have uniform access to the shared address space. Each processor has 16 full-duplex channels connected to 16 Weaver chips (Figure 1). The Weaver chips serve three functions: provide a common interface between the four processors within each node, access to a large L3 cache (8 MB total) and memory directory, and the DDR2 memory interface.

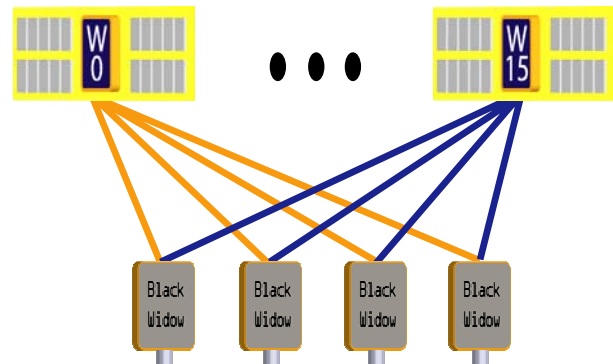


Figure 1. Node Architecture

Each BlackWidow processor has four network injection ports that connect to an independent *slice* of the network (described in Section 4). The BlackWidow compute blade is packaged with eight BW processors and 32 memory daughter cards (32 Weaver chips) on each physical blade. However, the two nodes on each blade are completely independent of each other.

4. Network

The Cray BlackWidow system uses a radix-64 folded-Clos [4] (fat tree) network. This design, using a high-radix router [5, 6] with many narrow channels, makes more efficient use of modern high-bandwidth ASIC technology than previous low-radix networks built from routers with a few wide channels. The high-radix fat tree topology [6] provides a number of benefits:

- Low hop count and latency

- Low cost for a given global bandwidth, or alternatively, high global bandwidth for a given cost
- High degree of flexibility in configuring a bandwidth profile; can *taper* the network at any level
- Many alternative routes for fault tolerance
- Ease of routing and configuration

The building block for the network is a radix-64 router (YARC¹). The YARC chip has 64 full-duplex ports with each port supporting multiple lanes (bits).

The system packaging hierarchy spans *blades*, *router modules*, *chassis*, and *cabinets*. The scalable building block for the network is a chassis. Each chassis holds eight blades with eight *endpoints* (an endpoint is either a BW processor, or a StarGate protocol bridge chip) per blade, and four rank-1 router modules with two YARCs per module (Figure 2). These components are divided into two 32-endpoint rank 1 subtrees.

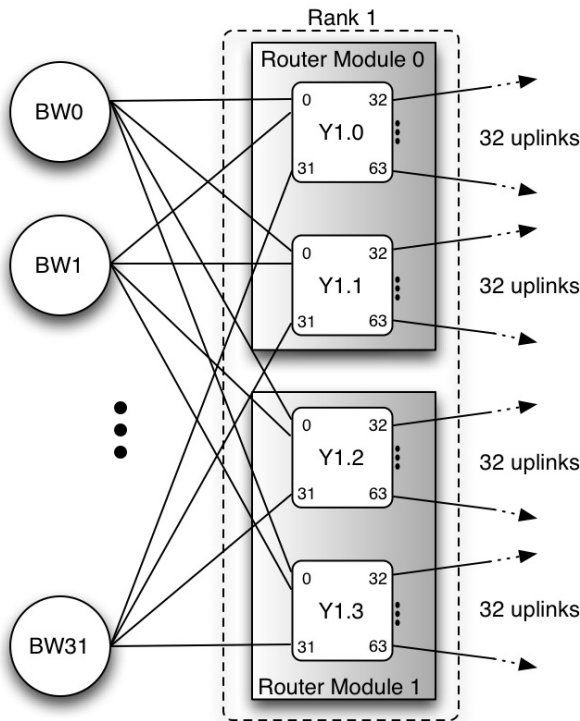


Figure 2. 32-processor Rank 1 network (figure used with permission from [6])

Within a subtree, each processor has four injection ports into the network, with each connecting to a different slice of the network. Each slice is a completely separate network with its own set of YARC router chips. The remainder of this section focuses on a single slice of the network.

¹ YARC stands for “yet another router chip”

Each YARC within the subtree provides 32 *downlinks* that are routed to the processor ports via the chassis mid plane, and 32 *uplinks* (or *side-links*) that are routed to eight cable connectors. Cables are used to connect the subtrees to expand the network.

Each *cabinet* holds two chassis (128 endpoints) organized as four 32-endpoint R1 sub-trees. Machines with up to 256 endpoints, eight R1 subtrees, can be connected by directly cabling the R1 subtrees to one another using *side-links* to create a rank 1.5 (R1.5) network.

To scale beyond 256 endpoints, the uplink cables from each rank 1 subtree are connected to rank 2/3 routers. A rank 2/3 router module packages four YARC router chips. The four R2/3 routers are logically distinct, being co-located on the module solely for packaging convenience. Up to sixteen R2/3 router modules are packaged into a stand-alone router cabinet.

5. Hardware Supervisory System (HSS)

The BlackWidow maintenance system consists of three main controller types (see Figure 3). Each type of controller incorporates a single board computer (SBC) that contains a microprocessor, SDRAM and flash memory. The HSS system monitors the power and cooling environment at multiple levels and can power down individual components in case of a failure.

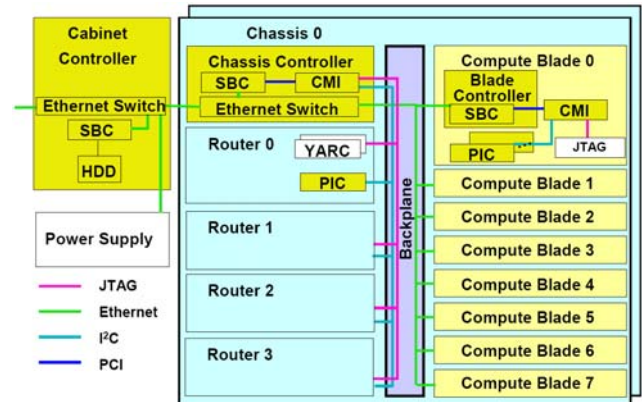


Figure 3. BlackWidow Maintenance System

Cabinet Controller

- Monitors and controls incoming power
- Monitors and controls cabinet cooling
- Contains a Ethernet switch for communication distribution to other controllers
- IDE hard disk drive

Chassis Controller

- Ethernet Switch
- Configuration Management Interface (CMI)

Blade Controller

- Configuration Management Interface (CMI)
- Sensor controllers to monitor power, cooling, etc.
- JTAG channels to communicate with logic chips

At the blade level the HSS system also monitors memory-mapped registers (MMRs) inside of the logic chips. The HSS system uses these MMRs to control maintenance, configuration, and degrade features. Detailed status of errors and other information is also available in the MMRs. The MMRs can also be accessed by the operating system.

The memory controller inside of the Weaver chip contains several features to achieve maximum performance and reliability of the DDR2 DRAMs. A very comprehensive training sequence can be controlled by the HSS system to compensate for timing differences in the memory components. The memory controller also contains several resiliency features: memory scrubbing of soft errors in DRAM, auto spare-bit insertion to recover from persistent single-bit errors, and auto and distributed refresh options.

6. Reliability and Resiliency

In order for systems to be truly scalable they must have very strong reliability and resilience features. The BlackWidow system has several reliability and resiliency features that enable the system to ride through some of the most common hardware failures. Components and chassis are designed for ease of serviceability with redundancy built in wherever possible. Some of these features are described below.

Fault detection, diagnoses, and recovery

A comprehensive error detection and logging system has been designed to protect and monitor all components throughout the system. The Hardware Supervisory System (HSS) is used to manage error reporting and error recovery.

All memories throughout the system have error detection. Most memories in the system are protected with single-bit error correction.

Reference timeouts are used to detect lost references in the system. Self-cleansing of the data paths features are designed to prevent cascading errors into the network.

Hardware firewalls for fault containment

The processor is designed with firewalls to protect other processors from being corrupted by a processor that has failed. Hierarchical boundaries are established between kernel groups using these firewall features; this protects the rest of the system even if a kernel is corrupted.

Graceful network degradation

All communication channels have self-diagnoses features to help prevent failures of the channels. These channels are CRC protected and when an error is detected a retransmission of the failed data transfer is initiated. The channel keeps track of retries and if the channel starts failing too often the channel can be auto degraded to avoid using the failing data paths (bit lanes). The channels have multiple bit lanes and can continue to operate with some number of bit lanes disabled. More severe errors will also be detected and the channel can take itself completely down if necessary. All channels can be restarted by the HSS system without interrupting the operating system.

Hot swappable boards

The hardware supports hot swapping of all major components in the system although software support of this may not be available initially. Compute blades, router modules, and bridge blades all can be individually powered down by the HSS system and removed from the system without affecting the operation of other components. Power supplies can be hot swapped without affecting the operation of the system and requires no software intervention.

Re-configurable routing tables

In the network, routing tables can be re-configured to avoid defective paths. The HSS system detects a failing link and notifies the operating system. The routing tables can then be modified to avoid routing any more packets through the bad link.

Redundant Components

Where ever possible redundant components are used to increase the resiliency of the system. All major power components are N+1 and are designed to disable the output in case of failure to minimize the effect on the power system. Some cooling components are redundant. The HSS system will detect a failing power supply or fan

and report to the maintenance system. Service then can be scheduled at a later date.

7. Packaging

A BlackWidow compute blade, shown in Figure 4, contains two BlackWidow compute nodes. The two nodes implemented on a BlackWidow blade are completely independent of one another, each belonging to a different subtree within the chassis. Each blade contains 8 BlackWidow chips (CPU and network access links), 32 memory daughter cards, voltage regulator modules (VRM), DC-to-DC converters, and a midplane connector.

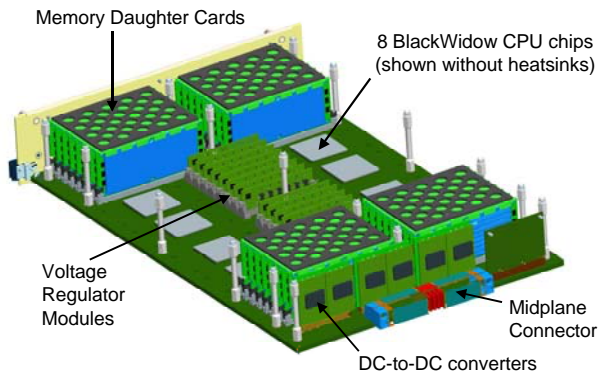


Figure 4. BlackWidow compute blade

Each memory daughter card on the compute blade contains 20 DDR2 SDRAM memory chips (1 Gbit initially) and one Weaver control chip. The 20 memory chips include 16 data and 4 ECC chips to implement full single error correction double error detection (SECDED). Larger memory will be available by *stacking* 1 Gbit memory chips initially, and 2 Gbit chips later as these components become available.

Bridge blades each contain eight Stargate chips. Each Stargate chip provides one interface between the BlackWidow YARC network [6] and the Cray XT [2] network. From the YARC network perspective, each Stargate chip functions as a YARC endpoint. Thus both BlackWidow blade types contain eight endpoints and can be plugged into a common chassis.

The BlackWidow chassis houses up to eight compute and/or bridge blades, which are inserted vertically. Rank-1 router modules are connected orthogonally to the compute blades through a mid-plane. There are 2 or 4 router modules per chassis depending on customers network bandwidth requirements. Figure 5 shows a chassis with four router modules installed.

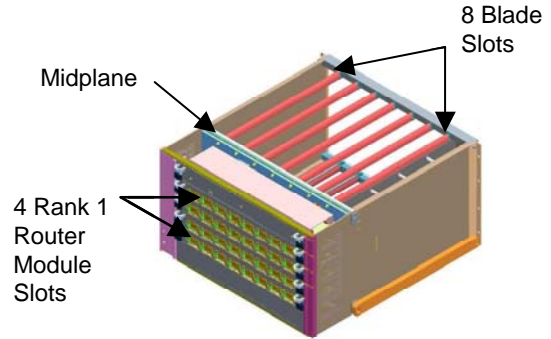


Figure 5. BlackWidow Chassis

There are two chassis per compute cabinet, for a total of 16 blades per cabinet or 128 endpoints. Figure 6 shows the base air-cooled cabinet.

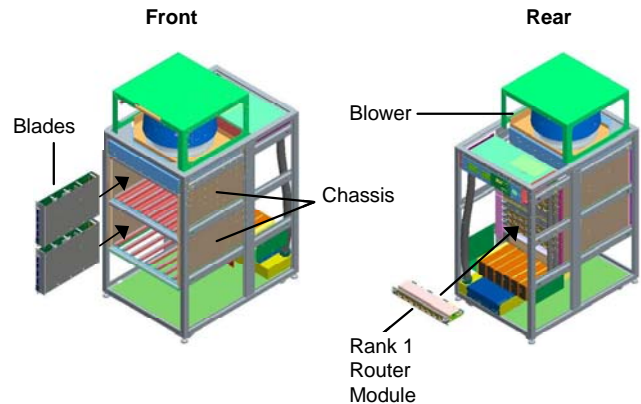


Figure 6. Air-cooled compute cabinet

Using cooling technologies similar to other Cray Inc systems, all the components in the BlackWidow cabinet are air-cooled. The blower located at the top of the cabinet draws air into the bottom of the cabinet. Air is blown across the blades and exhausts out the rear of the cabinet.

For customers that require more efficient heat removal, a liquid cooling kit is available for the compute cabinet. The kit includes water coils that are installed at the base of the cabinet. With this kit installed air circulates within the cabinet and heat is transferred to a customer-supplied water source.

Router cabinet houses up to 16 Rank 2/3 router modules. The router cabinet, shown in Figure 7, has the same footprint as the compute cabinet. All router cabinets are air-cooled. Like the compute cabinet, a blower draws air into the bottom of the cabinet, blows air across the modules, and exhausts air out the top of the cabinet.

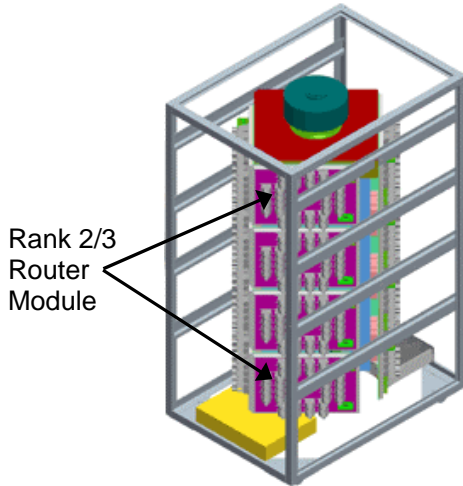


Figure 7. Router cabinet

Conclusion

The BlackWidow project uses a mixture of new and proven technologies to achieve outstanding performance. The combination of BlackWidow and Cray XT systems provides customers with a heterogeneous computing environment. The leading-edge chip technologies used enables many of the features outlined in this paper. Architectural features enable the system to scale, maintain a high level of reliability, and increase the scalar and vector performance. Packaging and cooling technologies that are used helped to reduce the over cost of the system.

Acknowledgments

The author would like to thank the BlackWidow development team for their contributions to this paper and for their dedication to the overall development of the BlackWidow system.

About the Author

Brick Stephenson is the Director of Engineering for the BlackWidow Project, Cray Inc. He has been with Cray for 25 years serving in various hardware-engineering roles. He can be reached at Cray Inc., 1020 Lowater Road, Chippewa Falls, Wisconsin 54729, E-Mail: bas@cray.com

References

- [1] Cray X1 <http://www.cray.com/products/x1>
- [2] Cray XT3 <http://www.cray.com/products/xt3>
- [3] C. Leiserson. Fat-trees: Universal networks for hardware efficient supercomputing. *IEEE Transactions on Computers*, C-34(10):892-901. October 1985
- [4] C. Clos, A Study of Non-Blocking Switching Networks. *The Bell System Technical Journal*, 32(2): 406-424, March 1953.
- [5] J. Kim, W.J. Dally, B. Towles, and A.K. Gupta, Microarchitecture of a high-radix router. In *ISCA'05: Proceedings of the 32nd International Symposium on Computer Architecture*, pp 420-431, Madison, WI. June 2005
- [6] S. Scott, D. Abts, J. Kim, and W.J. Dally, The BlackWidow High-Radix Clos Network. In *ISCA'06: Proceedings of the 33rd International Symposium on Computer Architecture*, Boston, MA. June 2006