



Architecture of Hitachi SR-8000

Uwe Küster

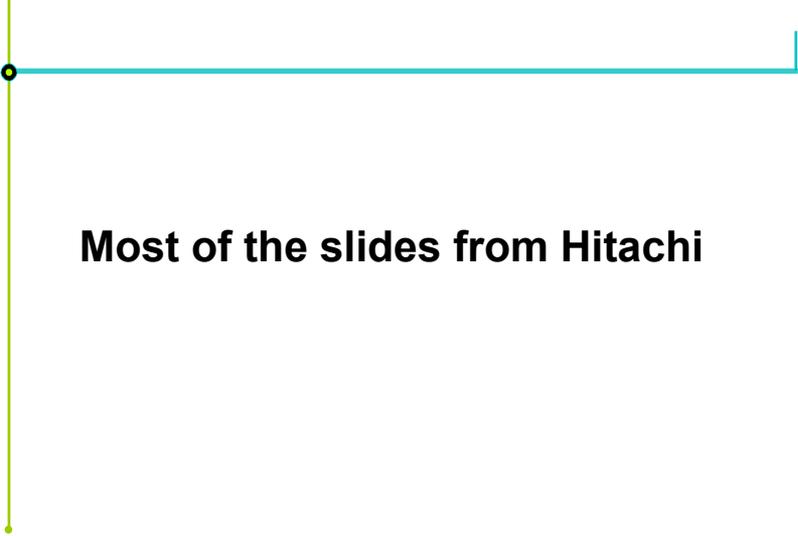
University of Stuttgart
High-Performance Computing-Center Stuttgart (HLRS)
www.hlrs.de



Uwe Küster
Slide 1

Höchstleistungsrechenzentrum Stuttgart

H L R I S 



Most of the slides from Hitachi



Uwe Küster
Slide 2

Höchstleistungsrechenzentrum Stuttgart

H L R I S 

the problem

- modern computer are data movers
- Floating Point Performance increases with 60% per year
(Moore's law)
- memory access is the bottleneck
 - ⊕ bandwidth is too small
 - ⊕ latencies are decreasing slowly
but increasing in clocks



Uwe Küster
Slide 3

Höchstleistungsrechenzentrum Stuttgart

H L R I S 

Latency from processor to various locations

	<u>seconds</u>	<u>clocks</u>	<u>future development</u>
clock rate	1-4 nsec		decreasing
register		1	
L1 cache		1-2	
L2 cache		4-8	
L3 cache		8-20	
memory	30-200 nsec	30-200	slowly decreasing
remote memory		500-700	slowly decreasing
shmem	~3 μsec		
MPI	6-20 μsec		slowly decreasing
disk	~8 msec		
tape	~30 sec		



Uwe Küster
Slide 4

Höchstleistungsrechenzentrum Stuttgart

H L R I S 

Behandlung von Latenzzeiten

- hiding of latency is necessary
- different possibilities:
 - ↑ Prefetching and Preloading from own and remote memory
 - ↑ Hardware-Multithreading
 - ↑ decreasing the number of transactions by pipelining
 - ↑ prozessor in memory (IRAM)



prefetch, example

```
do n=1,nmax
!   prefetch b(n+100), c(n+100) ! compiler generated
  a(n) = b(n) + c(n)
enddo
```



handling of latencies

compiler can hide latency by prefetching variables

but only in a semantical clear situation

modern processors support prefetching

explicit loops allow prefetching because future addresses can be calculated directly

branches are dangerous for prefetching if wrong data are prefetched and waste the bandwidth:

avoid unnecessary branches

calls break compiler optimization:

avoid frequent calls

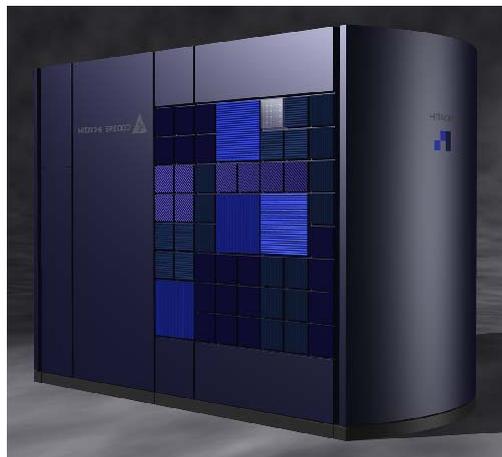


Uwe Küster
Slide 7

Höchstleistungsrechenzentrum Stuttgart

H L R I S 

16 processor SR8000

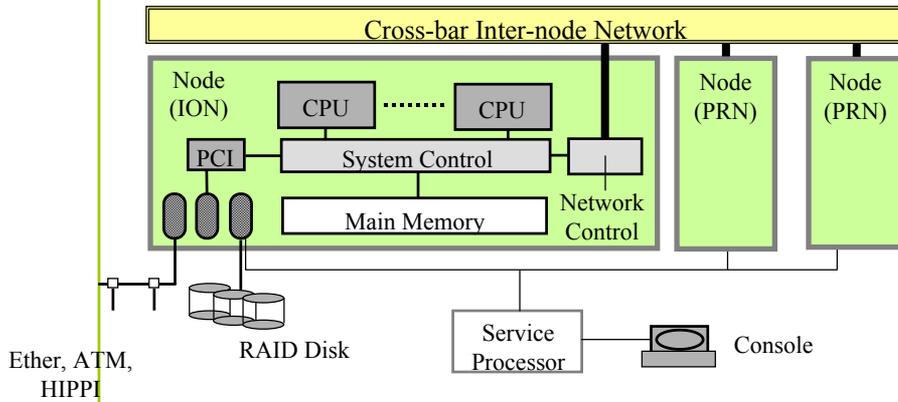


Uwe Küster
Slide 8

Höchstleistungsrechenzentrum Stuttgart

H L R I S 

System Architecture SR 8000

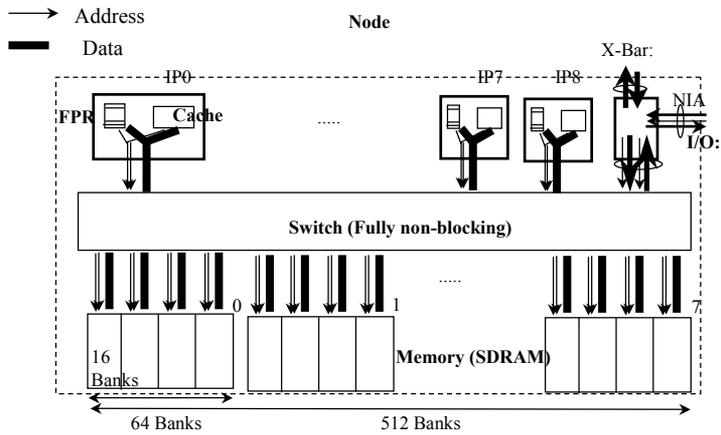


Uwe Küster
Slide 9

Höchstleistungsrechenzentrum Stuttgart

H L R I S

Hitachi SR-8000 node

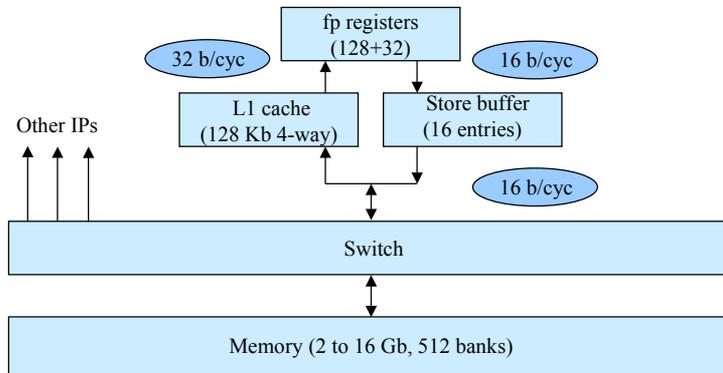


Uwe Küster
Slide 10

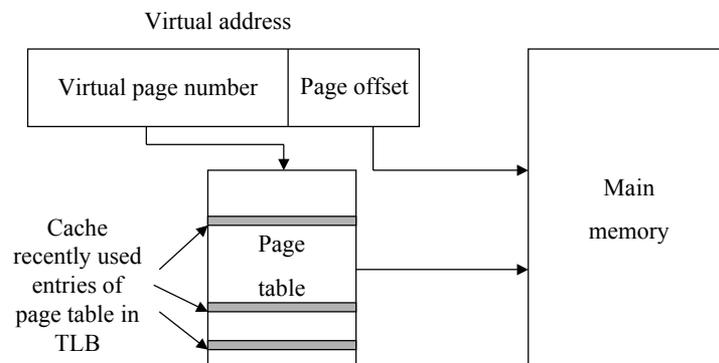
Höchstleistungsrechenzentrum Stuttgart

H L R I S

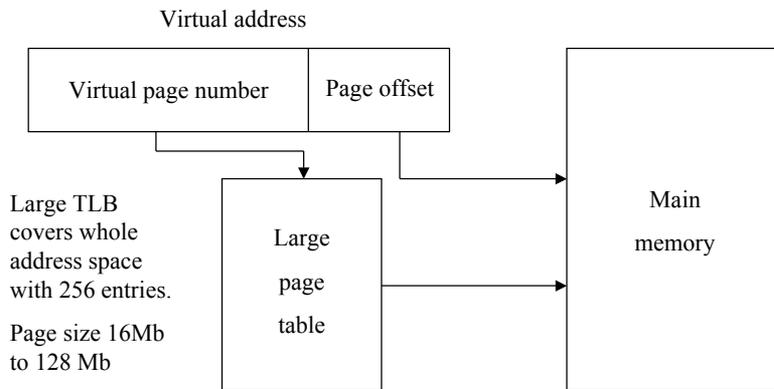
SR-8000: Memory Hierarchy



Address Translation



Large TLB



Uwe Küster
Slide 13

Höchstleistungsrechenzentrum Stuttgart

H L R I S

Topology SR-8000

3D Torus
8X8X8 (=512) processors largest size



Uwe Küster
Slide 14

Höchstleistungsrechenzentrum Stuttgart

H L R I S

Parameter debis-SFR Hitachi SR-8000

see Simmendinger: <http://thales.rz.st.dlr.de/cgi-bin/fom?file=16>

Peak performance	128 Gflop/s
Processors	16x8=128
Frequency	250 MHz
Memory	16x8=128 GB
Disk	440 GB
Inter-Node Bandwidth Hardware	1 GB/s
Inter-Node Bandwidth MPI/RDMA	950 MB/s
Inter-Node Bandwidth MPI	770 MB/s
Inner-Node Bandwidth Memory	32 GB/s/Node
L1 Cache	direct 4-way assoc
L1 Cache size	8x128KB/Node
L1 Cacheline size	128 B
Bandwidth Register - L1 Cache	32 B/cycle
Bandwidth L1 Cache - Memory	16 B/cycle
Bandwidth Register - Memory	8 B/cycle



Uwe Küster
Slide 15

Höchstleistungsrechenzentrum Stuttgart



Hitachi SR 8000:

- nodes with 8+1 processors
- processors have rotating registers (pseudo vector registers)
- 8 processors are cooperating for data parallel code
- or for shared memory parallelism



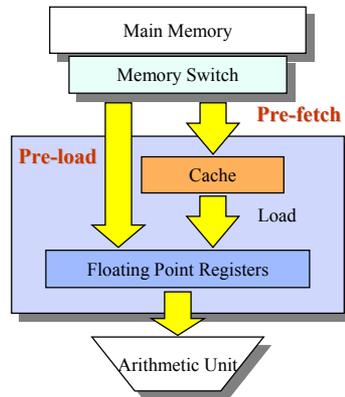
Uwe Küster
Slide 16

Höchstleistungsrechenzentrum Stuttgart

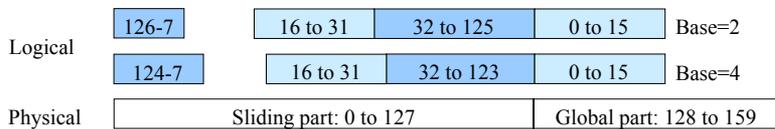


CPU Architecture

- 16 bytes/cycle memory BW
- 128 Kbyte L1 cache
- Pre-fetch ; 16 open transactions
- pre-load bypassing cache
- 160 FP registers
- 2 FP pipelines
- 4 flops/cycle



Slide Window Registers



- Registers for all instructions
- Registers for extended instructions only
- Fixed registers: 4, 8, 16, 32 (16 illustrated)
- Fixed + sliding = 128
- see also in IA-64



Instruction Set Extensions

- Load and store with extended registers
- Floating point arithmetic with extended registers
- Slide window control
- Pre-fetch and pre-load
- Thread start-up and finish
- Predicate instructions



Programming Models

Hardware	Programming Model	Example
Single CPU	Pseudo-vector processing	Vector application
Single Node	Independent processing on each IP	Compilation, parallel make
	Message passing	MPP application
	DO loop distribution with COMPAS	Vector application
Multiple Nodes	Parallel processing of independent blocks of code	
	Message passing	MPP application
Multiple Nodes	COMPAS and message passing	Vector parallel application



SR8000 Parallel Programming Models

Pseudo-Vector Processing:	PVP
COoperative MicroProcessors in single Address Space:	COMPAS
Message Passing:	MPI

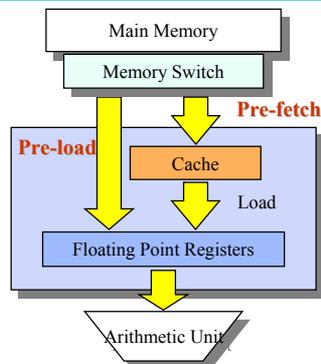


Uwe Küster
Slide 21 Höchstleistungsrechenzentrum Stuttgart

HLRS

Pre-fetch and Pre-load

- Pre-fetch: load cache line from memory to cache
- Pre-load: load one word from memory to register
- 16 streams



Uwe Küster
Slide 22 Höchstleistungsrechenzentrum Stuttgart

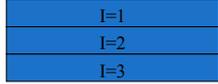
HLRS

Software Pipelining

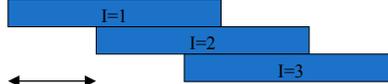
No SWPL



Infinite resource

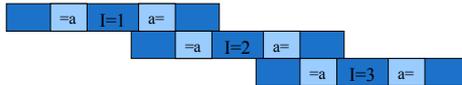


Finite resource



Initiation interval

Recurrence



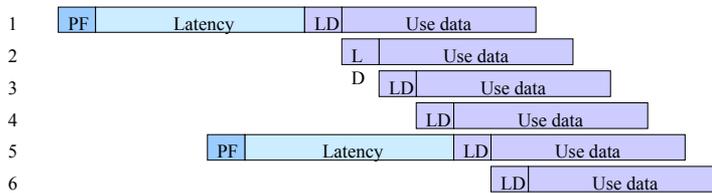
Resources:

registers, f.p. units,
instruction issue,
memory bandwidth
etc



Pre-fetch

Iteration

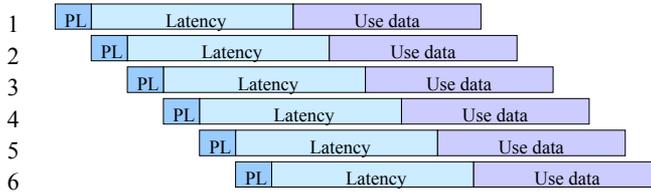


- Pre-fetch 128 bytes to cache
- Follow by LD to register



Pre-load

Iteration

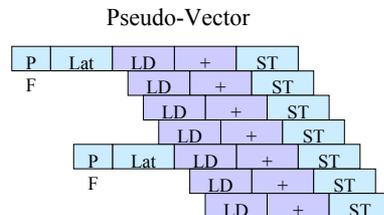
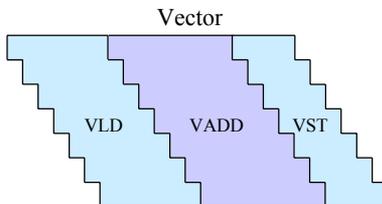


- Pre-load 8 bytes to register
- LD not required

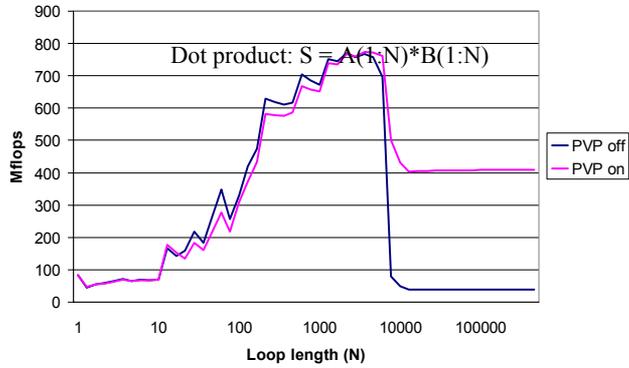


Pseudo-vector Processing

$$A(:) = A(:) + N$$



Effect of PVP

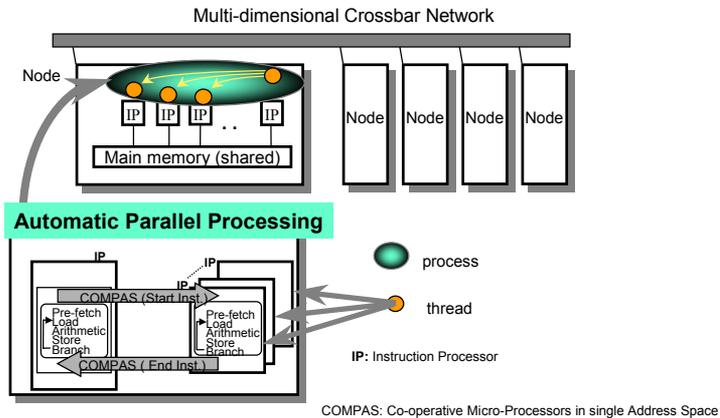


Uwe Küster
Slide 27

Höchstleistungsrechenzentrum Stuttgart

H L R I S

COMPAS



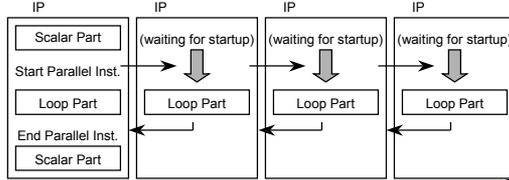
Uwe Küster
Slide 28

Höchstleistungsrechenzentrum Stuttgart

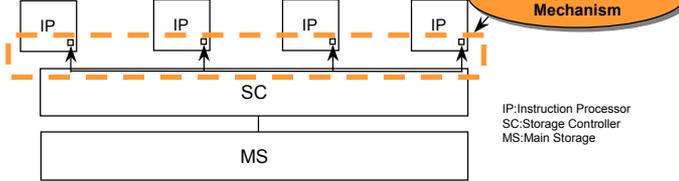
H L R I S

Hardware Support

Software



Hardware Support

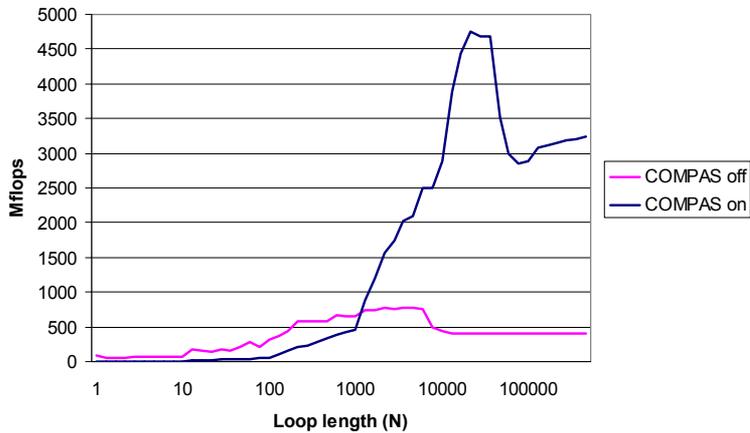


Uwe Küster
 Slide 29

Höchstleistungsrechenzentrum Stuttgart



Effect of COMPAS: Dot product: $S = A(1:N) * B(1:N)$

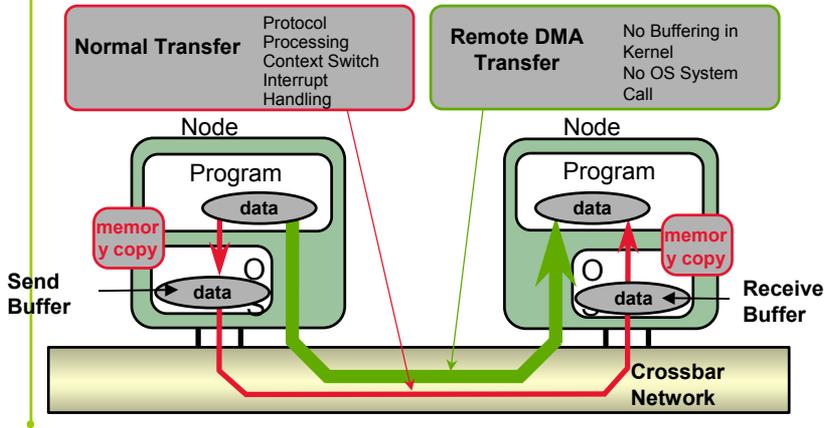


Uwe Küster
 Slide 30

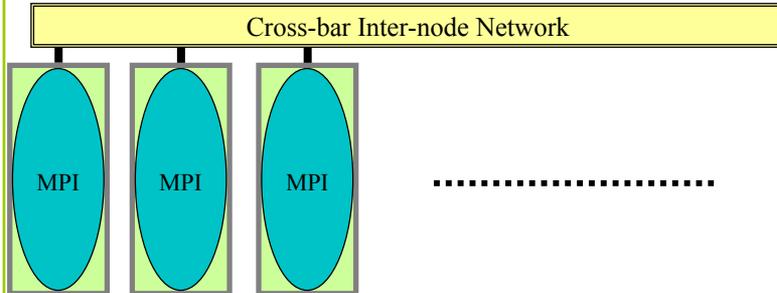
Höchstleistungsrechenzentrum Stuttgart



Remote DMA



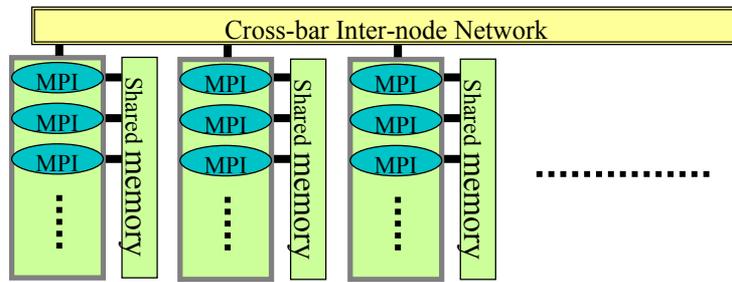
Inter-node MPI



One MPI process per node; RDMA transfer possible



Intra-node MPI



One MPI process per IP; RDMA transfer not possible



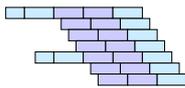
Uwe Küster
Slide 33

Höchstleistungsrechenzentrum Stuttgart

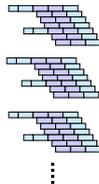
H L R I S 

SR8000: Parallelism on all levels

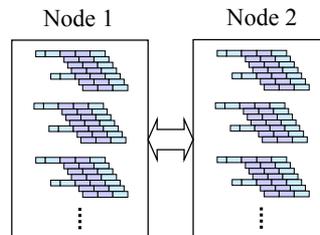
Instruction level
(PVP)



Multi-thread
(COMPAS)



Message passing
(MPI)



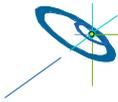
Uwe Küster
Slide 34

Höchstleistungsrechenzentrum Stuttgart

H L R I S 

Key Features of SR8000

- High performance RISC CPU with PVP
- High performance node with COMPAS
- High sustained memory bandwidth
- High scalability with fast network
- Low energy and space requirements



Uwe Klüster
Slide 35

Höchstleistungsrechenzentrum Stuttgart

H L R I S 