



SGI Message-Passing Status and Plans

Karl Feind

kaf@sgi.com

SGI Parallel Communication
Team

SUMMIT 2001

43rd CUG Conference on High Performance Computing & Visualization



MPT Themes

sgi™

- Performance
- Platforms and Interconnects
- Standards



SUMMIT 2001

43rd CUG Conference on High Performance Computing & Visualization

Performance Features in **sgi** MPT

- Low latency and high bandwidth.
- Fetchop-assisted fast message queuing
- Fast fetchop tree barriers
- MPI and SHMEM one-sided communication
- Large SSI support
- Automatic NUMA placement
- Optimized MPI collectives
- Internal MPI statistics reporting
- Integration with PCP
- Single copy send/recv transfers
- Runtime MPI tuning



SUMMIT 2001

43rd CUG Conference on High Performance Computing & Visualization

Communication on NUMALink



MPI Performance on 400 MHz Origin 3000)

send/rcv latency	5.5 usec
put/get latency	1 usec
Peak bandwidth (pcopy)	280 Mbytes/sec
Peak bandwidth (BTE)	600 Mbytes/sec
All communicate peak bandwidth per transfer	170 Mbytes/sec



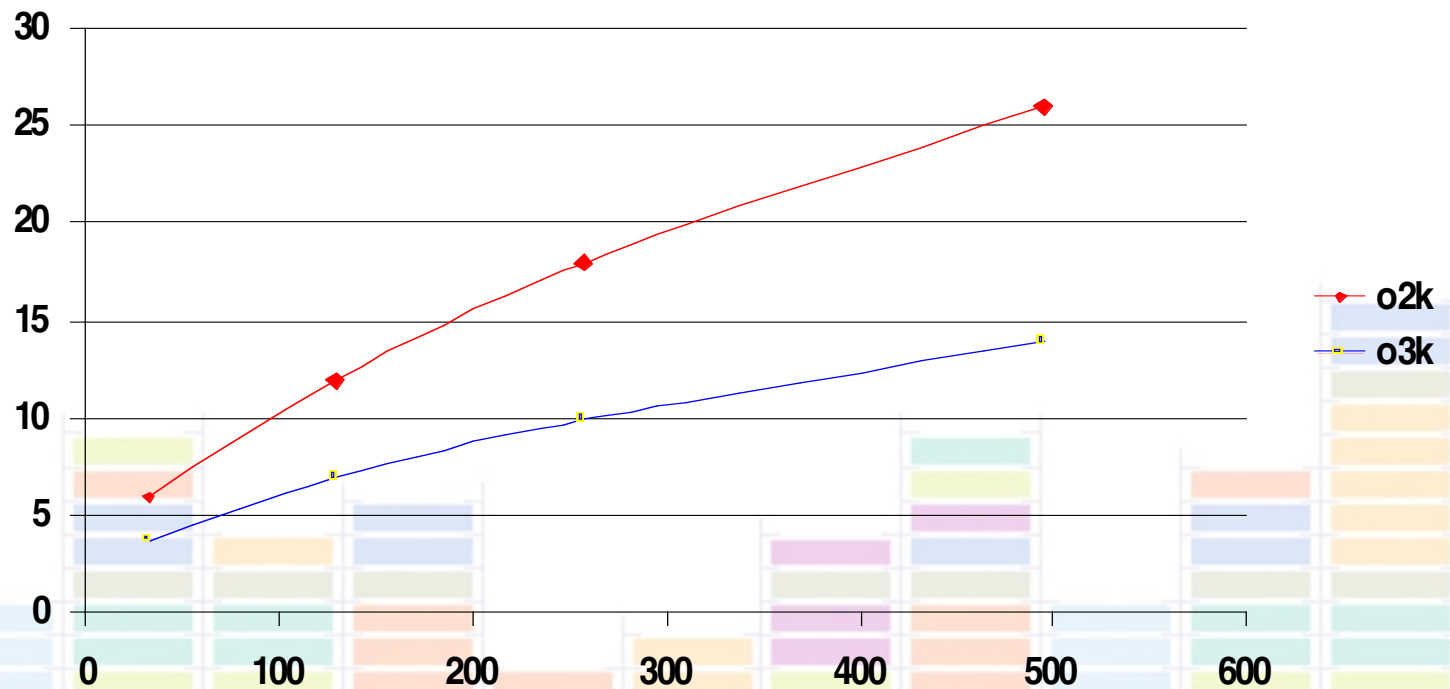
SUMMIT 2001

43rd CUG Conference on High Performance Computing & Visualization

Barrier Synchronization Time on O2K and O3K



Time (μsec)



MPI_BAR_DISSEM=on

Number of processes

SUMMIT 2001

43rd CUG Conference on High Performance Computing & Visualization



Assign MPI Ranks to Physical CPUs



- **Environment Variable Syntax**
 - **setenv MPI_DSM_CPULIST 0-15**
 - **setenv MPI_DSM_CPULIST 0,2,4,6**
- **Also SMA_DSM_CPULIST**
- **Maps ranks 0 - N onto the physical CPUs in the order specified**
- **Useful only on quiet systems**
- **Easier than using dplace command**



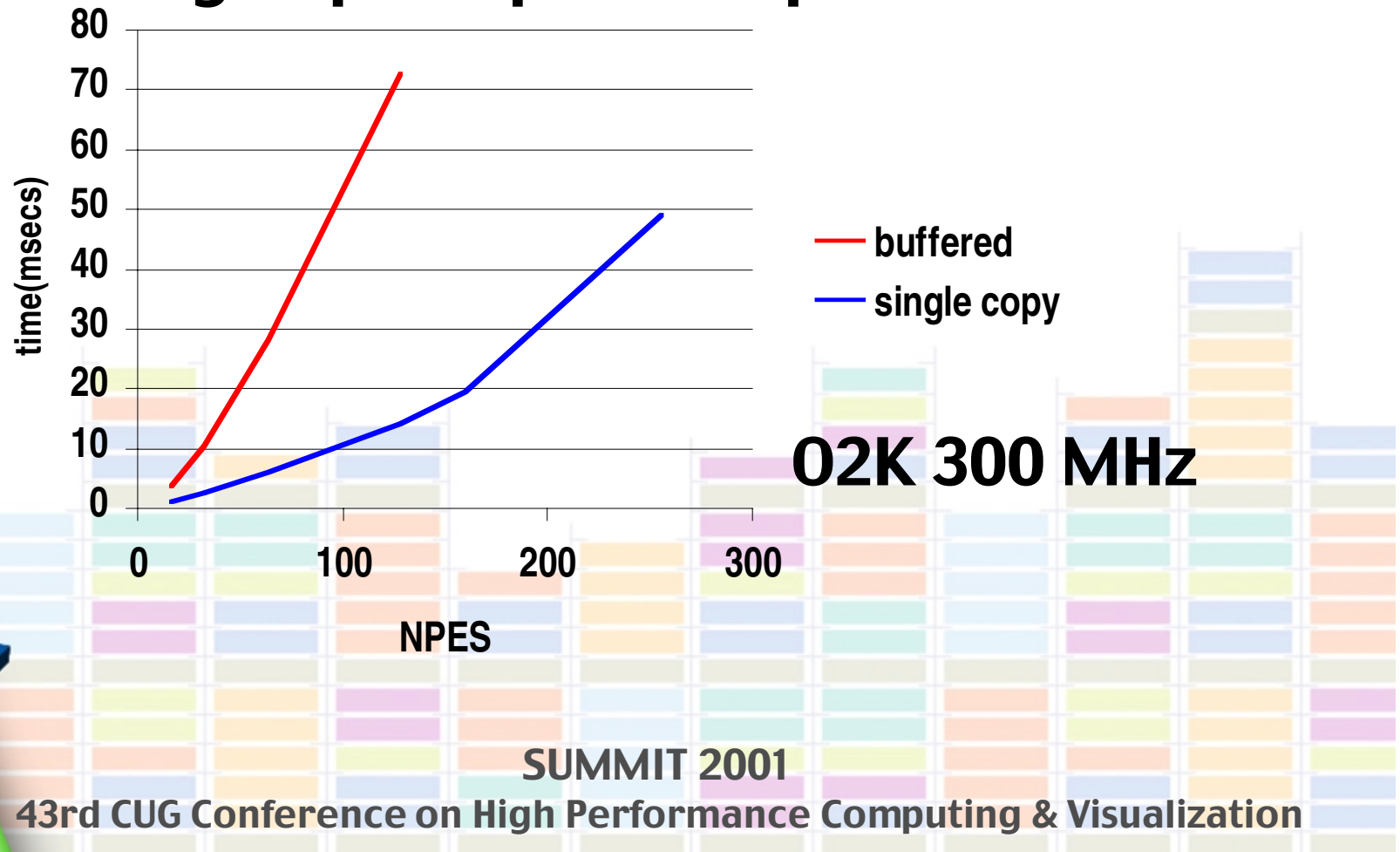
SUMMIT 2001

43rd CUG Conference on High Performance Computing & Visualization

Single Copy Speed-up



Alltoallv type communication pattern
using explicit point to point calls



Single Copy Send/Receive



- **Send buffer must be globally accessible (ABI 64 required)**
 - common block
 - symmetric heap (shmalloc/SHPALLOC)
 - global heap (ALLOCATE statement with SMA_GLOBAL_ALLOC set)
- **Set MPI_BUFFER_MAX to 2048**
- **Best used with MPI_Isend or MPI_Bcast**
- **Best if sender does not immediately wait for send completion.**
- **Little payoff below 8 Kbyte messages**

SUMMIT 2001

43rd CUG Conference on High Performance Computing & Visualization



Reducing Run-Time Variability



- **Recommended algorithm for workload manager launch of MPI jobs in SSI:**
 - Batch scheduler creates a cpuset
 - Batch scheduler launches mpirun into cpuset
 - MPI job is confined within cpuset during execution by virtue of fork/exec/cpuset semantics
 - MPI job performs automatic NUMA placement within the cpuset
 - When MPI job completes, the batch scheduler destroys the cpuset.

- **NOTE**

- Use exclusive cpusets or restrict interactive use of the system.

SUMMIT 2001

43rd CUG Conference on High Performance Computing & Visualization



Platforms and Interconnects: MPI



■ MIPS

- Single kernel NUMALink
- Partitioned NUMALink (available June 2001)
- GSN (libst 2.0 work planned in July 2001)
- Myrinet
- Sockets
- HIPPI

■ SNIA

- Single kernel NUMALink (prototype working)
- Partitioned NUMALink (prototyping in late 2001)
- Myrinet (prototype running on IA64)
- Sockets (prototype working)

SUMMIT 2001

43rd CUG Conference on High Performance Computing & Visualization



Platforms and Interconnects: SHMEM



■ MIPS

- Single kernel NUMALink
- Partitioned NUMALink (planned Dec 2001)

■ SNIA

- Single kernel NUMALink (planned Sep 2001)
- Partitioned NUMALink (planned Nov 2001)



SUMMIT 2001

43rd CUG Conference on High Performance Computing & Visualization

Platforms and Interconnects: PVM



- **MIPS**

- PVM support is retired after MPT 1.6 (2002)

- **SNIA**

- SGI will not provide PVM on SNIA



SUMMIT 2001

43rd CUG Conference on High Performance Computing & Visualization

MPI-2 Features Planned



- **MPI I/O enhancements: MPI_Wait integration**
- **MPI-2 datatypes: replacements for deprecated MPI-1 datatypes**
- **Expanded one-sided communication**
- **Process spawn**



SUMMIT 2001

43rd CUG Conference on High Performance Computing & Visualization

SGI Message-Passing References



- “relnotes mpt” gives information about new features and how to install MPT
- “man mpi” tells about all environment variables
- “man shmem” tells about the SHMEM API
- *MPI Programmer’s Manual* viewable with *insight* and on the web at <http://techpubs.sgi.com>
- MPT web page:
<http://www.sgi.com/software/mpt/>



SUMMIT 2001

43rd CUG Conference on High Performance Computing & Visualization