# Extreme Scale Computing

William Gropp
**University of Illinois at Urbana-Champaign**
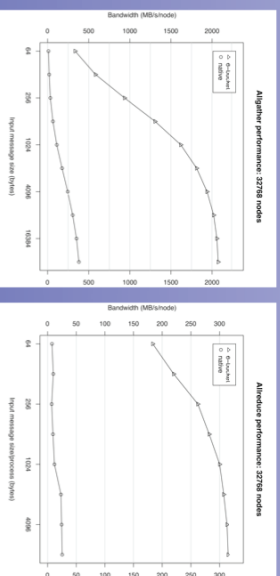
## Overview

Research Interests:
Programming models for parallel computing
Scalable numerical methods
Extreme scale computing
Current projects
MPI for petascale
Scalable Algorithms
Exascale Computing
"Blue Waters", the NSF Track 1 Petascale System
Software Projects
MPICH2 – A High Performance, Scalable, Portable Implementation of MPI
pnetCDF – A parallel version of the netCDF 3 file format and utilities
PETSc – A Scalable numerical library for solving large systems of linear and nonlinear equations
Two examples follow

**Paul Sack**

## Enhancing Performance of MPI Collective Operations

Collective communication and computation are important operations
Implementations have often minimized the data moved and are optimal in this sense
Interconnection networks, particularly at extreme scale, are complex and have limited bisection bandwidth compared with a complete graph
Goal: Minimum *time* collective operations
Idea: Additional Data Motion can *reduce* contention in network, by avoiding contention



Allgather performance, 32768 nodes



Allreduce performance, 32768 nodes

## Addressing Scalability Challenges of Algebraic Multigrid

Illinois: **Hormozd Gahvari** and Luke Olson
LLNL: Martin Schulz and Ulrike Meier Yang
IBM: Kirk Jordan

Goal: Access and enhance the scalability of an "optimal" numerical method for extreme scale systems

### Algebraic Multigrid

- **Algebraic multigrid (AMG), allows for the fast solution of large problems by using coarse-grid approximations that involve far fewer points:**

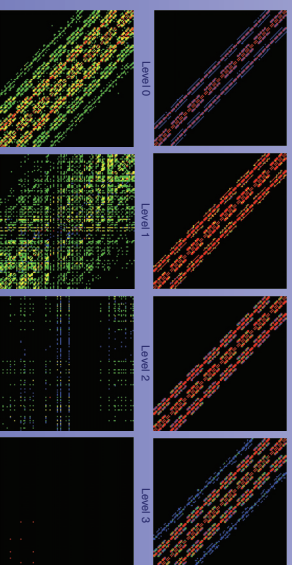- **Work per unknown remains constant, which is great for supercomputing**

### Scalability Challenges

- **AMG scales well on IBM Blue Gene machines, but has difficulties on multicore clusters**

- **Quick comparison between two machines, Hera (multicore Linux cluster at LLNL) and Intrepid (IBM Blue Gene/P at ANL) is highly illustrative**

- **Results here are for one V-cycle on a 3D 7-point Laplace problem on 128 processors with 62,500 points and one MPI process per core:**

| Level | Unknowns | Hera | Intrepid |
|-------|----------|------|----------|
| 0 | 8,000,000 | $2.30 \times 10^{-3}$ s | $5.19 \times 10^{-3}$ s |
| 1 | 614,521 | $8.28 \times 10^{-1}$ s | $1.42 \times 10^{-3}$ s |
| 2 | 120,607 | $6.52 \times 10^{-1}$ s | $4.38 \times 10^{-3}$ s |
| 3 | 13,643 | $1.24 \times 10^{-1}$ s | $1.53 \times 10^{-3}$ s |
| 4 | 1,509 | $1.87 \times 10^{-1}$ s | $1.27 \times 10^{-3}$ s |
| 5 | 212 | $2.02 \times 10^{-1}$ s | $1.72 \times 10^{-3}$ s |
| 6 | 23 | $2.75 \times 10^{-1}$ s | $5.36 \times 10^{-5}$ s |
| 7 | 3 | $5.44 \times 10^{-1}$ s | $2.22 \times 10^{-5}$ s |
| total | | $9.18 \times 10^{-2}$ s | $7.55 \times 10^{-2}$ s |

Hera:
-Slowdown on levels 3 → 5
-Level 5, with only 212 points, is nearly as slow as level 0!

Intrepid:
-Slowdown on level 5.

- **This is driven by increasing amounts of interprocessor communication:**



Level 0, Level 1, Level 2, Level 3, Level 4, Level 5, Level 6, Level 7

## Performance Model

- **To help us understand what we are seeing, develop performance model for AMG solve cycle**

- **Baseline model: α-β (latency-bandwidth) with parameters**
  - $P$ – number of processes
  - $C_i$ – number of grid points in level i
  - $l_i$ – maximum grid index
  - $s_i$, $\hat{s}_i$ – average number of nonzeros per row in solve and interpolation operators, respectively
  - $p_i$, $\hat{p}_i$ – maximum number of sends per active process in solve and interpolation operators, respectively
  - $n_i$, $\hat{n}_i$ – maximum number of elements sent per active process in solve and interpolation operators, respectively
  - $t_i$ – time per floating-point operation on level i

- **Runtime at each level is sum of:**



- **Take architectural features into account with penalties:**
  - Distance of communication: add γ term (charge distance × network diameter to each message)
  - Lower effective bandwidth: multiply β by Hardware Bandwidth/MPI Bandwidth
  - Multicore latency penalty: multiply α by (σ·P)/P   (σ = cores per node, P_i = no. active processes at level i)
  - Multicore distance penalty: multiply γ by (σ·P)/P

- **Results spotlight impact of architecture on performance:**



Intrepid, 8192 Processes

Intrepid:
- Not too much degradation on coarse grids
- Only limited bandwidth penalty applies



Hera, 1024 Processes

Hera:
- Massive degradation on coarse grids
- All penalties apply – lots of network contention and switching delays from multicore nodes (16 cores/node)

- **Future machines will be more multicore....**

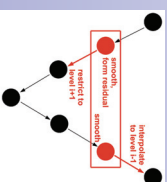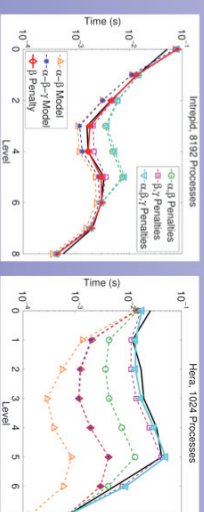This work was performed under the auspices of the U.S. Department of Energy by Lawrence Livermore National Laboratory under Contract DE-AC52-07NA27344.