### System Software Issues for the Future

Alok Choudhary, Professor **Director: Center for Ultra-Scale** Computing and Information Security Dept. of Electrical & Computer Engineering And Kellogg School of Management **Northwestern University** 

choudhar@ece.northwestern.edu

## HPC Foundation



# Runtime Systems – Challenges

#### User specifies how **Complex non-portable optimization space** streaming/ Small/large configuration s/w layer Terabytes Mair Regular/irregular Memory Local/remote **MEMs** Low Power DRAM user burdened Petabytes High Performance Disk Ineffective interfaces Non-communicating layers Holographic Memory Massive Arrays of Idle Disks Exabyte **Tape Silos**

#### **Emerging Storage Hierarchy**

| User application        | Access patterns: shared files, individual files, data partitioning, check-<br>pointing, data structures, inter-data relationship  |
|-------------------------|---|
| HDF5 pnetCDF            | Data types (byte-alignment), data structures (flexible dimensionality),<br>hierarchical data model  |
| MPI-IO                  | Collectives, independents<br>I/O hints: access style (read_once, write_mostly, sequential, random,),<br>collective buffering, chunking, striping  |
|                         | Caching, fault tolerance, read-ahead, write-behind, I/O load balance, wide-area, heterogeneous FS support, thread-safe  |
| Client-side file system | Open mode (O_RDONLY, O_WRONLY, O_SYNC), file status, locking,<br>flushing, cache invalidation<br>Machine dependent: data shipping, sparse access, double buffering  |
|                         | application-aware caching, pre-fetching, file grouping, "vector<br>of bytes", flexible caching control, object-based data<br>alignment, memory-file layout mapping, more control over<br>hardware, Shared file descriptors, |
| Server-side file system | Read-ahead, write-behind, metadata management, file striping, security, redundancy  |
|                         | Group locks, flexible locking control, scalable metadata<br>management, zero-copying, QoS, Shared file descriptors,   |
| Storage system          | Access base on : file blocks, objects Scheduling, aggregation<br>Active storage: data filtering,object-based/hierarchical<br>storage management, indexing, mining, power-management   |

# Decouple "What" from "How"

#### Current

Goal



## Caching Example: Direct Access Caching at Compute Nodes (BG)



April 20, 2006

### Coherence Control at I/O Nodes



### Data at I/O Nodes

#### Logical partitioning view of a file

block 0 block 1 block 2 block 3 block 4

#### **Distributed metadata**



#### Cache pages at I/O nodes



April 20, 2006