



*System Design Considerations*  
*of*  
*Big Data Processing*

**Xian-He Sun**

Illinois Institute of Technology

Chicago, Illinois

[sun@iit.edu](mailto:sun@iit.edu)



# Addressing the Big Data Challenges

---

Trends indicate that the “data tsunami” and “memory-wall” will continue, waiting for miracle is not a answer

Big-Data problem is a HPC problem:

- ❑ **Data access & Interface**

Need **rethinking** from the **data-centric** point-of-view in:

- **Understanding** the system, application, and algorithm relevant to data access
- **Optimizing** current systems
- **Developing** new system architectures
- **Developing integrated solutions**
  - ❑ Algorithm, programming model, system, architecture, co-design
  - ❑ In situ application-aware data access optimization



Understanding

# MEMORY SYSTEM BEHAVIOR



# Improve via Memory Hierarchy

Capacity  
Access Time, Bandwidth

## CPU Registers

<8KB  
<0.2~0.5 ns, 500~800 GB/s/core

## Cache

<50MB  
1-10 ns, 50~150GB/s/core

## Main Memory

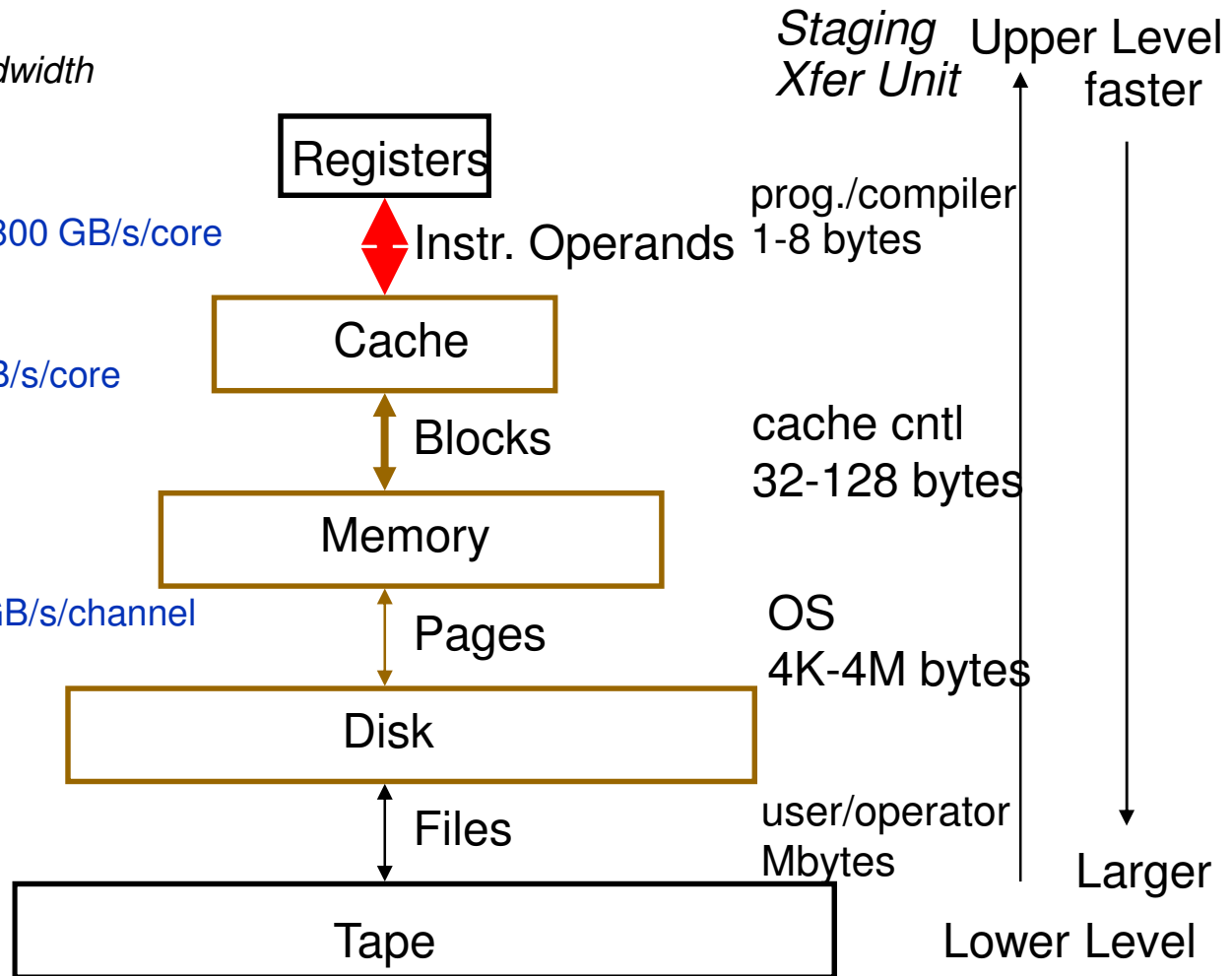
Giga Bytes  
50ns-100ns 5~10GB/s/channel

## Disk

Tera Bytes, 5 ms  
100~300MB/s

## Tape

Peta Bytes or  
infinite  
sec-min





# Solution: Memory Hierarchy & Parallelism

Multi-core  
Multi-threading  
Multi-issue



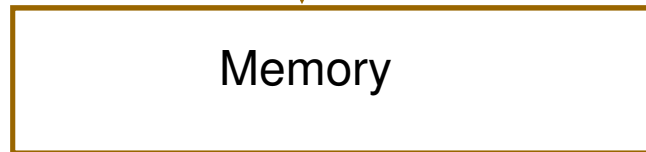
Out-of-order Execution  
Speculative Execution  
Runahead Execution

Multi-banked Cache  
Non-blocking Cache  
Multi-level Cache



Pipelined Cache  
Data Prefetching  
Write buffer

Multi-channel  
Multi-rank  
Multi-bank



**Input-Output (I/O)**

*Parallel File System*

Disks

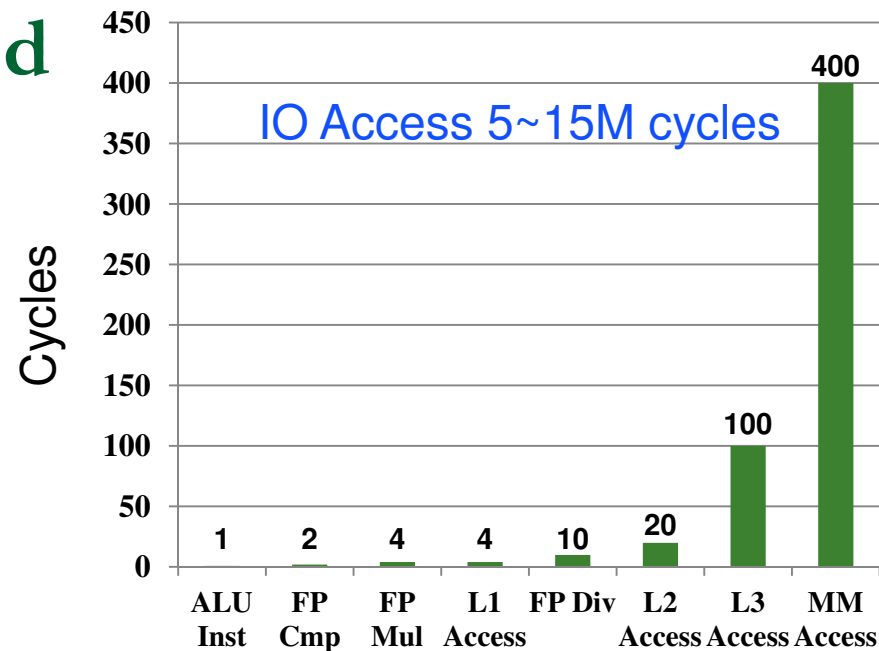


# Assumption of Current Solutions

- ❑ Memory Hierarchy: **Locality**
- ❑ Concurrency: **Data access pattern**
  - Data stream

**Extremely Unbalanced  
Operation Latency**

**Performances vary  
largely**





# Existing Memory Metrics

- ❑ Miss Rate(MR)
  - {the number of miss memory accesses} over {the number of total memory accesses}
- ❑ Misses Per Kilo-Instructions(MPKI)
  - {the number of miss memory accesses} over {the number of total committed Instructions  $\times$  1000}
- ❑ Average Miss Penalty(AMP)
  - {the summary of single miss latency} over {the number of miss memory accesses}
- ❑ Average Memory Access Time (AMAT)
  - $AMAT = \text{Hit time} + MR \times AMP$
- ❑ **Flaw of Existing Metrics**
  - Focus on a single component or
  - A single memory access

**Missing memory parallelism/concurrency**



# The Introduction of APC

- Access Per Cycle (APC)
  - $APC = A/T$
- APC is measured as the number of memory accesses per cycle
  - Measures the overall memory system performance
  - Each memory level has its own APC value
- Benefits of APC
  - Separate memory evaluation from CPU evaluation
  - A better understanding of memory system as a whole
  - A better understanding of the match between computing capacity and memory system performance

X.-H. Sun and D. Wang, "APC: A Performance Metric of Memory Systems",  
ACM SIGMETRICS Performance Evaluation Review, Volume 40 , Issue 2, 2012.





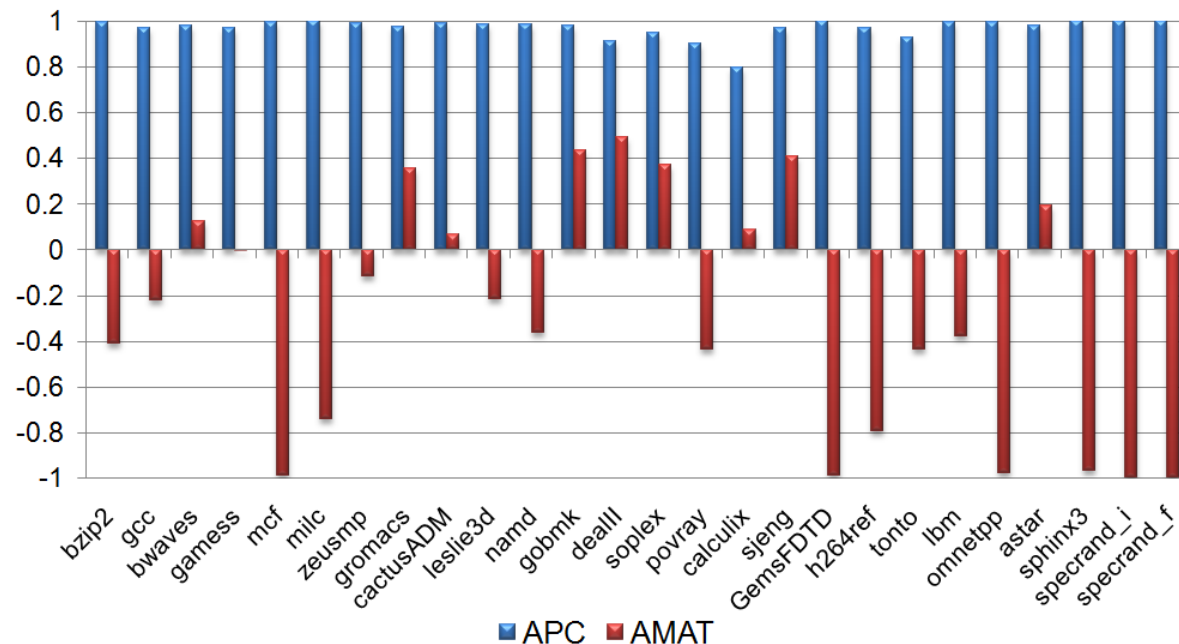
# APC Measurement

- The difficulty is measuring the total cycle  $T$ 
  - Hundreds of memory accesses co-exist the memory system
  
- Measure  $T$  based on the *overlapping mode*
  - When there are several memory accesses co-existing during the same clock cycle,  $T$  only increases by one
  - Measure the concurrence
  - Measure the concurrence at each level





## APC & IPC: Changing Cache Parallelism



- Changing the number of MSHR entries (1→2→10→16)
- APC still has the dominant correlation, with average value of 0.9656
- AMAT does not correlate with IPC for most applications
  - APC record the CPU blocked cycles by MSHR cycles
  - AMAT cannot records block cycles, it only measure the issued memory requests



# Exhausted Testing

- With different benchmarks, and with different configurations
- With advanced cache technologies
  - Non-block cache
  - Pipelined cache
  - Multi-port cache
  - Hardware prefetcher
- With single core or multicore
- **APC always has the highest CC values among all the memory metrics**

D. Wang, X.-H. Sun "Memory Access Cycle and the Measurement of Memory Systems",  
IEEE Transactions on Computers, 2013



# APC Applications

- Identify the contribution of concurrency
- Quantitatively **define data intensiveness**
- Find the lowest level that has a dominating correlation with IPC
- Provide a mean to study the **matching** between memory organization and microprocessor architecture,
- Provide a mean to study the **matching** between memory organization and a given application
- Co-Design
- **C-AMAT= 1/APC (Architecture Analysis)**

D. Wang and X.-H. Sun, "Concurrent Average Memory Access Time",  
Illinois Institute of Technology Technical Report (IIT/CS-SCS-2012-05), 2012



Understanding

# APPLICATION BEHAVIOR



# Data Access is Application Dependent

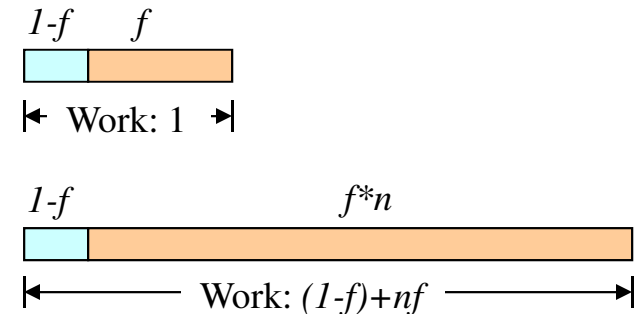
- Conventional algorithm analysis
  - Floating point operation
  
- Data-centric algorithm analysis
  - Floating point operation
  - Memory requirement
  - Data reuse rate
  - Data access/movement pattern





# The Memory-bounded Speedup

- Tacit assumption in Amdahl's law
  - The problem size is **fixed**
  - The speedup emphasizes **time reduction**
- Gustafson's Law, 1988
  - Fixed-time speedup model



$$\begin{aligned}
 \text{Speedup}_{\text{fixed-time}} &= \frac{\text{Sequential Time of Solving Scaled Workload}}{\text{Parallel Time of Solving Scaled Workload}} \\
 &= (1 - f) + nf
 \end{aligned}$$

- Sun and Ni's law, 1990
  - Memory-bounded speedup model

$$\begin{aligned}
 \text{Speedup}_{\text{memory-bounded}} &= \frac{\text{Sequential Time of Solving Scaled Workload}}{\text{Parallel Time of Solving Scaled Workload}} \\
 &= \frac{(1 - f) + fG(n)}{(1 - f) + fG(n)/n}
 \end{aligned}$$



# Contribution of Memory-bounded Speedup

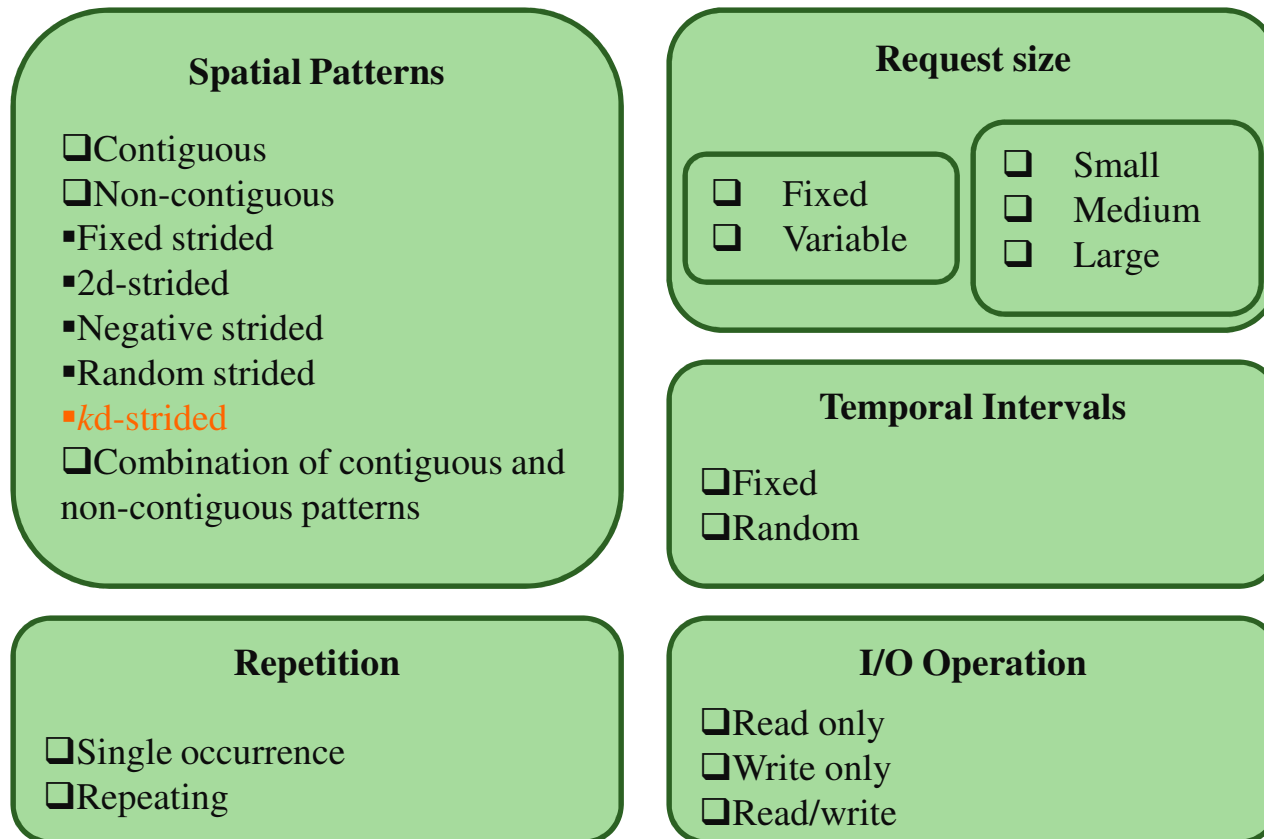
- Data-centric thinking
- Where the memory-bound function  $W = G(M)$  provide
  - $W$ , the work in floating point operation
  - $M$ , the memory requirement
  - $G$ , the data reuse rate
  - Enough for memory hierarchy, but not parallelism
- Need to find the data access/movement patterns for data access concurrency
  
- Dense Linear Algebra,  $M$  memory,  $M^{3/2}$  work
- FFT,  $M$  memory,  $O(M \log(M))$  work
- $G(pM) > pW$ , can lead to large increase in execution time
  - (ex) 10K x 10K matrix factorization: 800MB, 1 hr in uniprocessor  
with 1024 processors, 320K x 320K matrix, 32 hrs





# Data Access Signature: Patterns and Notation

- Comprehensive access pattern classification
- Implemented for I/O and memory access (MPI datatype)





## I/O Trace Signature

- Description of a sequence of I/O accesses in a pattern
- Form:  $\{I/O \text{ operation, init position, dimension, } ([\text{offset pattern}], \text{request size pattern}), \text{pattern of number of repetitions}\}, [\dots], \# \text{ of repetitions}\}$

## Pattern Signature

- provides a simple description that explains the nature of a pattern
- Form:  $\{I/O \text{ operation, } \langle \text{Spatial pattern, Dimension} \rangle, \langle \text{Repetitive behavior} \rangle, \langle \text{Request size} \rangle, \langle \text{Temporal Intervals} \rangle\}$

## I/O Pattern Detection

- Developed a pattern detection tool
- Five pattern detectors for finding patterns among initial positions, offsets, request sizes, temporality, and repetitions
- Outputs I/O Signature that can be used for prefetching, data layout, and data reorganization



# IOSIG: An I/O Characterization Tool

Website: [www.cs.iit.edu/~scs/iosig/](http://www.cs.iit.edu/~scs/iosig/)

*Goal:* To provide a better understanding of parallel I/O accesses and information to be used for optimization techniques.

*Two steps:*

**Trace collection**

- Collects parallel I/O calls of an application
- Does not require any code modification

**Trace analysis**

- Analyzes the collected information to give a clear understanding of I/O behavior of the application
- Handles trace files in any text-based format

Y. Yin, et.al, "Boosting Application-Specific Parallel I/O Optimization Using IOSIG", in Proc. of IEEE/ACM CCGrid, 2012.



# More about IOSIG

- Additional contributions
  - Data access pattern categories
  - Local and global I/O signatures
- Applications

---

Prefetching	<b>SC'08</b>	S. Byna, Y. Chen, X.-H. Sun, R. Thakur, and W. Gropp <i>Parallel I/O Prefetching Using MPI File Caching and I/O Signatures</i>
Data Layout	<b>HPDC11</b>	H. Song, Y. Yin, Y. Chen, X.-H. Sun, <i>A Cost-intelligent Application-specific Data layout Scheme for Parallel File Systems</i>
Data Organization	<b>PDSW11</b>	J.He, H. Song, X.-H. Sun, Y. Yin, and R. Thakur Pattern-aware File Reorganization in MPI-IO
Data Replication	<b>IPDPS13</b>	Y. Yin, J. Li, J. He, X.-H. Sun, and R. Thakur, <i>Pattern-Direct and Layout-Aware Replication Scheme for Parallel I/O Systems</i>

---

## Data Compression

- Website: [www.cs.iit.edu/~scs/iosig/](http://www.cs.iit.edu/~scs/iosig/)



I/O System

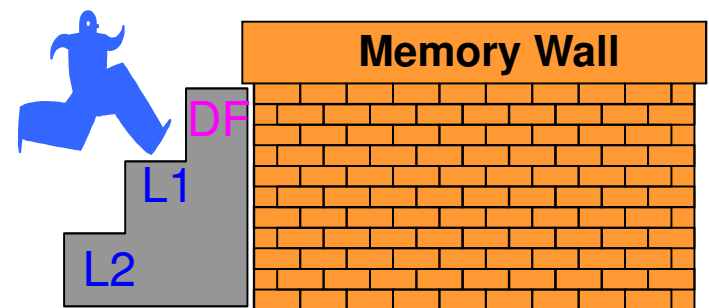
# OPTIMIZATION *FROM DATA POINT OF VIEW*



# In Situ Application-aware Optimization

---

- Data access is application dependent
- Dynamic, application-aware
- Data access pattern, feed back, control
- Integrating language, memory system, and hardware/software infrastructures
- Understanding design trade-off





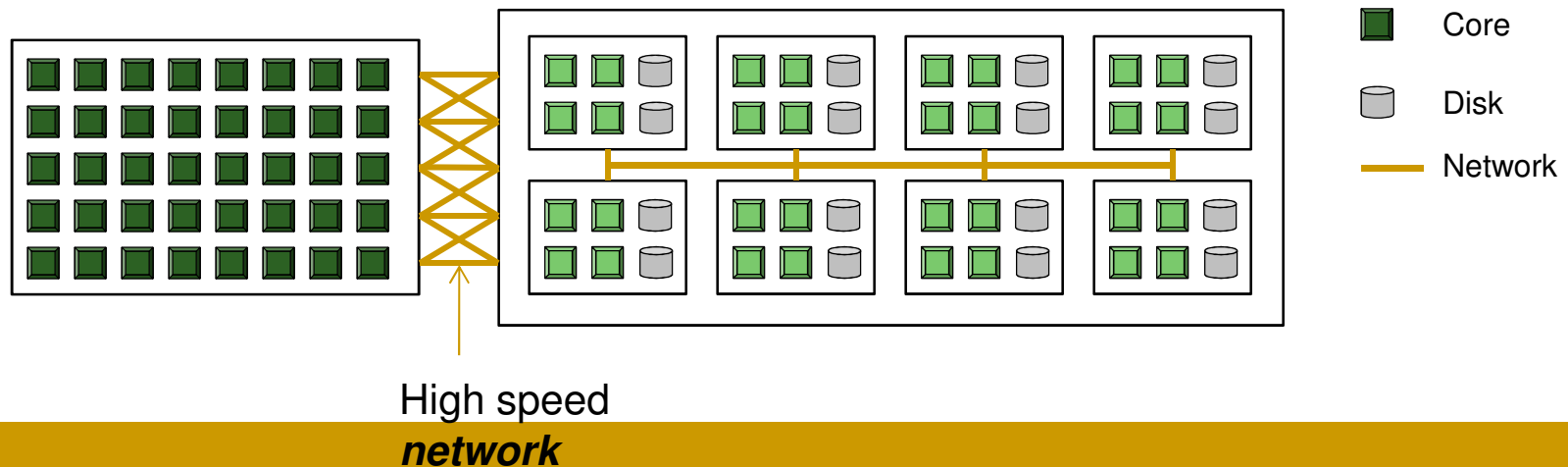
# Ours Method: Application-Aware I/O (1)

## Decoupled-Execution Paradigm:

- ❑ Handle computation- and data- intensive phases separately
- ❑ One interface-Two systems, transparent to users
- ❑ Integration, scheduling, optimization

**Supercomputer** or  
**many-core computing system**  
for execution of computing  
intensive part of an application

**Data cloud** or **storage cluster**  
for execution of data  
intensive part of an application





# Application-Aware I/O Optimization (2)

## Interoperability between different file systems

- Enable MPI Apps to access data-intensive file systems
- HPC-Cloud, Data-Cloud



H. Jin, X.-H. Sun, et. al, "CHAIO: Enabling HPC Applications on Data-Intensive File Systems", **ICPP2012**.





# Conclusion: Big Data Rethinking

---

- Power & fault-tolerant depending on data handling
- Big Data requires the rethinking of program model, system, algorithm, and architecture
- Big Data is a community effort
- Understand the memory system and application data behavior is a must
- Integrated data-access system design is the first step

COMPUTING SCIENCE IS  
STILL FULL OF EXCITEMENT  
STILL FULL OF EXCITEMENT