

# A Compiler-Guided Approach for Reducing Disk Power Consumption by Exploiting Disk Access Locality

S. W. Son, G. Chen, and M. Kandemir  
Dept. of CSE

The Pennsylvania State University

The 4th Annual International Symposium on  
Code Generation and Optimization (CGO-4)  
March 28, 2006

# Outline

- n Motivation
- n Related Work
- n TPM vs. DRPM
- n Our Approach
  - n Single vs. Multi-processor
- n Experimental Evaluation
- n Conclusion

# Motivation

- n High-end cluster/server systems consume a significant amount of power
- n Disk subsystem is one of major contributors to overall power budget

Dell PowerEdge

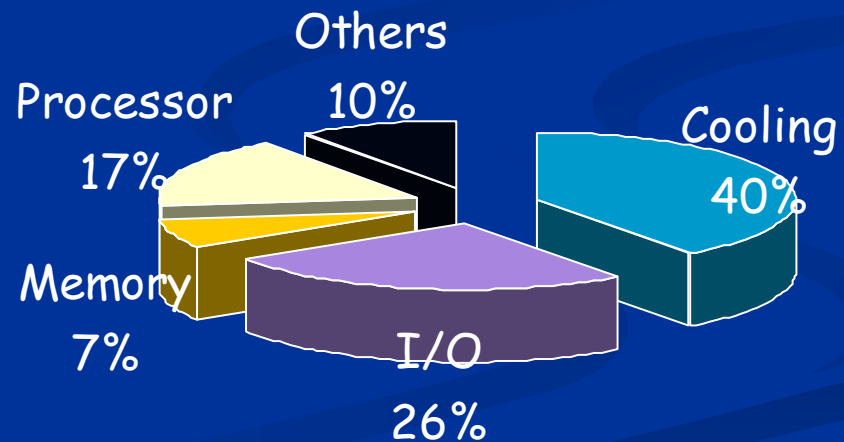


2.1 kW

IBM eServer  
z990



5.3/15.8 kW

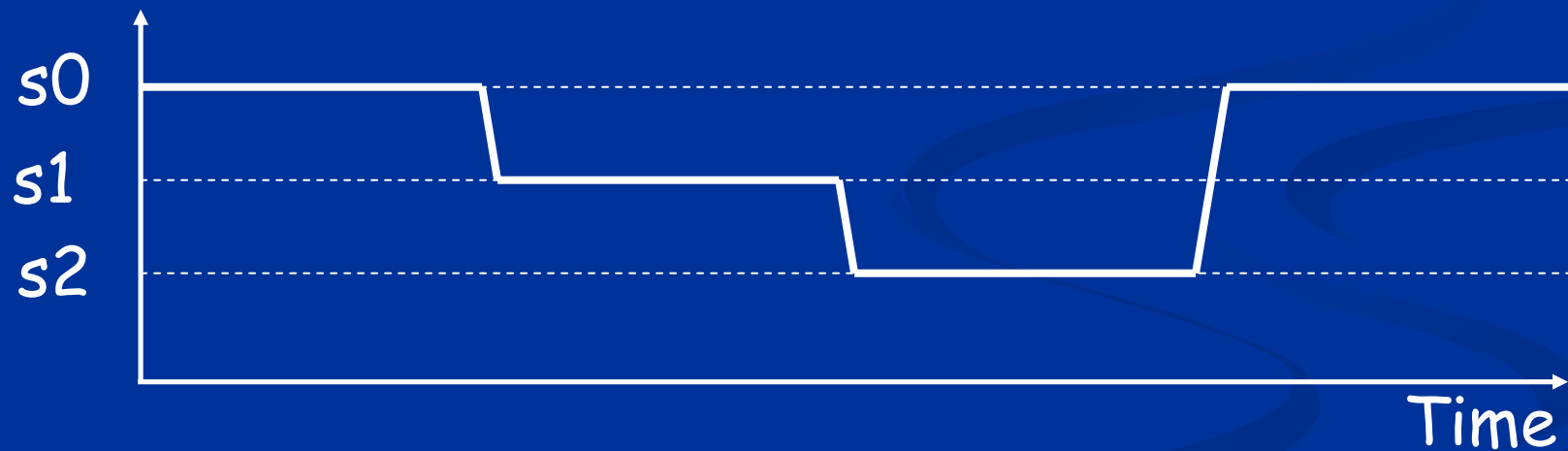
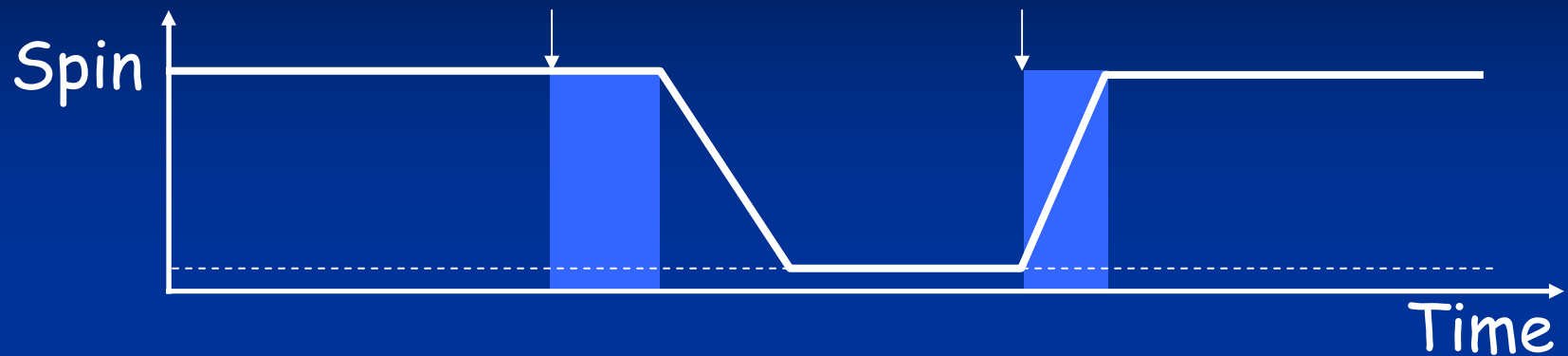


\*Source: Mike Rosenfield, ACEED, February 2003.

# Related Work

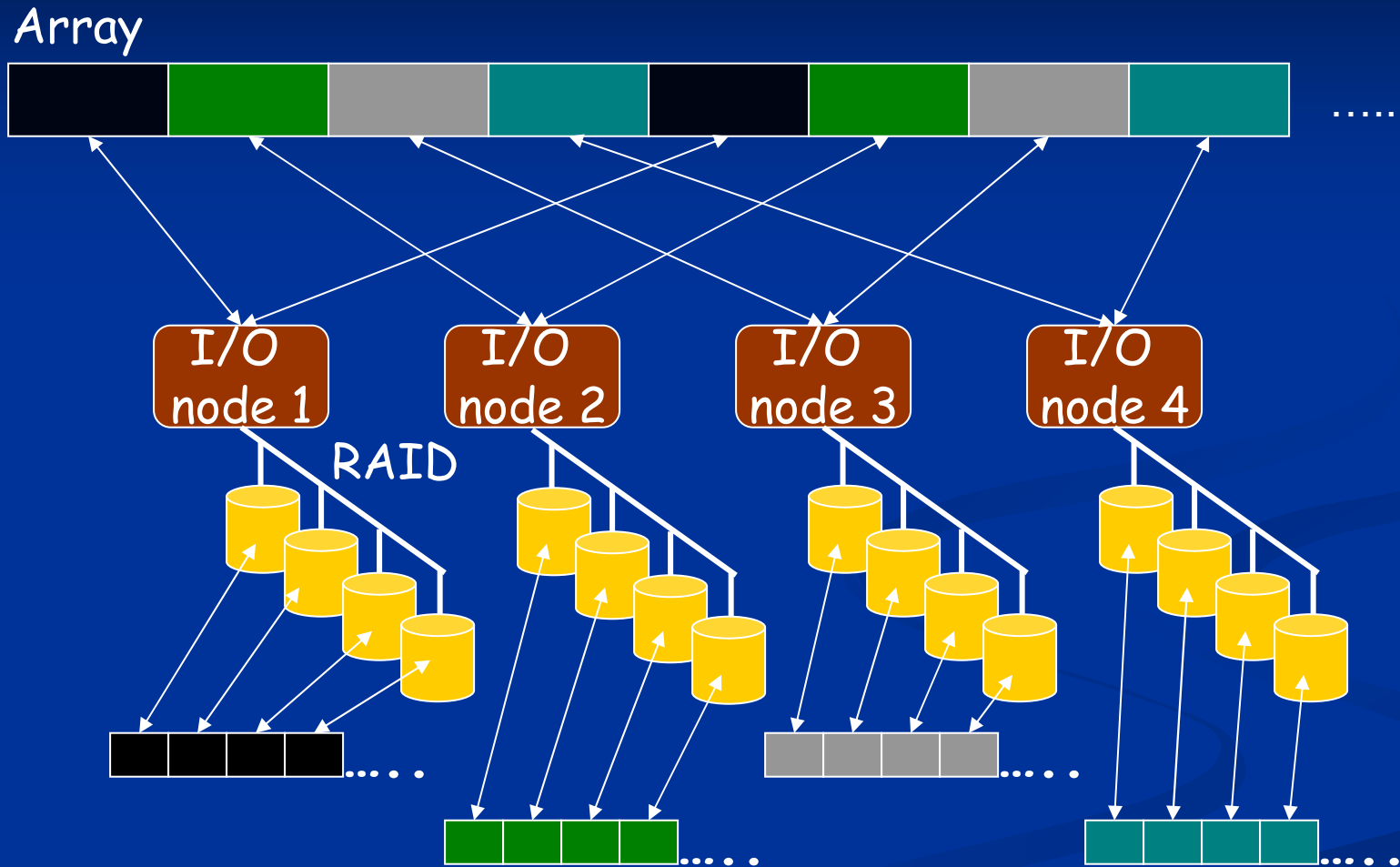
- n Data Locality Optimizations
  - n Target caches and main memories
  - n Lots of prior studies on tiling, shackling, and fusion
- n H/W Based Approaches
  - n Spin up/down (TPM) vs. multi-speed disk (DRPM)
- n OS Based Approaches
  - n PA-LRU, PB-LRU, and PDC
- n Compiler Directed Approaches
  - n Compiler-directed proactive power management, energy-aware disk layout optimization
  - n **Our approach**: compiler-directed code restructuring and generation for increasing disk idle periods

# TPM vs. DRPM



- TPM disk: F. Douglis et al., "Thwarting the Power-Hungary Disk", in USENIX'94
- DRPM disk: S. Gurusurthi et al., "DRPM: Dynamic Speed Control for Power Management in Server Class Disks", in ISCA'03

# Two-Level File Striping

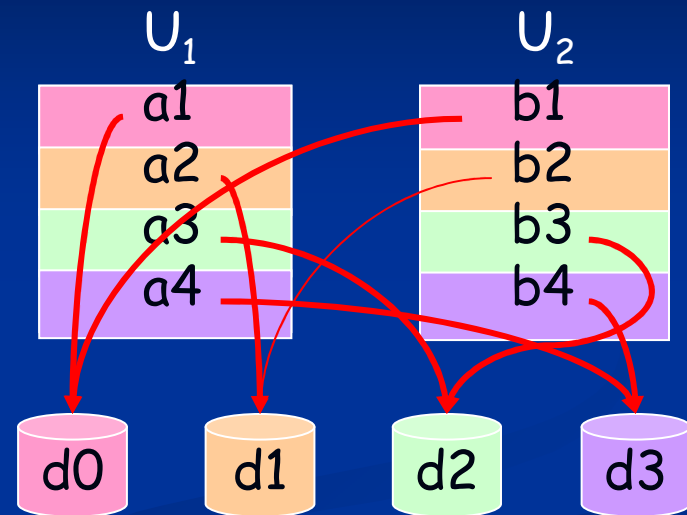


File Striping in Parallel File System, e.g., PVFS, GPFS

# Single Processor Execution

```

L1: for i = 1..N-1
        for j = 1..N-1
            ...U1[i-1][j+2]...
L2: for i = 1..N-1
        for j = 1..N-1
            ...U2[i+1][2j-1]...
L3: for i = 1..N-1
        for j = 1..N-1
            ...U1[2i][j+1]...
    
```



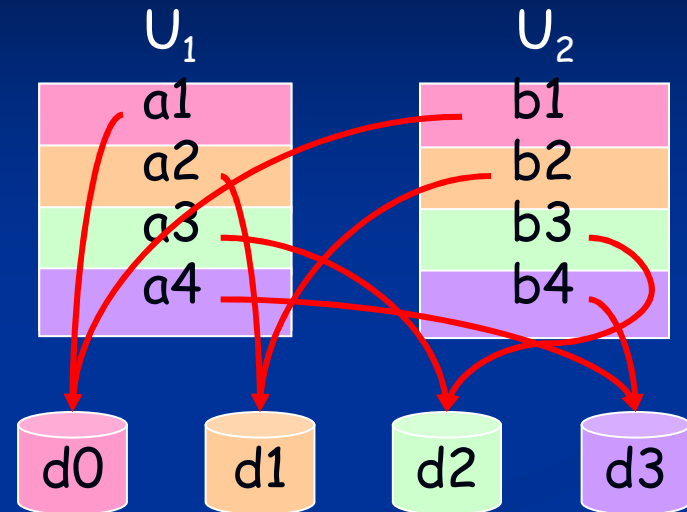
L1	x	x		
L2	x	x	x	
L3	x	x	x	x

•  $U_1$  and  $U_2$  are of same size,  $2N \times 2N$

# Single Processor Execution

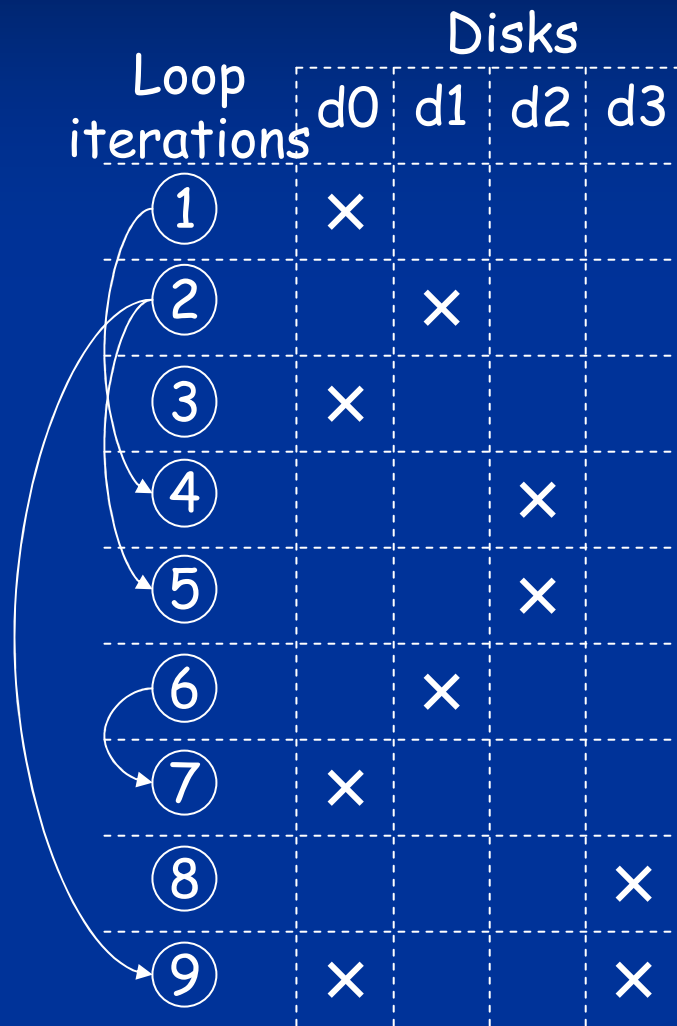
```

for ii = 1..4 {
L1: for i = max(N/2*(ii-1)-1,1)
      ..min(N/2*ii,N-1)
      for j = 1..N-1
        ..U1[i-1][j+2]...
L2: for i = max(N/2*(ii-1)-1,1)
      ..min(N/2*ii-2,N-1)
      for j = 1..N-1
        ..U2[i+1][2j-1]...
L3: for i = max(N/4*(ii-1),1)
      ..min(N/4-1,N-1)
      for j=1..N-1
        ..U1[2i][j+1]...
}
    
```

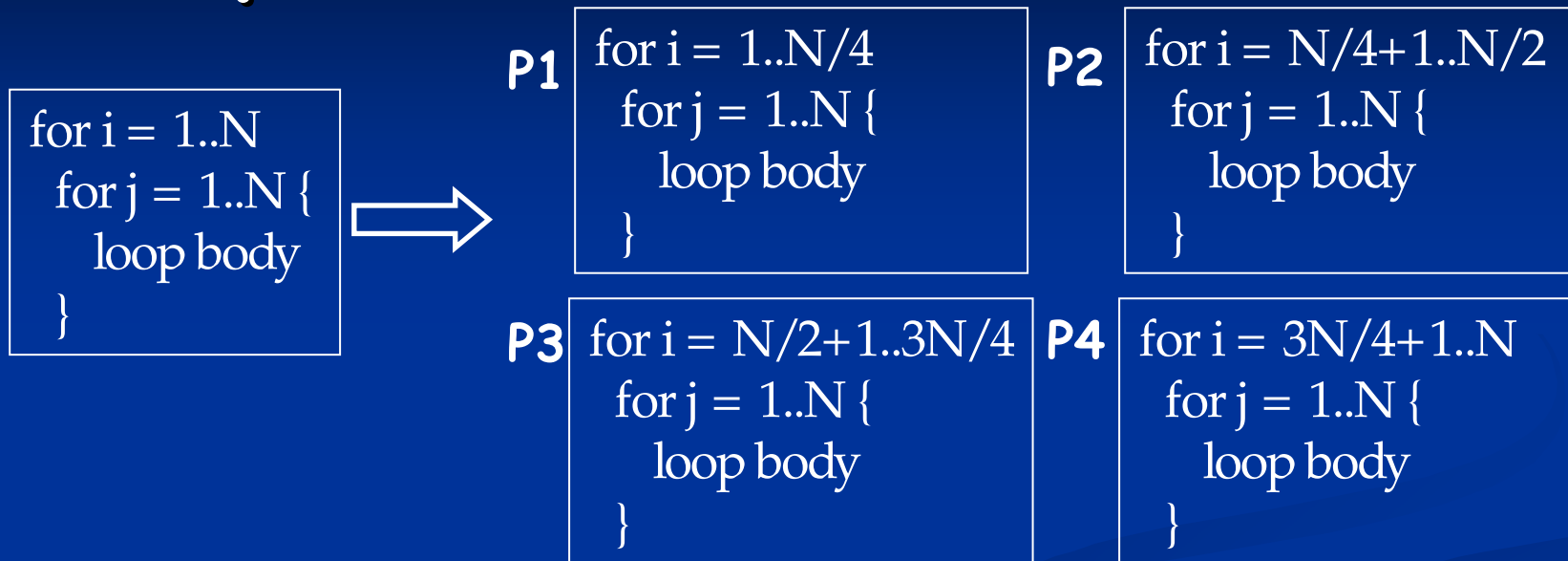


L1	X			
L2	X			
L3	X			
L1		X		
L2		X		
L3		X		
L2			X	
L3			X	
L3				X

# Single Processor Execution with Data Dependences

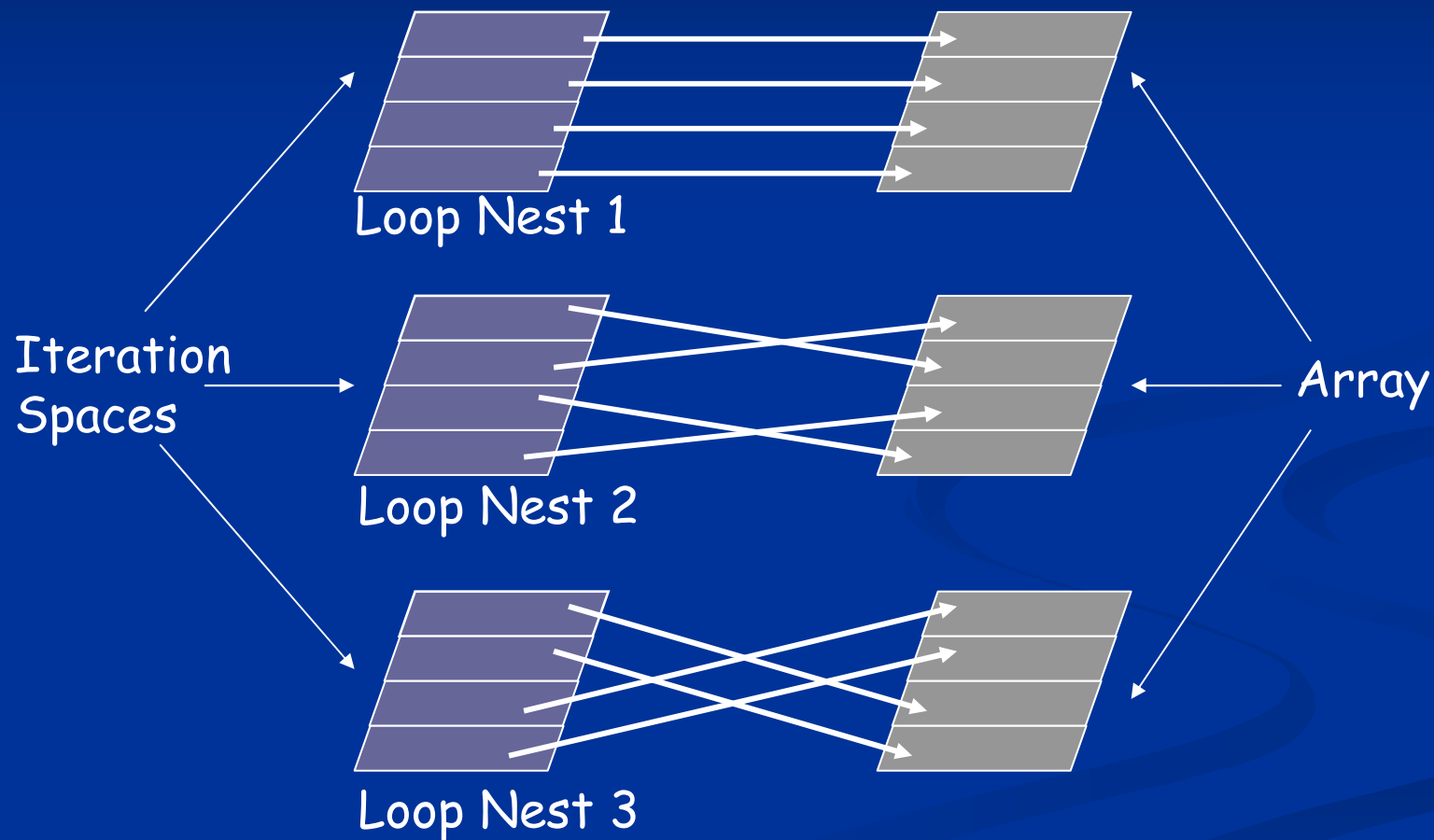


# Loop Based Code Parallelization



- n Coarse grain parallelism: outermost loop
- n Data dependence constraints
- n Parallelize each loop nest individually
  - n Drawback: not capturing data locality between different loop nests

# Multi-Processor Execution - Motivating Example



# Loop-Based vs. Reuse-Aware

LOOP-BASED



Loop Nest 1



Loop Nest 2



Loop Nest 3

REUSE-AWARE



# Reuse-Aware Parallelization

- n Main goal: improving disk locality for each processor – loop iterations accessing the same disk-resident array should be executed by the same processor
- n Methodology: inter-loop-nest disk reuse aware workload assignment
- n Constraint: inter-loop-iteration data dependences

# Mathematical Model

- $n$  Assign loop iterations to each processor based on disk access patterns

For processor  $s$  ( $1 \leq s \leq p$ ) and loop nest  $k$  ( $1 \leq k \leq n$ ).

$Z_{s,j}$ : data elements of array  $Z_j$  assigned to processor  $s$ .

$Q_{s,j,k}$ : the set of loop iterations from loop nest  $k$  accessing  $Z_{s,j}$ .

$\Rightarrow$  Assign  $Q_{s,j,k}$  to processor  $s$ .

# Issues in Reuse-Aware Parallelization

1. Array element partitioning
  - n Disk access patterns, disk reuse, and Parallelism of the program
2. Partial array access
  - n A loop nest accesses a portion of the array
3. Multiple arrays
  - n Multiple arrays accessed by the same processors share the same disk-resident array
  - n Loop iteration mapping must consider all the arrays to improve disk reuse locality

# Issue1: Array Element Partitioning

- n Use loop based parallelization to extract maximum parallelism
- n Given iteration set  $I_s$  executed by processor  $s$ , and the array reference function set  $R_s$ , determine the accessed array region
$$D_s = \{\vec{d} \mid \exists \vec{I} \in I_s, \exists R \in R_s \text{ such that } R(\vec{I}) = \vec{d}\}.$$
- n  $D_s$  can be different for different loop nests
  - n A unification step to obtain a globally acceptable array partitioning

# A Unification Example

P1

```
L1: for i = 1 to N
      for j = 1 to N
          ...U[R1(i, j)]...
```

P2

```
L2: for i = 1 to N
      for j = 1 to N
          ...U[R2(i, j)]...
```

P3

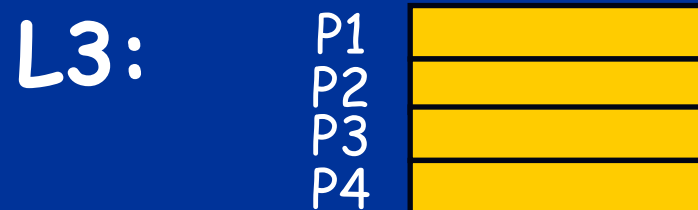
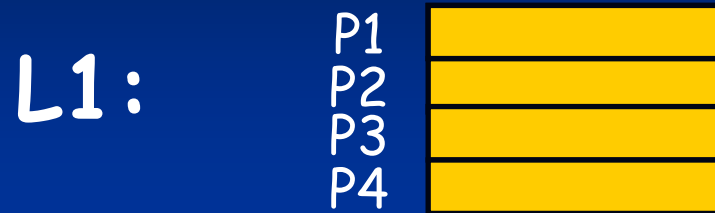
P4

```
L3: for i = 1 to N
      for j = 1 to N
          ...U[R3(i, j)]...
```

U[N][N]

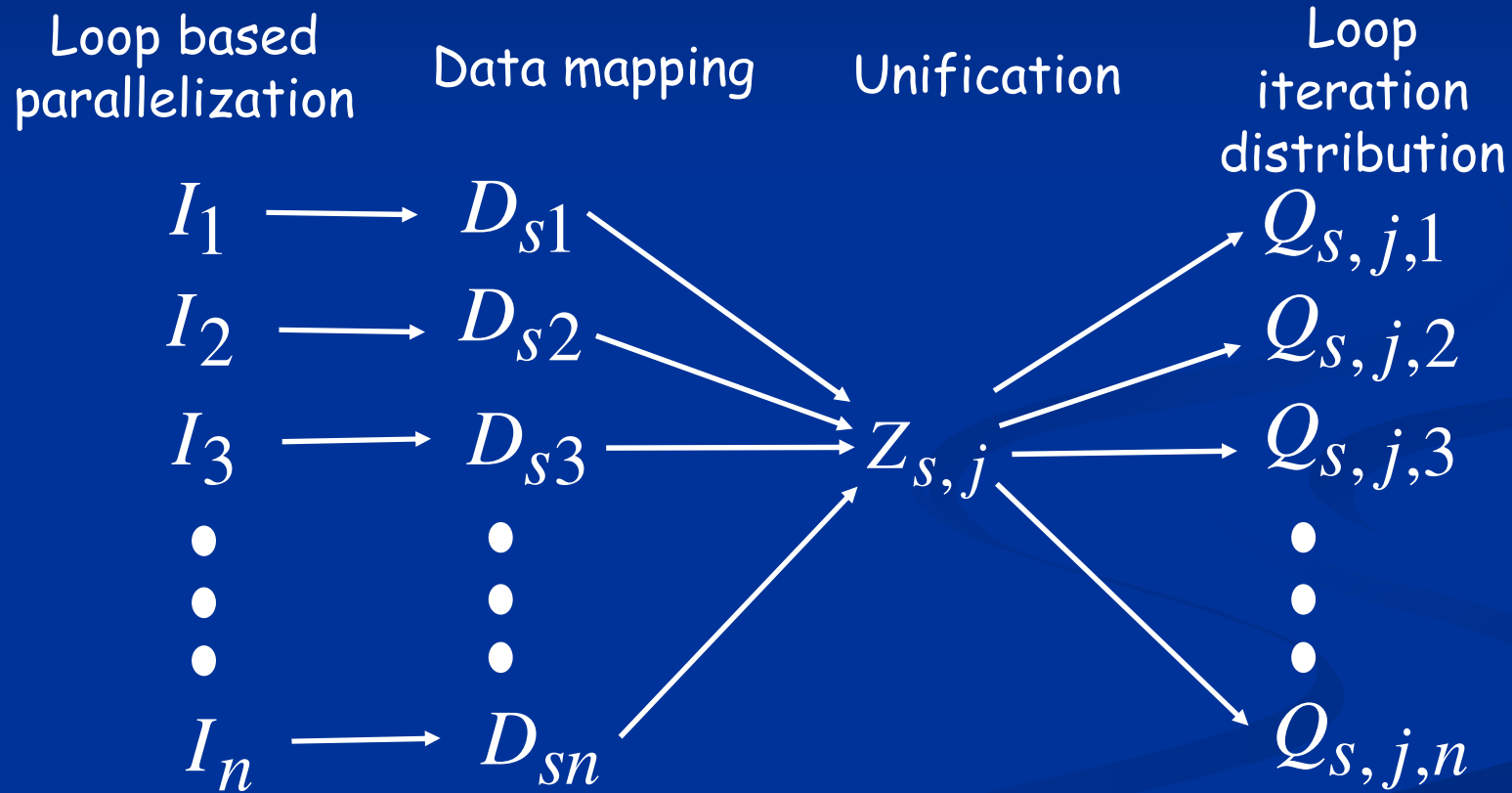


# A Unification Example



Row-block  
array  
partitioning is  
selected

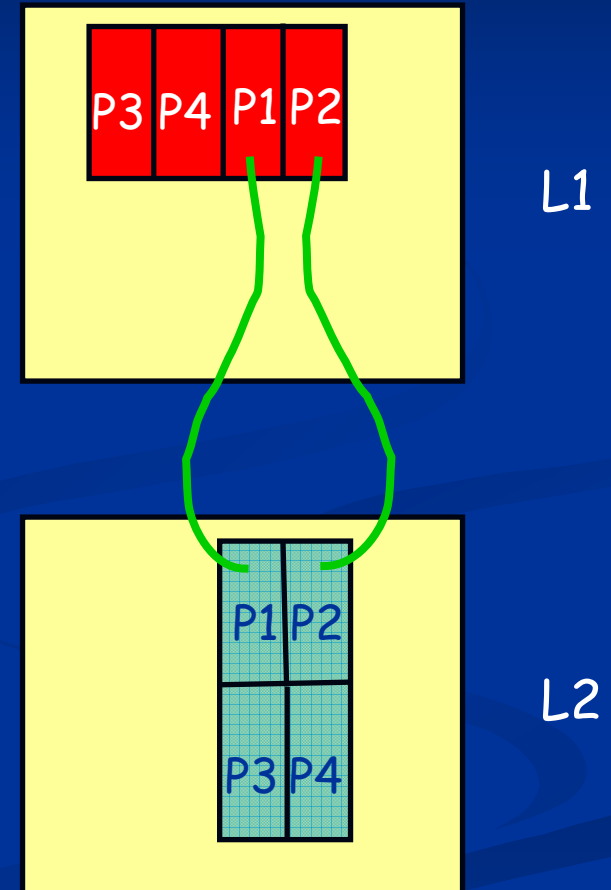
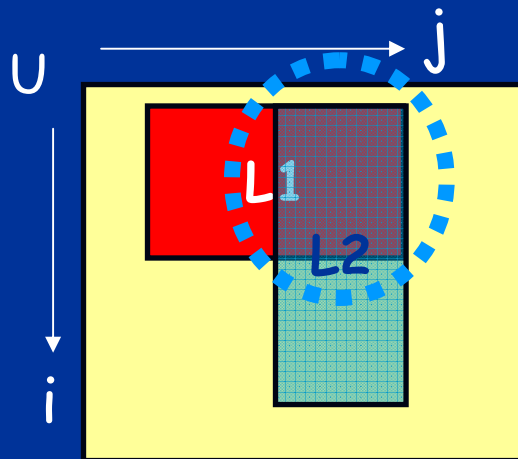
# Data Mapping and Loop Iterations Assignment



# Issue2: Partial Array Access Pattern

**L1:** for  $i = 10$  to  $300$   
for  $j = 100$  to  $900$   
... $U[i][j]$ ...

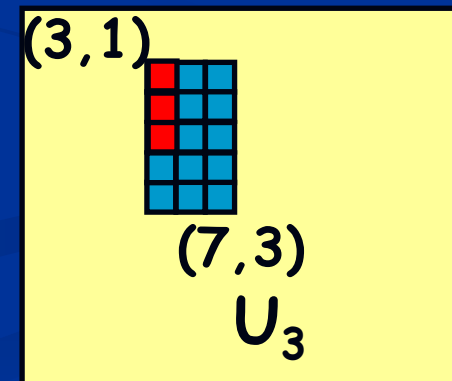
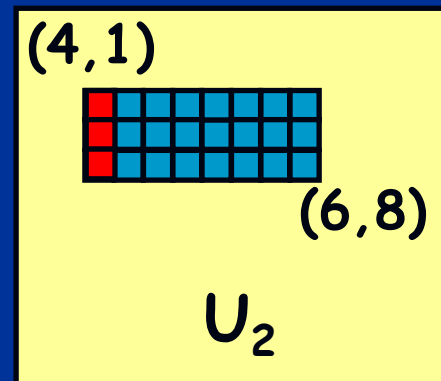
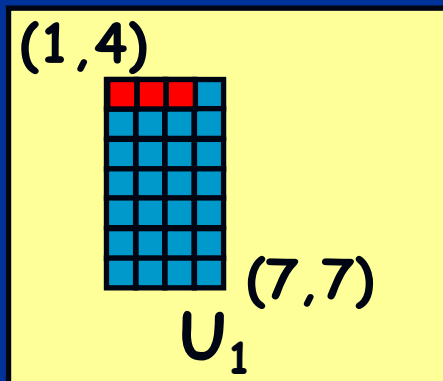
**L2:** for  $i = 10$  to  $600$   
for  $j = 500$  to  $900$   
... $U[i][j]$ ...



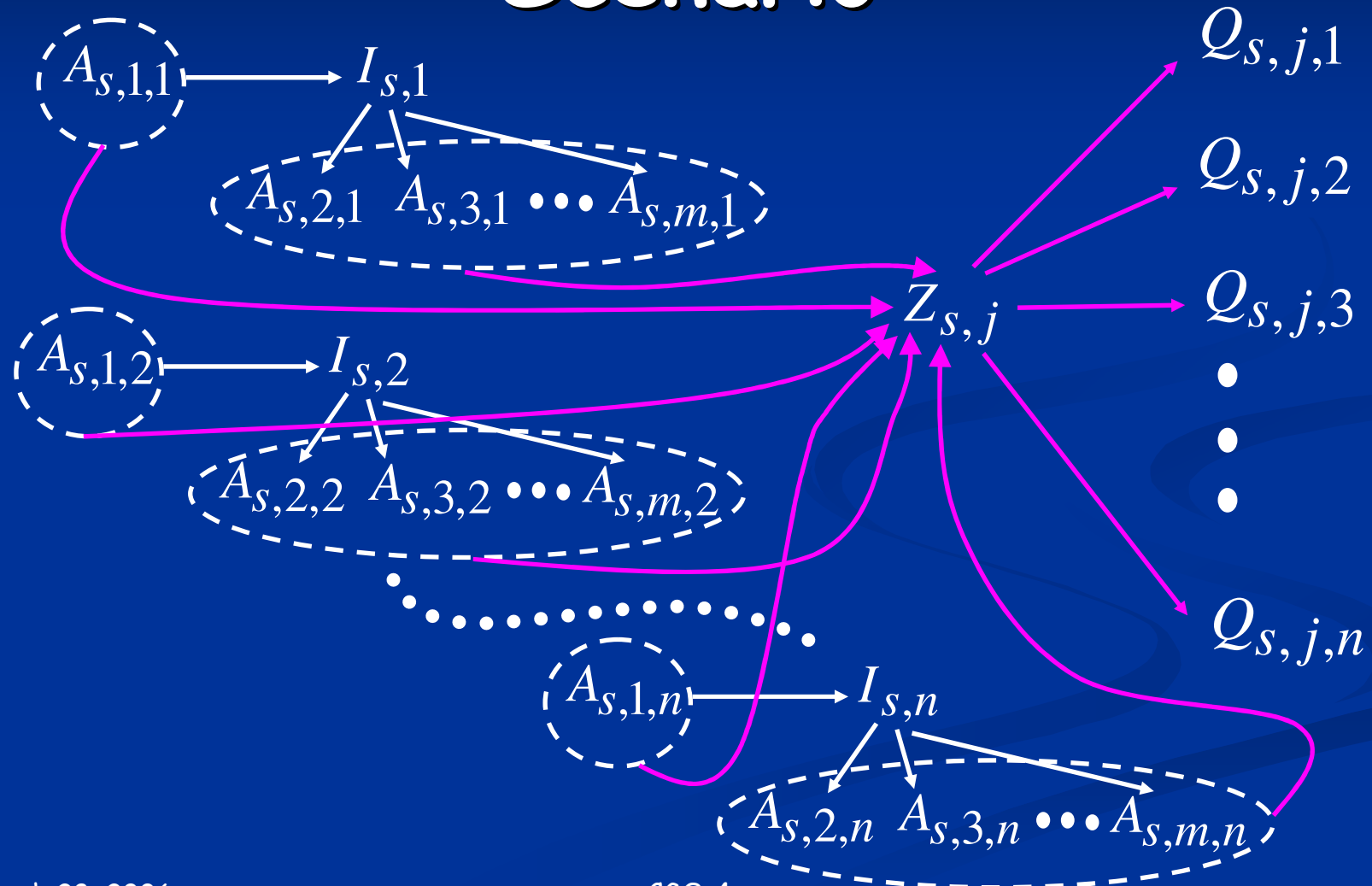
# Issue3: Multiple arrays

- n Two data items accessed by the same loop iteration are said to exhibit *affinity*
- n *Affinity class*: a set of data elements that exhibit affinity

```
for i = 1..M-2
  for j = 4..N
    ..U1[i][j]..U2[j][i]..U3[i+2][j-3]
```



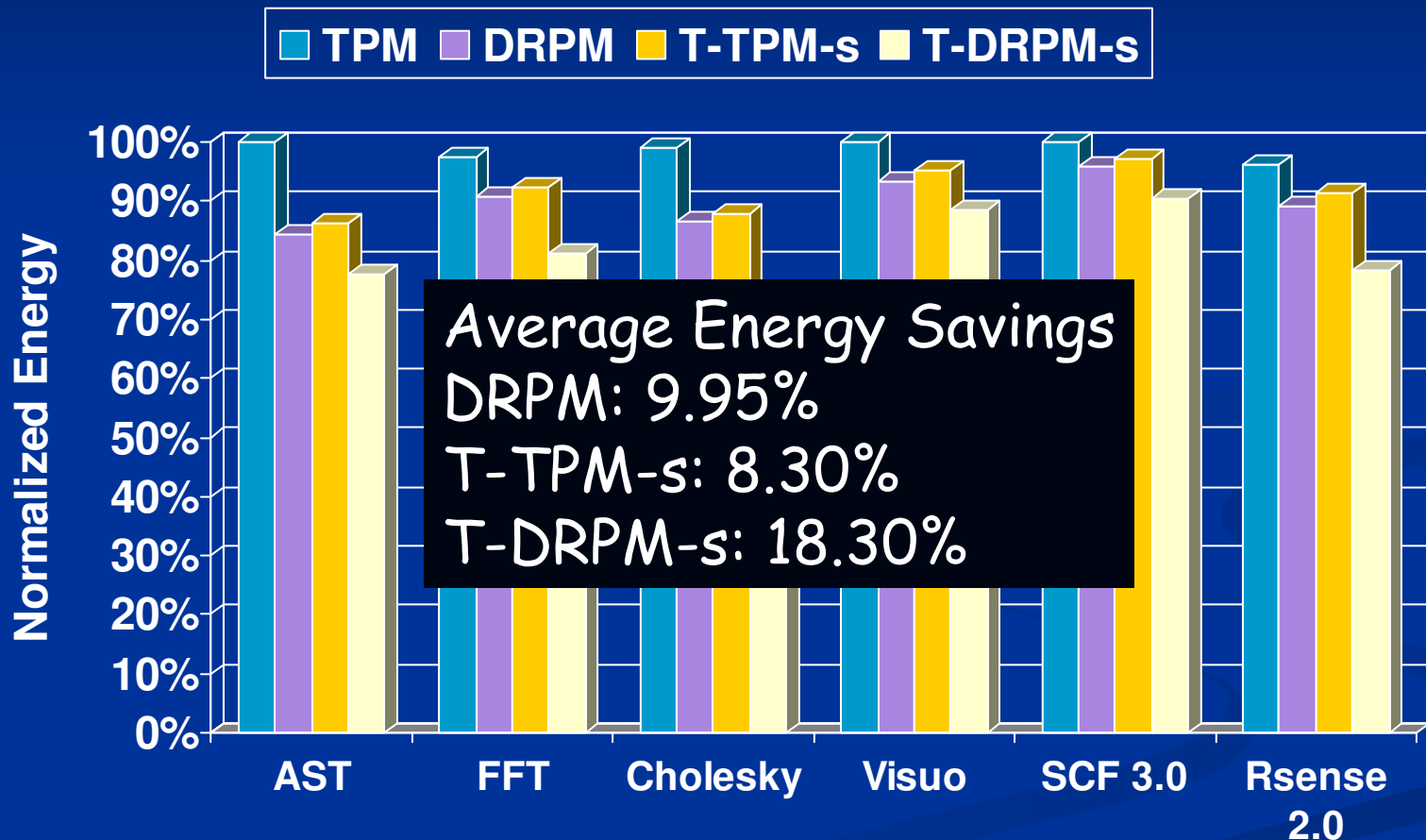
# Determining the Workload of Processor under Multiple Arrays Scenario



# Experimental Platform

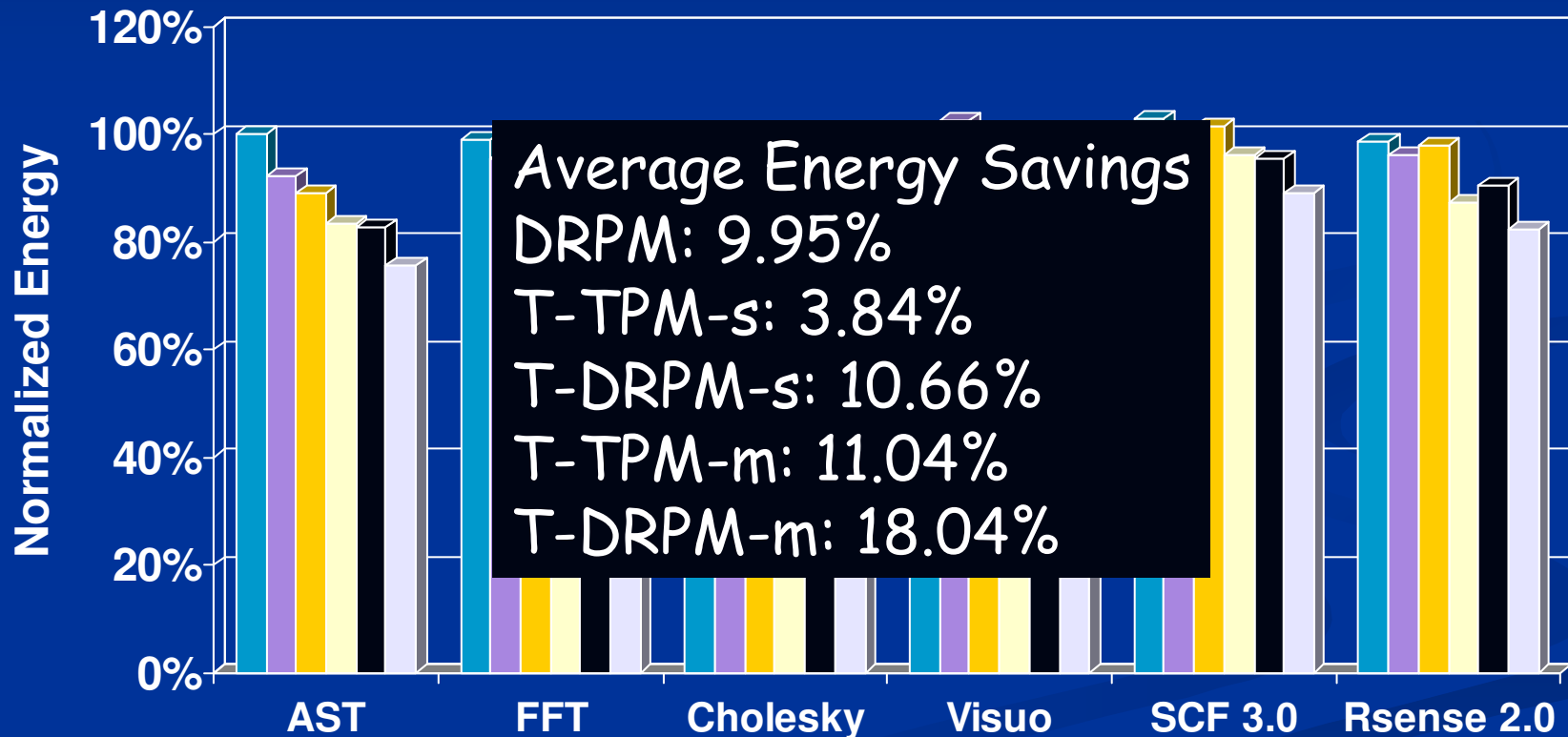
- n Trace generator is used to collect I/O traces
- n Omega Library is used to generate codes after loop iteration assignment
- n Disk energy simulator
  - n Based on IBM36Z15 disk power model
  - n TPM and DRPM energy model
- n Used six I/O intensive codes

# Energy Consumption - single processor execution



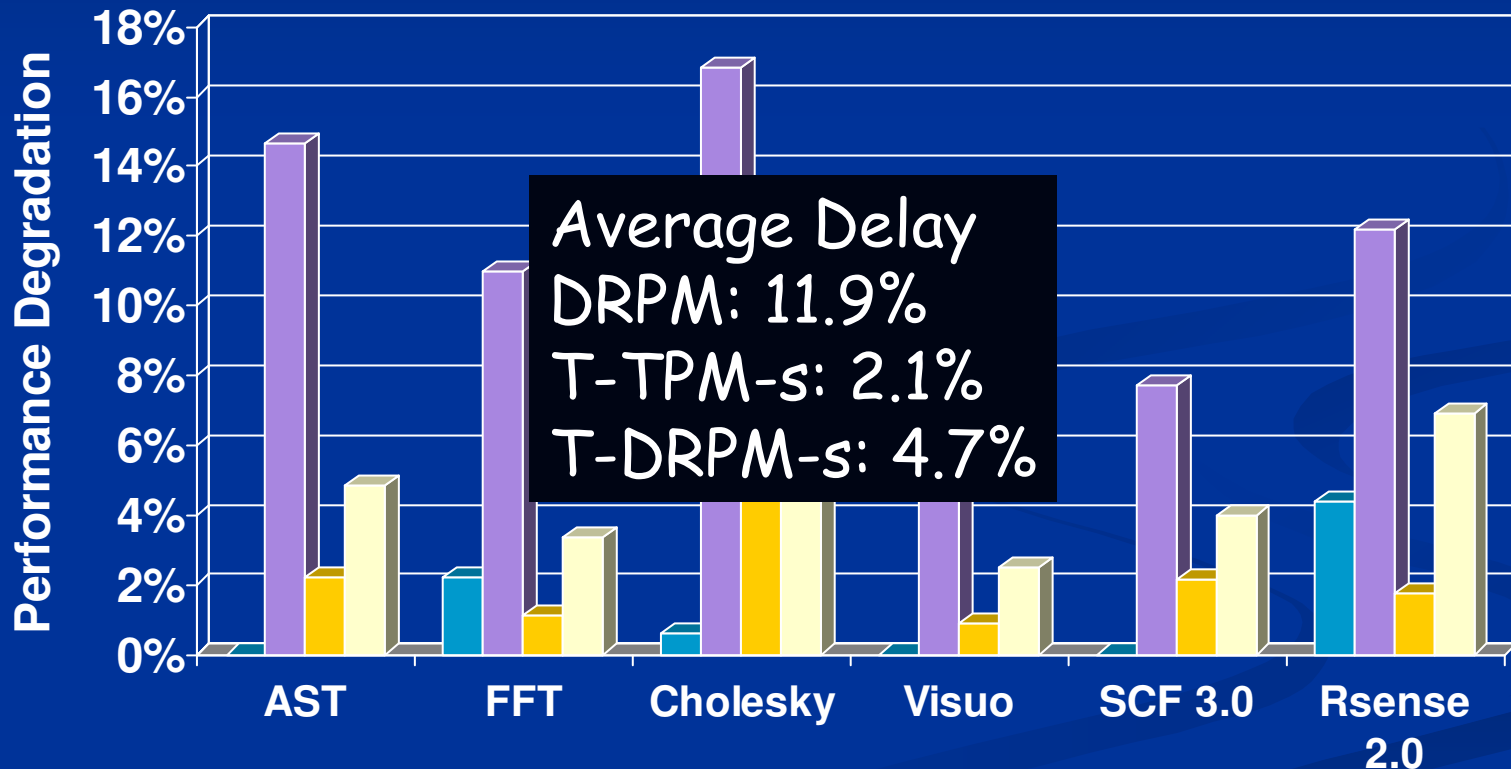
# Energy Consumption - multi-processor execution

■ TPM 
 ■ DRPM 
 ■ T-TPM-s 
 ■ T-DRPM-s 
 ■ T-TPM-m 
 ■ T-DRPM-m



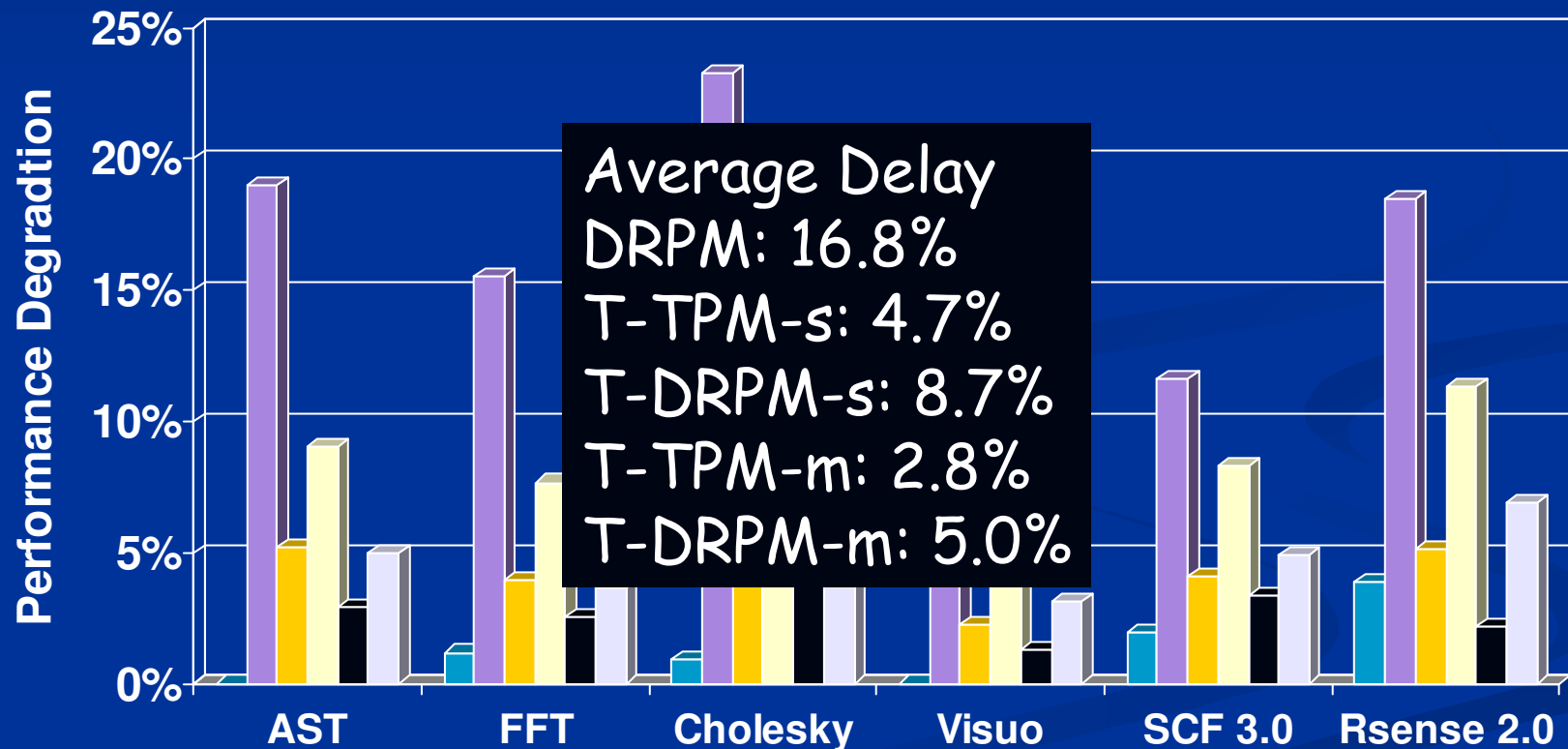
# Performance Degradation - single processor execution

TPM DRPM T-TPM-s T-DRPM-s



# Performance Degradation - multi-processor execution

TPM DRPM T-TPM-s T-DRPM-s T-TPM-m T-DRPM-m



# Conclusion

- n Disk subsystems is one of major contributor to overall power consumption of high-end server systems
- n We proposed a compiler-guided approach to increase the effectiveness of the previously-proposed disk power management schemes
- n Simulation with both single and multi-processor execution shows that our approach achieves more energy savings than hardware schemes

Thank you!

[sson@cse.psu.edu](mailto:sson@cse.psu.edu)