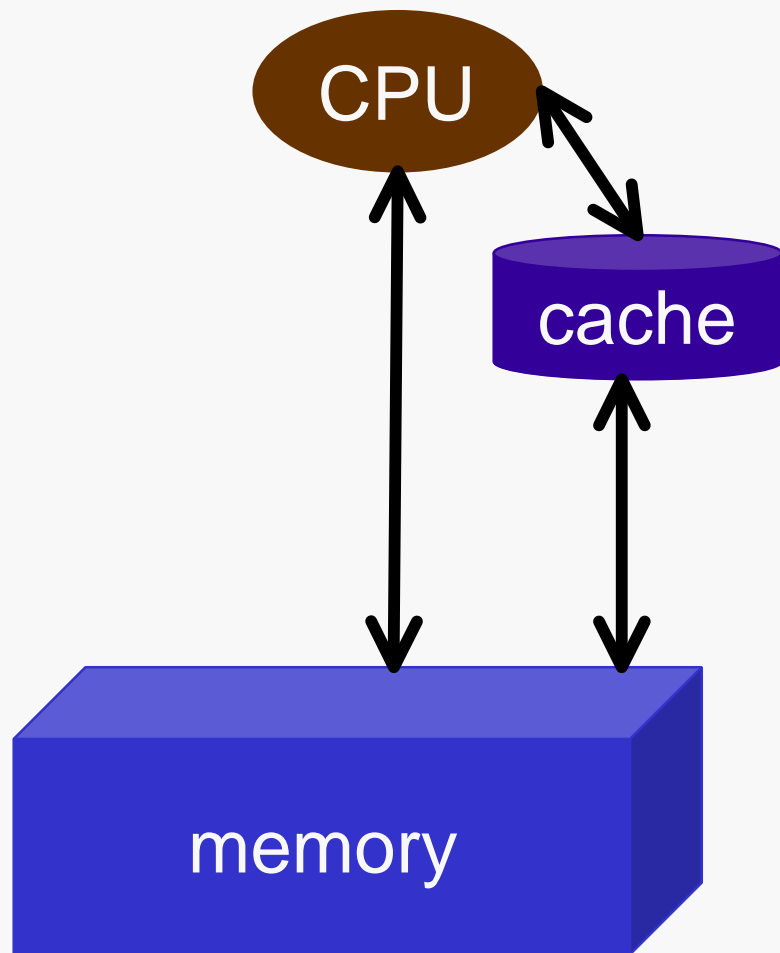


# **Techniques for Improving Global and Dynamic Cache Reuse**

**Chen Ding**

**CS252, Spring 2003**

# Memory Performance



- Problem

- high memory latency

- Improvement 1

- fast cache

- Improvement 2

- Data prefetching

Is there enough bandwidth?

# Bandwidth Bottleneck

- **Hardware trends**
  - CPU speed improved 6400 times in 20 years
  - Memory bandwidth improved 139 times
- **Software trends**
  - large data sets
  - dynamic content and computation
  - modularized programming

# Performance Model

- **Balance**
  - Callahan, Cocke, and Kennedy. JPDC 1988.
- **Machine balance**
  - max words per cycle divided by max flops per cycle
  - # load/store units divided by # floating-point units
- **Program balance**
  - # words accessed divided by # flops executed
  - total loads/stores divided by total floating-point ops
- **Consequences**
  - $MB = PB \rightarrow$  full utilization
  - $MB > PB \rightarrow$  memory idle
  - $MB < PB \rightarrow$  CPU idle

## Extended Model

- **Multi-level memory hierarchy**
  - Ding and Kennedy. IPDPS'00
- **Machine balance**
  - max bandwidth between each pair of consecutive levels divided by max flop
  - max flop, register and cache bandwidth
    - use machine spec
  - memory bandwidth
    - use simple kernels
- **Program balance**
  - # bytes loaded between each pair of consecutive levels for each flop executed
  - use hardware counters to measure # flops, loads/stores, and cache misses

## Balance Table

program/ machine	Program/machine balance		
	L1-Reg	L2-L1	Mem-L2
Convolution	6.4	5.1	5.2
Dmxdy	8.3	8.3	8.4
Mmiki (-O2)	24.0	8.2	5.9
Mmiki (-O3)	8.1	1.0	0.04
FFT	8.3	3.0	2.7
SP	10.8	6.4	4.9
Sweep3D	15.0	9.1	7.8
Origin2000	4.0	4.0	0.8

# Memory-Bandwidth Bottleneck

- Ratios of demand to supply

Applications	Ratio: demand/supply		
	Reg BW	Cache BW	Mem BW
Convolution	1.6	1.3	6.5
Dmxdpy	2.1	2.1	10.5
Mmiki (-O2)	6.0	2.1	7.4
FFT	2.1	0.8	3.4
SP	2.7	1.6	6.1
Sweep3D	3.8	2.3	9.8

- Memory bandwidth is least sufficient
- Maximal CPU utilization: 10% to 33%
- The imbalance is getting worse
- **Software solution: Better caching**

# Problem of Caching

- An example
  - 7 accesses
  - single-element cache
- Cache management
  - LRU: 1 reuse
  - Belady: 2 reuses
- Our approach
  - fuse computation on the same data: 4 reuses
  - group data used by the same computation

adbaabd

adbaabd

aaaddbb

Memory





## Computation Fusion

```
Function Initialize:  
  read input[1:N].data1  
  read input[1:N].data2  
End Initialize
```

```
Function Process:  
  //Step_A  
  A_tmp[1:N].data1  
    ← input[1:N].data1  
  
  //Step_B  
  B_tmp[1:N].data1  
    ← input[1:N].data2  
  ...  
End Process
```

```
//Fused_step_1  
for each i in [1:N]  
  read input[i].data1  
  A_tmp[i].data1  
    ← input[i].data1  
end for
```

```
//Fused_step_2  
for each i in [1:N]  
  read input[i].data2  
  B_tmp[i].data1  
    ← input[i].data2  
end for  
...
```

## Data Grouping

```
//Fused_step_1
for each i in [1:N]
  read input[i].data1
  A_tmp[i].data1
  ← input[i].data1
end for
```

```
//Fused_step_2
for each i in [1:N]
  read input[i].data2
  B_tmp[i].data1
  ← input[i].data2
end for
...
```

```
//Fused_step_1
for each i in [1:N]
  read group1[i].data1
  group1[i].data2
  ← group1[i].data1
end for
```

```
//Fused_step_2
for each i in [1:N]
  read group2[i].data1
  group2[i].data2
  ← group2[i].data1
end for
...
```

## Original

```
Function Initialize:  
  read input[1:N].data1  
  read input[1:N].data2  
End Initialize
```

```
Function Process:  
  A_tmp[1:N].data1  
    ← input[1:N].data1  
  
  B_tmp[1:N].data1  
    ← input[1:N].data2  
  
  ...  
End Process
```

## Transformed

```
for each i in [1:N]  
  read group1[i].data1  
  group1[i].data2  
    ← group1[i].data1  
end for  
  
for each i in [1:N]  
  read group2[i].data1  
  group2[i].data2  
    ← group2[i].data1  
end for  
  
...
```

- Computation fusion recombines all functions
- Data grouping reshuffles all data

## But How?

- **Programmers**
  - loss of modularity
  - data layout depends on function
- **Hardware/operating system**
  - limited scope
  - run-time overhead
- **Compilers**
  - global scope
  - off-line analysis/transformation
  - imprecise information

## Outline

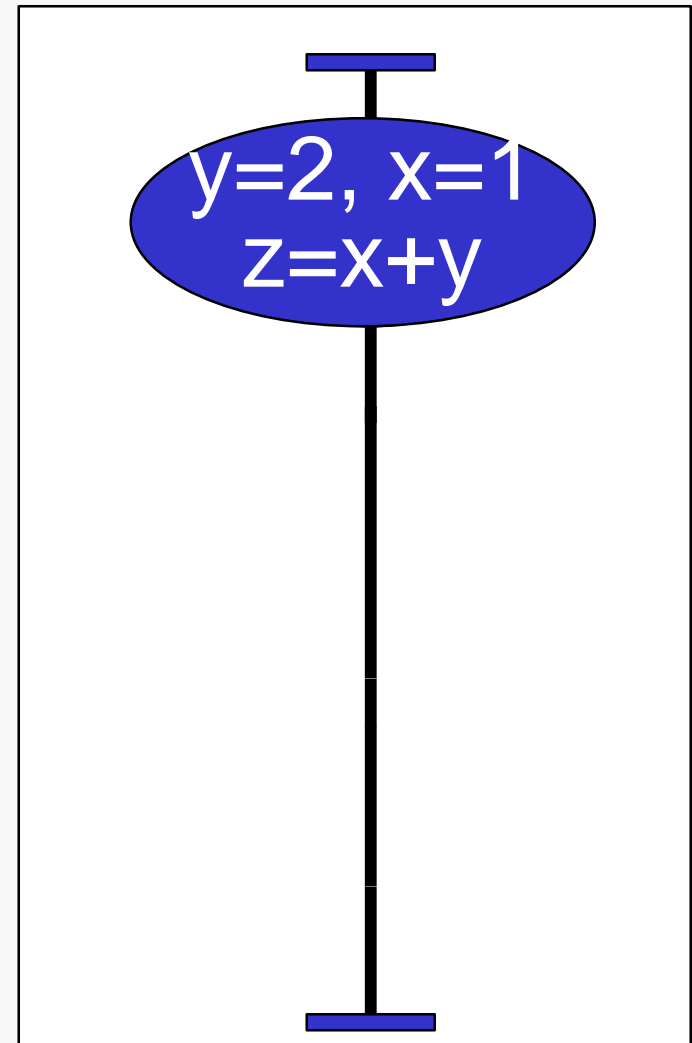
- **New compiler transformations**
  - global computation fusion
  - global data grouping
  - run-time adjustments
- **Evaluation**
  - regular programs
  - dynamic programs
- **Summary**

# **Global Computation Fusion**

**[IPDPS'01 Best Paper  
with Ken Kennedy]**

## Overall Fusion Process

for each statement in program  
  find its data sharing predecessor  
  try clustering them (fusion)  
  if succeed  
    *apply fusion recursively*  
  end if  
end for

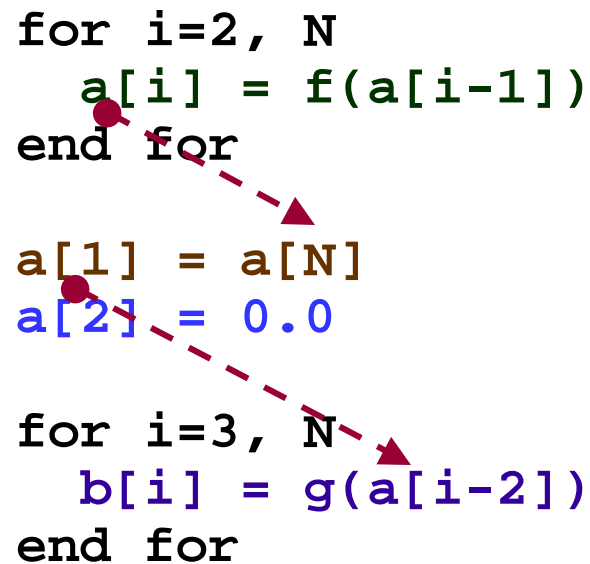


## Example Fusion

```
for i=2, N
  a[i] = f(a[i-1])
end for

a[1] = a[N]
a[2] = 0.0

for i=3, N
  b[i] = g(a[i-2])
end for
```



### Difficulties

- incompatible shapes
- data dependence

### Three cases of fusion

- between iteration & loop
  - embedding
- between iterations
  - interleaving + alignment
- otherwise
  - iteration reordering,
  - e.g. loop splitting



## Example Fusion

```
for i=2, N
  a[i] = f(a[i-1])
end for
```

```
a[1] = a[N]
a[2] = 0.0
```

```
for i=3, N
  b[i] = g(a[i-2])
end for
```

```
for i=2, N
  a[i]=f(a[i-1])
  if (i==3)
    a[2]=0.0
  else if (i==N)
    a[1]= a[N]
  end if
end for
```

```
for i=3, N
  b[i] = g(a[i-2])
end for
```

- loop embedding

## Example Fusion

```
for i=2, N
  a[i] = f(a[i-1])
end for
```

```
a[1] = a[N]
a[2] = 0.0
```

```
for i=3, N
  b[i] = g(a[i-2])
end for
```

```
for i=2, N
  a[i]=f(a[i-1])
  if (i==3)
    a[2]=0.0
  else if (i==N)
    a[1]= a[N]
  end if
end for
```

```
for i=4, N
  b[i] = g(a[i-2])
end for
```

```
b[3] = g(a[1])
```

- loop embedding, loop splitting,

## Example Fusion

```
for i=2, N
  a[i] = f(a[i-1])
end for
```

```
a[1] = a[N]
a[2] = 0.0
```

```
for i=3, N
  b[i] = g(a[i-2])
end for
```

```
for i=2, N
  a[i]=f(a[i-1])
  if (i==3)
    a[2]=0.0
  else if (i==N)
    a[1]= a[N]
  end if

  if (i>2 && i<N)
    b[i+1] = g(a[i-1])
  end if
end for
```

```
b[3] = g(a[1])
```

- loop embedding, loop splitting, interleaving+alignment

## More on Fusion

- **Multi-level fusion**
  - gives priority to fusion at outer levels
- **Optimal fusion**
  - minimal data sharing
  - an NP-hard problem
  - hyper-graph formulation of data sharing
  - a polynomial-time case

## Is Fusion Enough?

```
//Fused_step_1
for each i in [1:N]
  read input[i].data1
  A_tmp[i].data1
  ← input[i].data1
end for
```

```
//Fused_step_2
for each i in [1:N]
  read input[i].data2
  B_tmp[i].data1
  ← input[i].data2
end for
...
```

```
//Fused_step_1
for each i in [1:N]
  read group1[i].data1
  group1[i].data2
  ← group1[i].data1
end for
```

```
//Fused_step_2
for each i in [1:N]
  read group2[i].data1
  group2[i].data2
  ← group2[i].data1
end for
...
```

# **Data Regrouping**

**[LCPC'99, with Ken Kennedy]**

# Program Analysis

- **Computation phases**
  - amount of data accessed
- **Compatible arrays**
  - array size
  - order of access
- **Arrays are split into the smallest unit possible**
  - $A(3,N) \rightarrow A1(N), A2(N), A3(N)$
- **Regrouping is applied to individual set of compatible arrays**

## Example

<i>Computation phases</i>	<i>Arrays accessed</i>
<i>Constructing neighbor list</i>	<code>position</code>
<i>Smoothing attributes</i>	<code>position, velocity, heat, derivative, viscosity</code>
<i>Hydrodynamics 1</i>	<code>density, momentum</code>
<i>Hydrodynamics 2</i>	<code>momentum, volume, energy, cumulative totals</code>
<i>Stress 1</i>	<code>volume, energy, strength, cumulative totals</code>
<i>Stress 2</i>	<code>density, strength</code>



# Regrouping Algorithm

- **Requirements**
  - 1. regroup as many arrays as possible
  - 2. do not introduce useless data
- **Solution**
  - group two arrays if and only if they are always accessed together
- **Properties**
  - time complexity is  $O(N_{\text{arrays}} * N_{\text{phases}})$
  - compile-time optimal
  - minimal page-table working set
  - eliminate all useless data in cache

## Example

<i>Computation phases</i>	<i>Arrays accessed</i>
<i>Constructing neighbor list</i>	position
<i>Smoothing attributes</i>	position, velocity, heat, derivative, viscosity
<i>Hydrodynamics 1</i>	density, momentum
<i>Hydrodynamics 2</i>	momentum, volume, energy, cumulative totals
<i>Stress 1</i>	volume, energy, strength, cumulative totals
<i>Stress 2</i>	density, strength

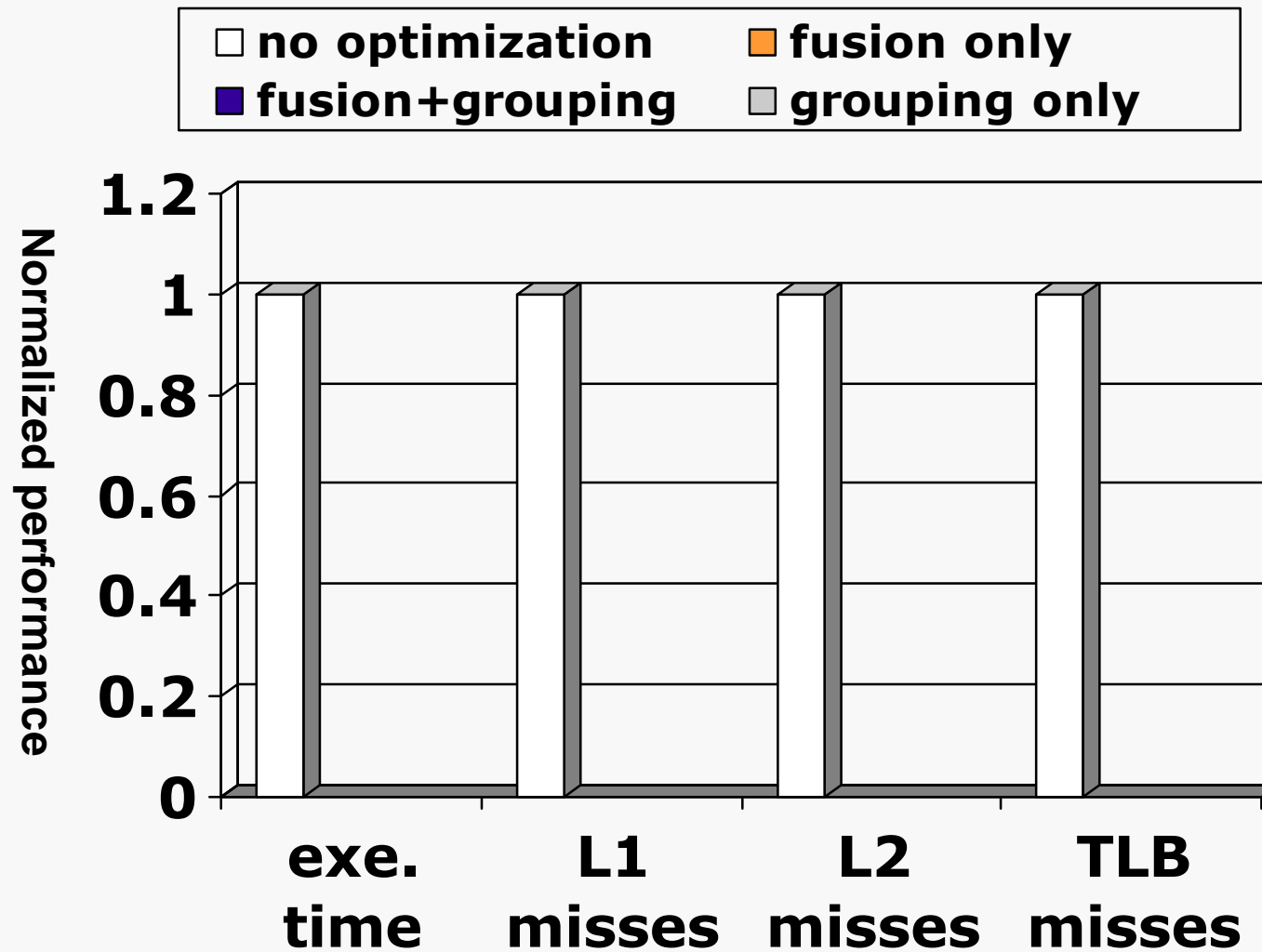
## More on Regrouping

- **Multi-level regrouping**
  - grouping array segments as well as array elements
  - consistent regrouping
- **Extensions**
  - allowing useless data
  - allowing dynamic remapping
  - NP-complete proofs
  - architecture-dependent data layout

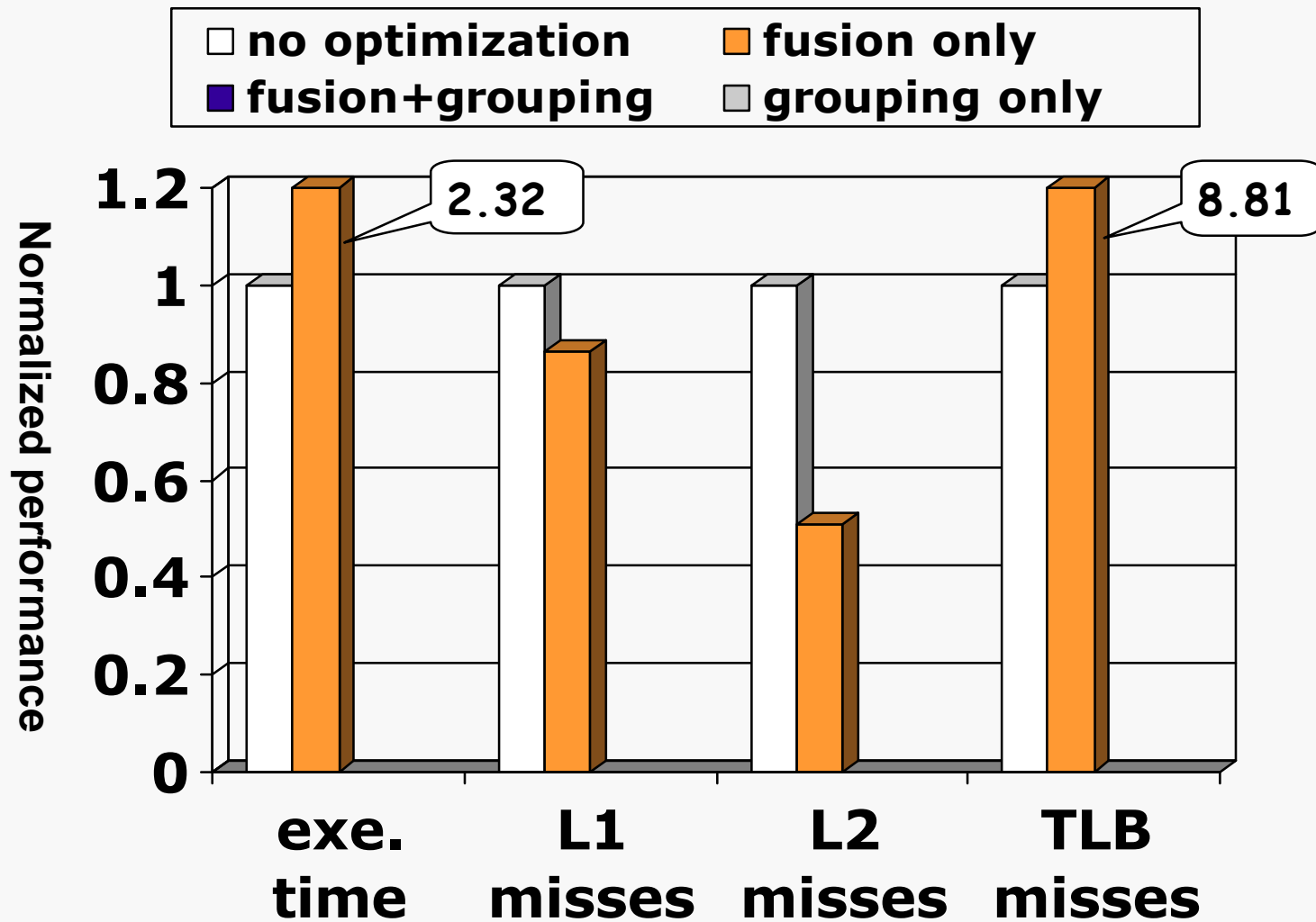
## ***NAS/SP***

- **Benchmark application from NASA**
  - computational fluid dynamics (CFD)
  - class B input, 102x102x102
  - 218 loops in 67 loop nests, distributed into 482 loops
  - 15 global arrays, split into 42 arrays
- **Optimizations**
  - fused into 8 loop nests
  - grouped into 17 new arrays, e.g.
    - {ainv[n,n,n], us[n,n,n], qs[n,n,n], u[n,n,n,1-5]}
    - {lhs[n,n,n,6-8], lhs[n,n,n,11-13]}

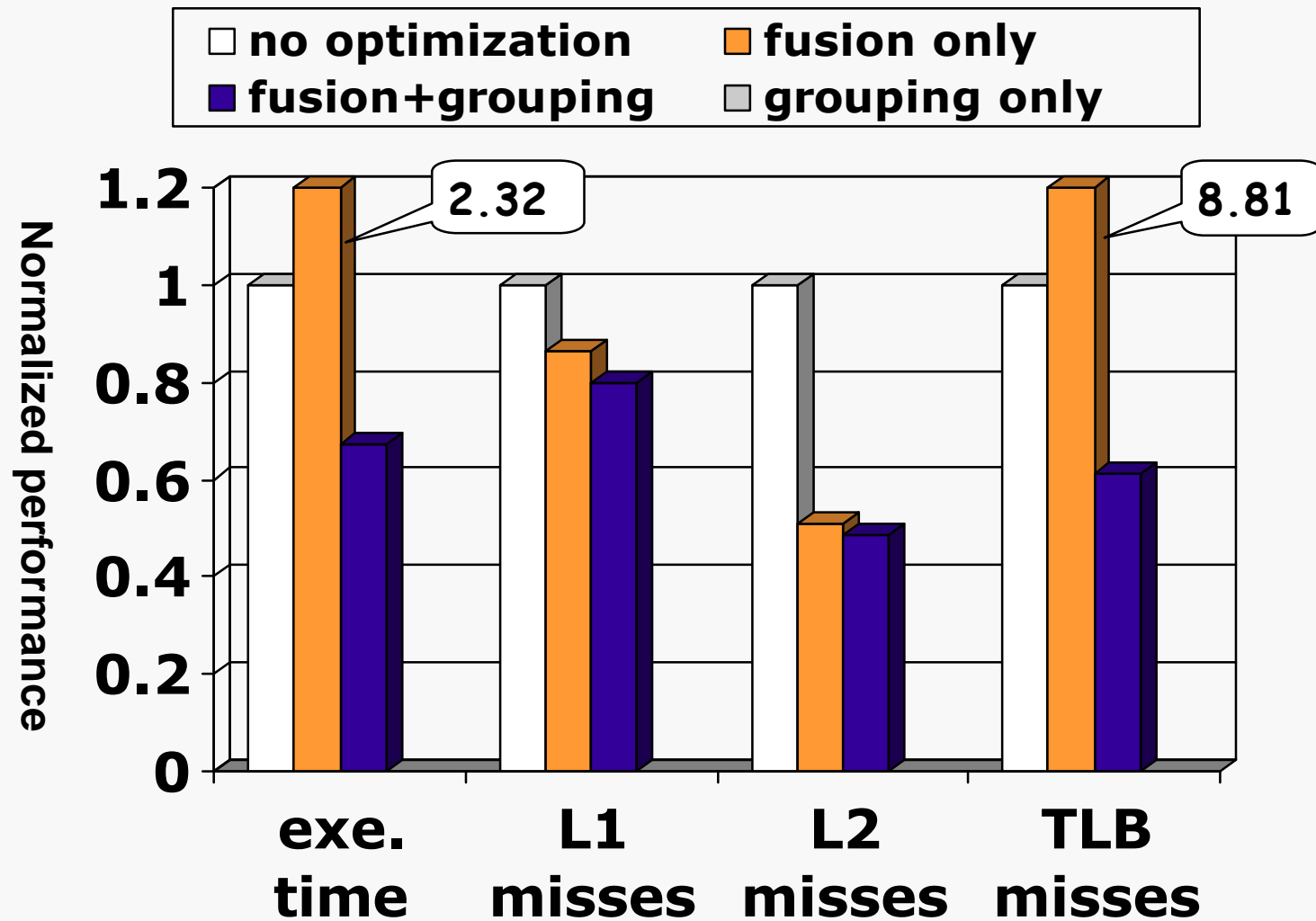
## NAS/SP



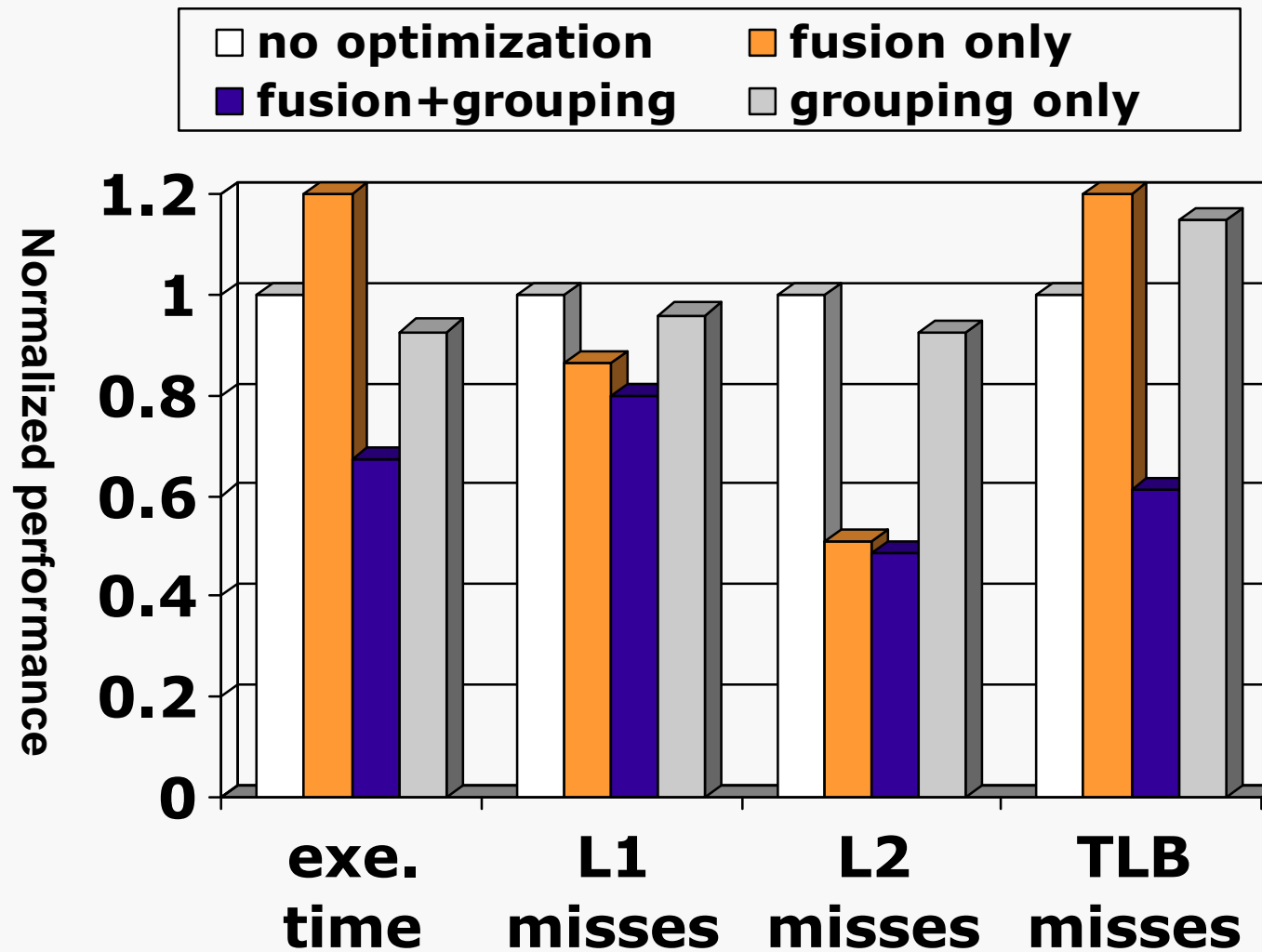
## NAS/SP



## NAS/SP



## NAS/SP





# **Dynamic Optimizations**

## **[PLDI'99, with Ken Kennedy]**

## Unknown Access

“Every problem can be solved by adding one more level of indirection.”

- Irregular and dynamic applications
  - Irregular data structures are unknown until run time
  - Data and their uses may change during the computation
- For example
  - Molecular dynamics
  - Sparse matrix
- Problems
  - How to optimize at run time?
  - How to automate?

## Example packing

original  
array

f[1]

f[2]

f[3]

...

data  
access

f[8], f[800], f[8], f[2], ...

transformed  
array

f[8]

f[800]

f[2]

...

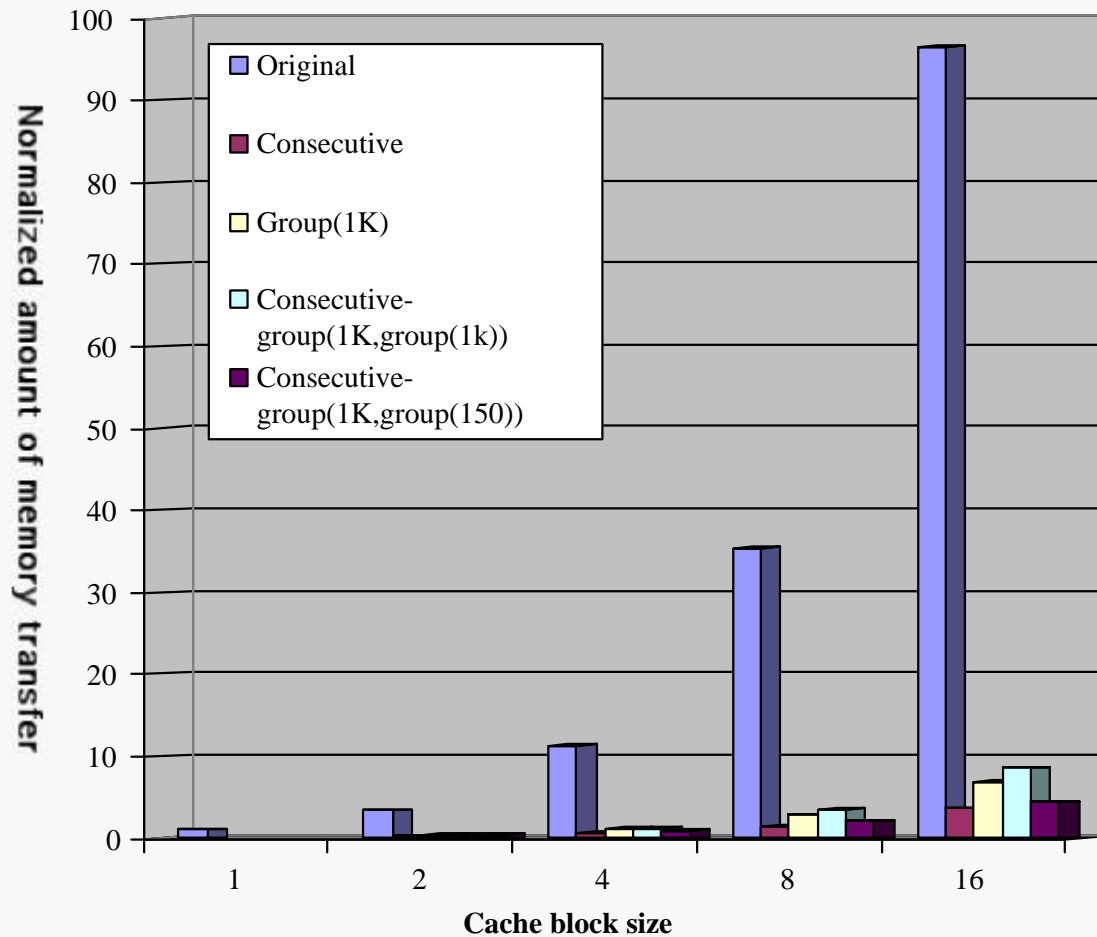
Software remapping:

$f[t[i]] \rightarrow f[\text{remap}[t[i]]] \rightarrow f[t'[i]]$

$f[i] \rightarrow f[\text{remap}[i]] \rightarrow f[i]$

# Example of Run-Time Optimization

Moldyn, 8K elements, 4K cache



## Dynamic Optimizations

- **Locality grouping & Dynamic packing**
  - run-time versions of computation fusion & data grouping
  - linear time and space cost
- **Compiler support**
  - analyze data indirections
  - find all optimization candidates
  - use run-time maps to guarantee correctness
  - remove unnecessary remappings
    - pointer update
    - array alignment
- **The first set of compiler-generated run-time transformations**

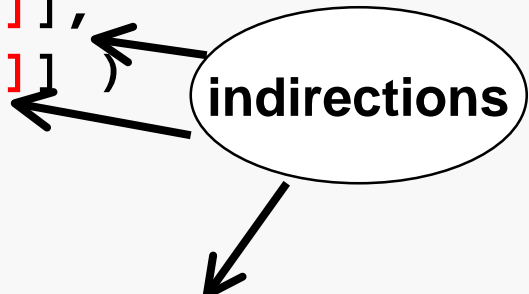
```
packing Directive: apply packing using interactions
```

```
for each pair (i,j) in interactions  
    compute_force( force[i], force[j] )  
end for
```

```
for each object i  
    update_location( location[i], force[i] )  
end for
```

```
apply_packing(interactions[*],force[*],inter_map[*])
for each pair (i,j) in interactions
  compute_force( force[inter_map[i]],
                 force[inter_map[j]] )
end for

for each object i
  update_location(location[i],force[inter_map[i]])
end for
```



The diagram illustrates the role of the 'indirections' variable. It is represented by an oval labeled 'indirections'. Three arrows originate from this oval: one points to the closing parenthesis of the 'compute\_force' function call, another points to the 'force[inter\_map[j]]' argument, and a third points to the 'force[inter\_map[i]]' argument in the 'update\_location' function call. This indicates that 'indirections' is used to store the index of the object being interacted with, which is then used to look up the force vector in the 'force' array.

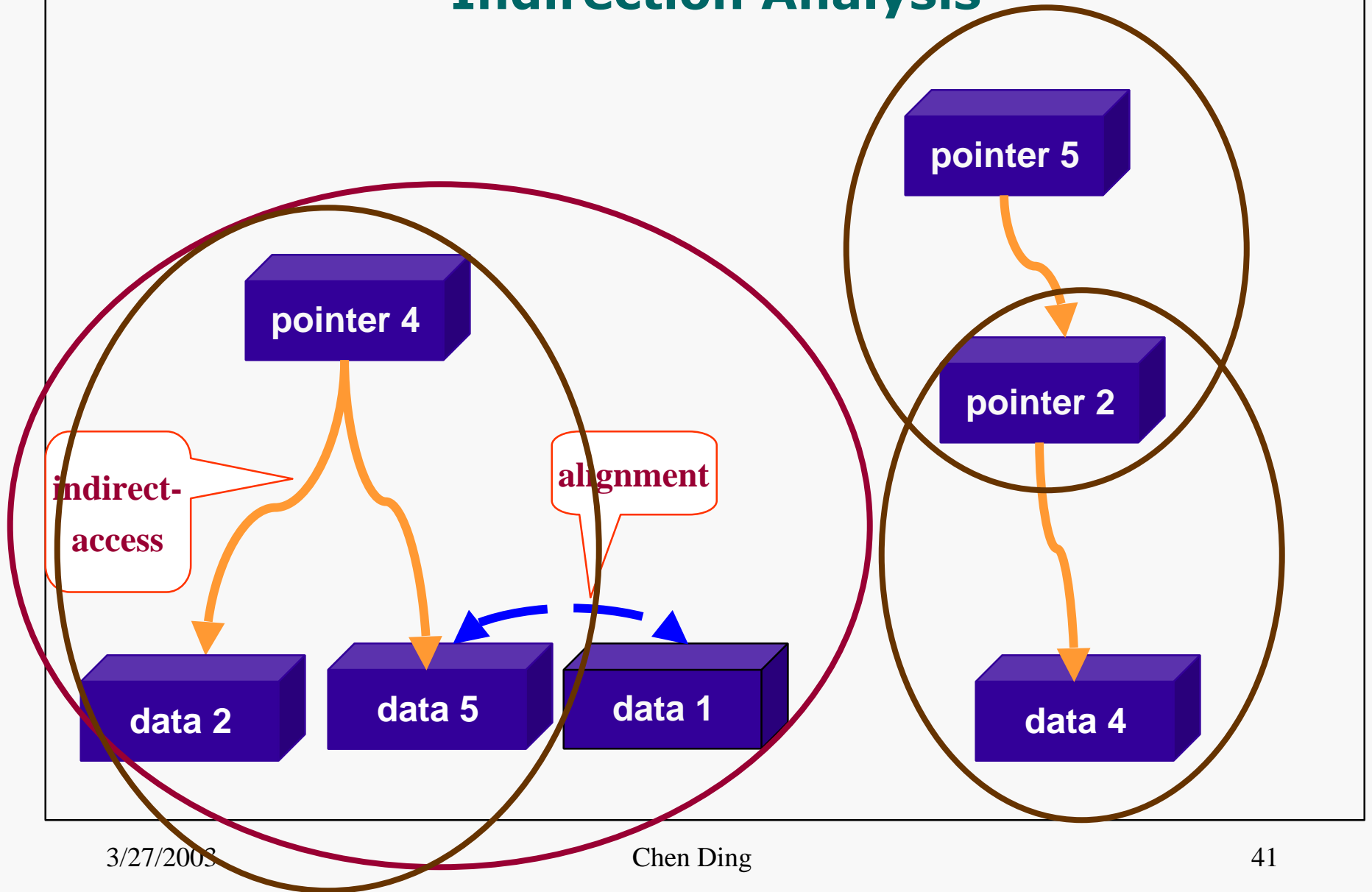
```
apply_packing(interactions[*],force[*],
              inter_map[*], update_map[*])
update_indirection_array(interactions[*],
                          update_map[*])
transform_data_array(location[*],update_map[*])

for each pair (i,j) in interactions
  compute_force( force[i], force[j] )
end for

for each object i
  update_location( location[i], force[i] )
end for
```



# Indirection Analysis

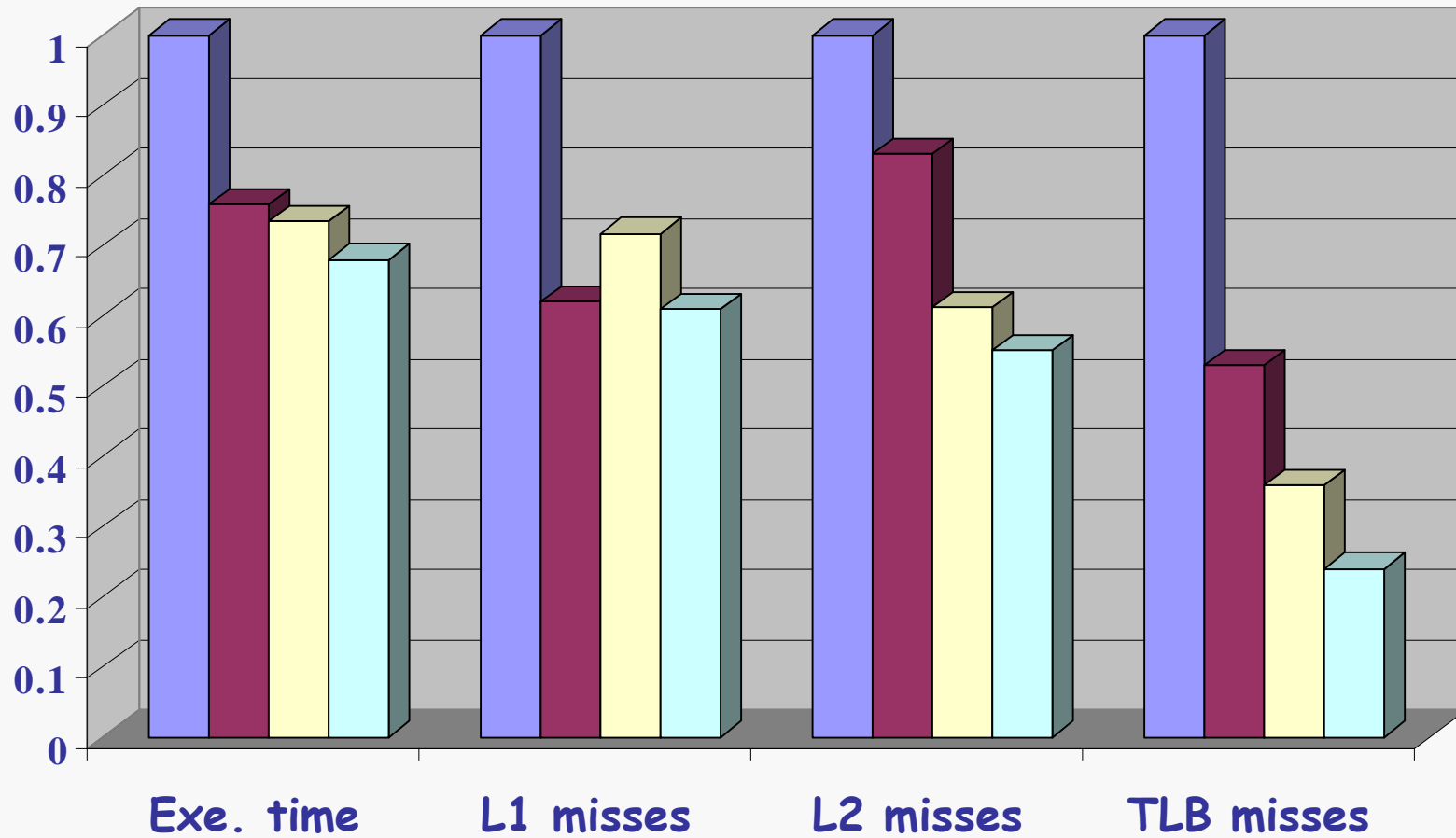


## DoD/Magi

- A real application from DoD Philips Lab
  - particle hydrodynamics
  - almost 10,000 lines of code
  - user supplied input of 28K particles
  - 22 arrays in major phases, split into 26
- Optimizations
  - grouped into 6 arrays
  - inserted 1114 indirections to guarantee correctness
  - optimization reorganized 19 more arrays
  - removed 379 indirections in loops
  - reorganized 45 arrays 4 times during execution

# Magi

original data regrouping base packing opt packing



## Overall Comparison

programs	L2 misses			TLB misses			Speedup over SGI
	NoOpt	SGI	New	NoOpt	SGI	New	
<b>Swim</b>	1.00	1.10	0.94	1.00	1.60	1.05	1.14
<b>Tomcatv</b>	1.00	0.49	0.39	1.00	0.010	0.010	1.17
<b>ADI</b>	1.00	0.94	0.53	1.00	0.011	0.005	2.33
<b>NAS/SP</b>	1.00	1.00	0.49	1.00	1.09	0.67	1.49
<b>Average</b>	1.00	0.88	0.59	1.00	0.68	0.43	1.52
<b>Moldyn</b>	1.00	0.99	0.19	1.00	0.77	0.10	3.02
<b>Mesh</b>	1.00	1.34	0.39	1.00	0.57	0.57	1.20
<b>Magi</b>	1.00	1.25	0.76	1.00	1.00	0.36	1.47
<b>NAS/CG</b>	1.00	0.95	0.15	1.00	0.97	0.03	4.36
<b>Average</b>	1.00	1.13	0.37	1.00	0.83	0.27	2.51

# Software Techniques Summary

main steps	sub-steps	example techniques (* studied in my work)
temporal reuse	global (multi-loop)	*loop fusion
	local (single loop)	blocking, register allocation
	dynamic	*dynamic partitioning
spatial reuse	global (inter-array)	*inter-array data regrouping
	local (intra-array)	loop permutation, array reshaping, combined schemes
	dynamic	*dynamic data packing
	cache interference	padding
latency tolerance	local (single loop)	data prefetching, instruction scheduling
program tuning & scheduling	global (whole program)	*balance model *bandwidth-based perf. tuning & prediction

## Summary of This Talk

- **Global transformations**
  - Combining both computation and data reordering at a large scale
- **Dynamic transformations**
  - Combining compile-time and run-time analysis and transformation
- **The first published compiler research**
  - splits and regroups global computation and data
  - for the whole program and at all times