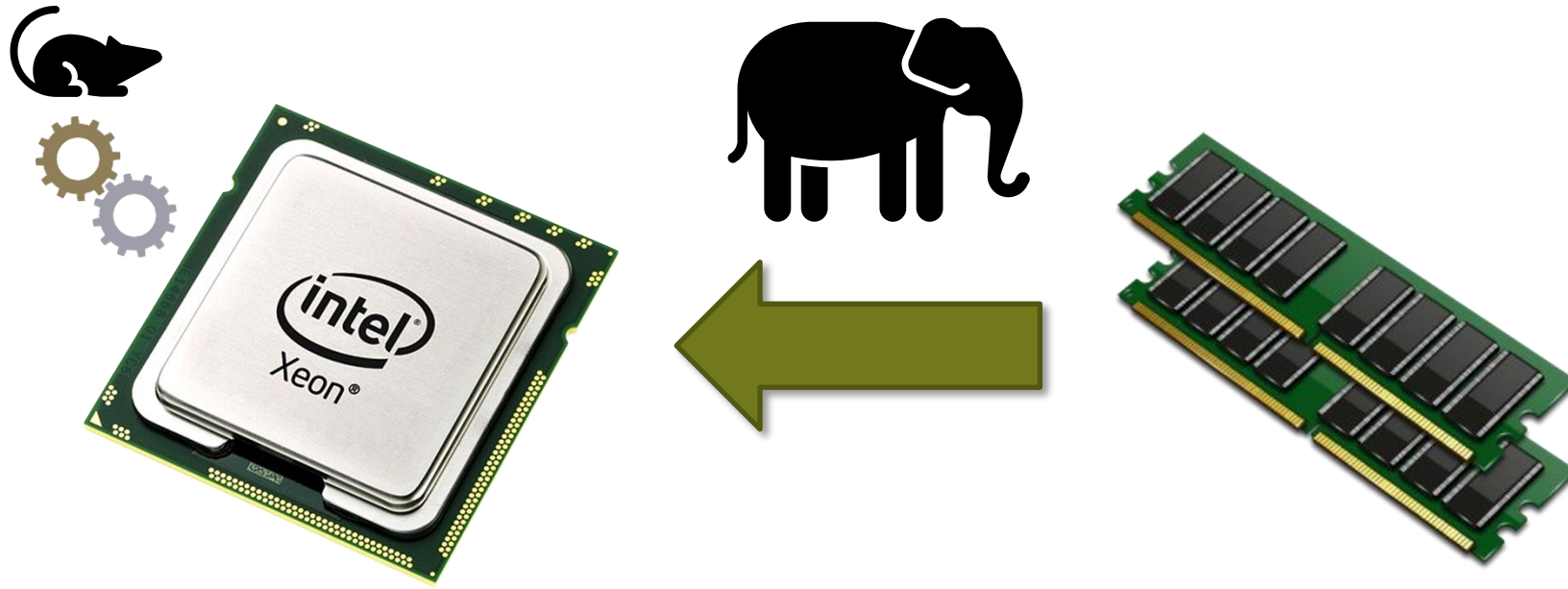# Boosting Performance Optimization with Interactive Data Movement Visualization

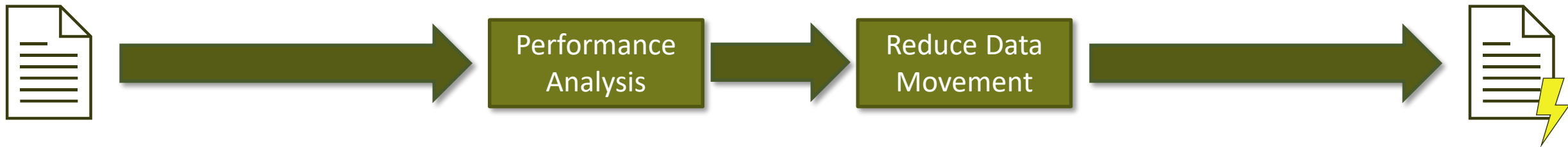**Philipp Schaad**[*], Tal Ben-Nun[*], Torsten Hoefler[*]

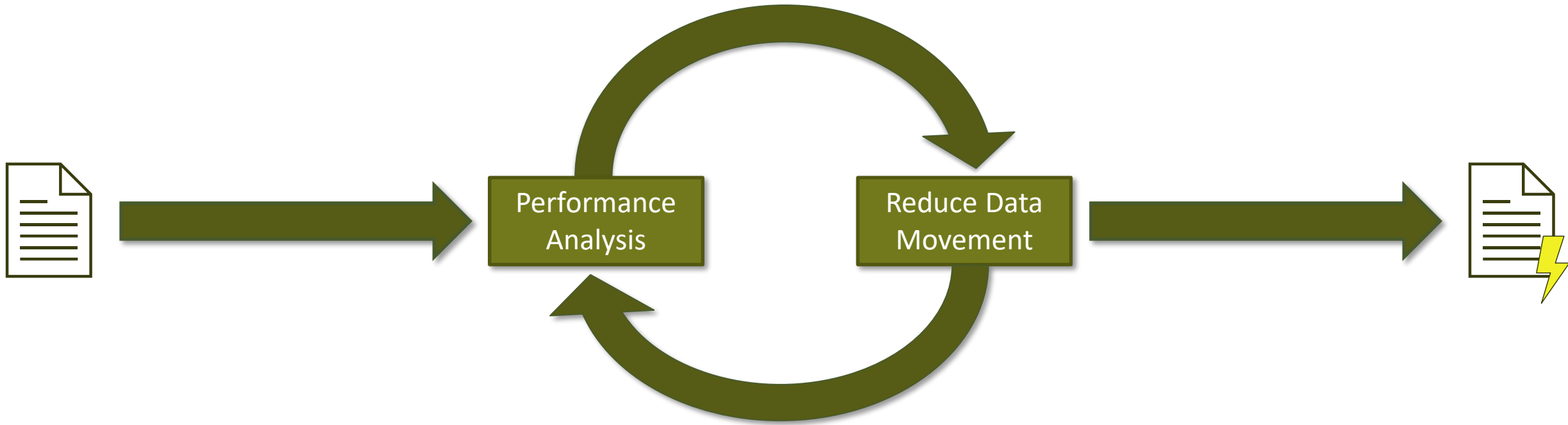[*] Scalable Parallel Computing Laboratory, ETH Zurich

# The Cost of Data Movement



Exploit *spatial locality* and *temporal locality*!

# Data Movement Optimization
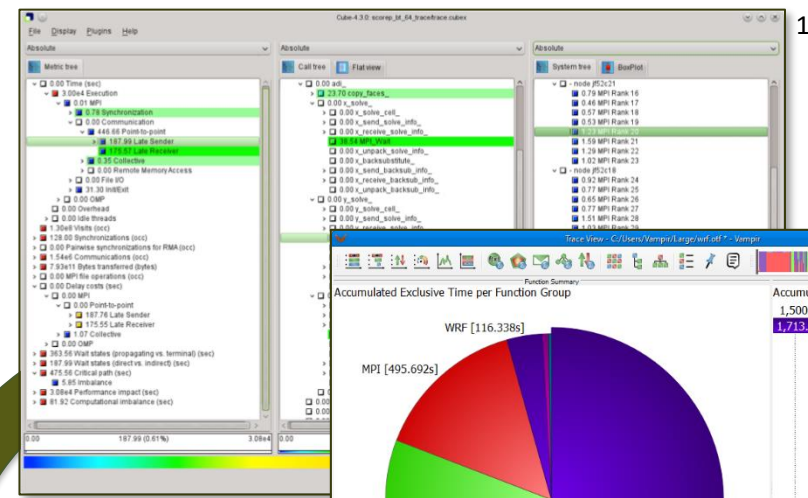
# Data Movement Optimization
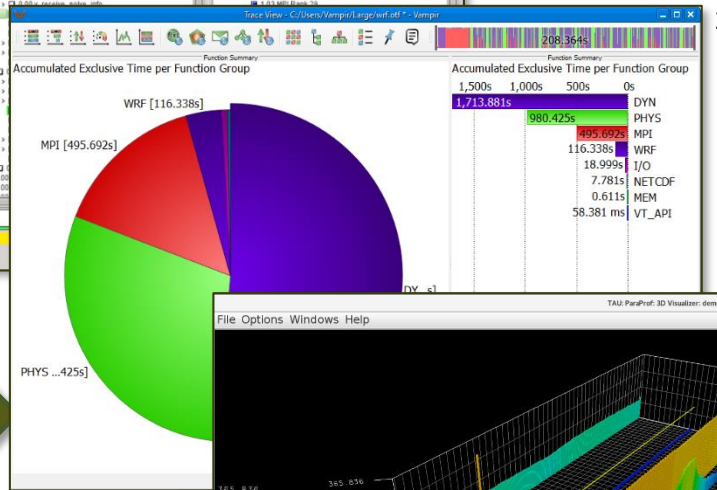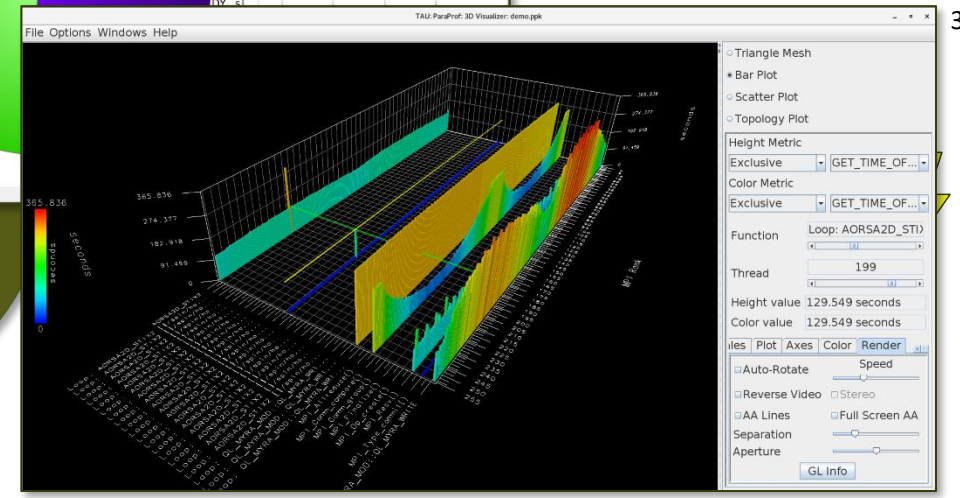
# Data Movement Optimization



PAPI
Intel Vtune
LIKWID
Perf

Performance Analysis

[1] Saviankou et al., Cube v4: From Performance Report Explorer to Performance Analysis Tool
[2] Nagel et al., VAMPIR: Visualization and Analysis of MPI Resources
[3] Bell et al., ParaProf: A Portable, Extensible, and Scalable Tool for Parallel Performance Profile Analysis

# Data Movement Optimization



Performance Analysis

PAPI
Intel Vtune
LIKWID
Perf

Requires Execution!

[1] Saviankou et al., Cube v4: From Performance Report Explorer to Performance Analysis Tool
[2] Nagel et al., VAMPIR: Visualization and Analysis of MPI Resources
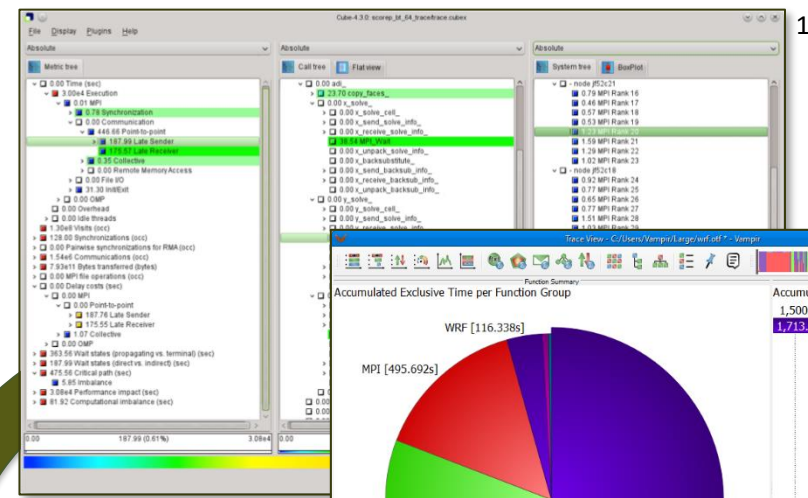[3] Bell et al., ParaProf: A Portable, Extensible, and Scalable Tool for Parallel Performance Profile Analysis
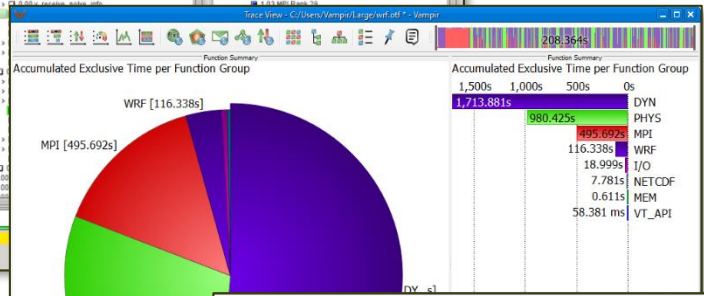
# Data Movement Optimization



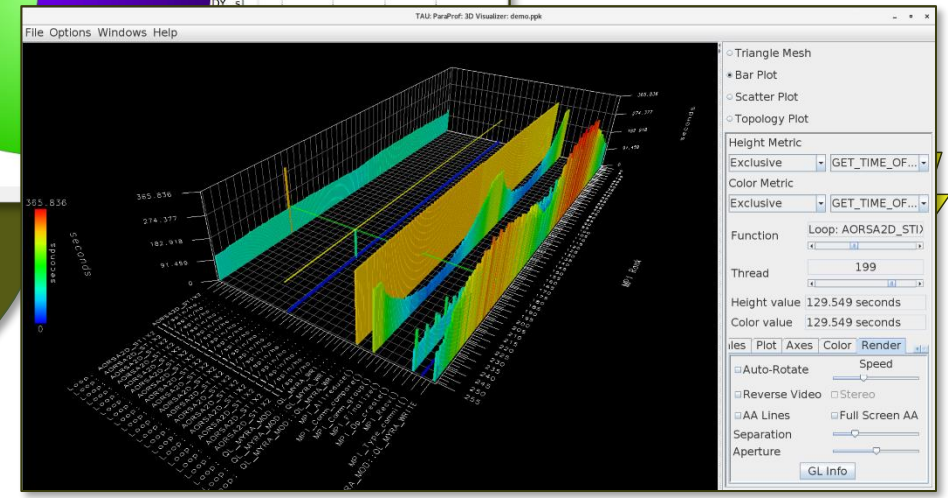**Performance analysis *without* program execution**
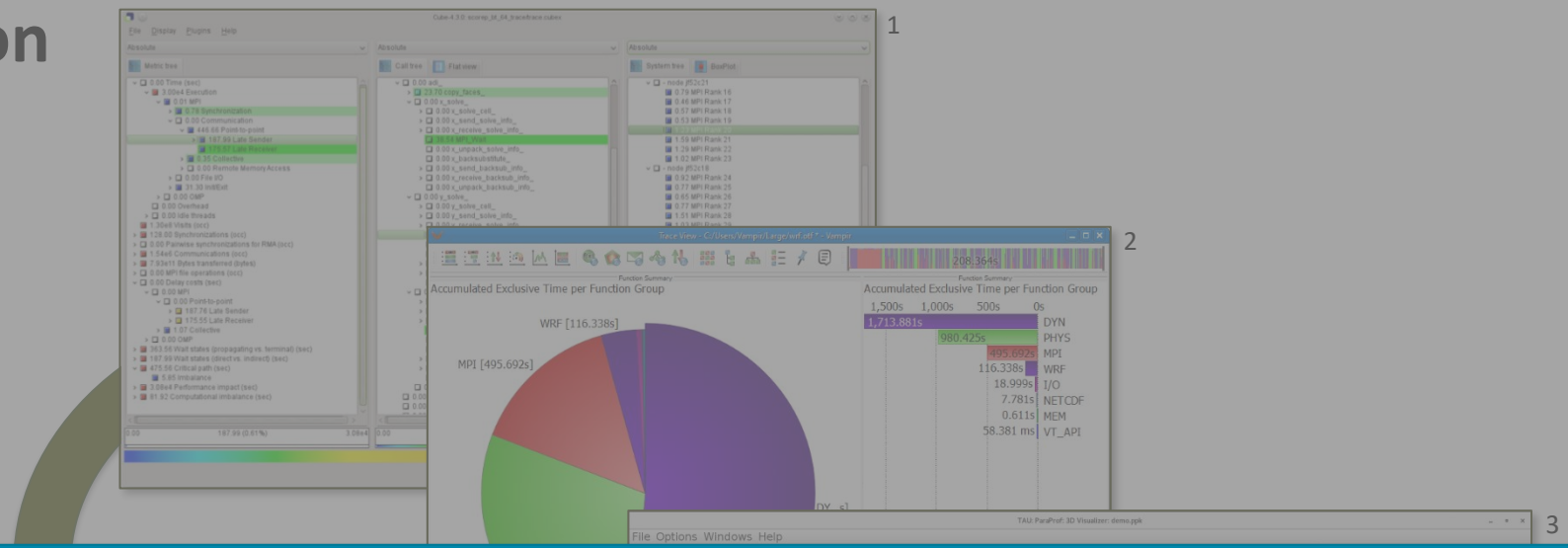
PAPI
Intel Vtune
LIKWID
Perf

Requires Execution!

[1] Saviankou et al., Cube v4: From Performance Report Explorer to Performance Analysis Tool
[2] Nagel et al., VAMPIR: Visualization and Analysis of MPI Resources
[3] Bell et al., ParaProf: A Portable, Extensible, and Scalable Tool for Parallel Performance Profile Analysis
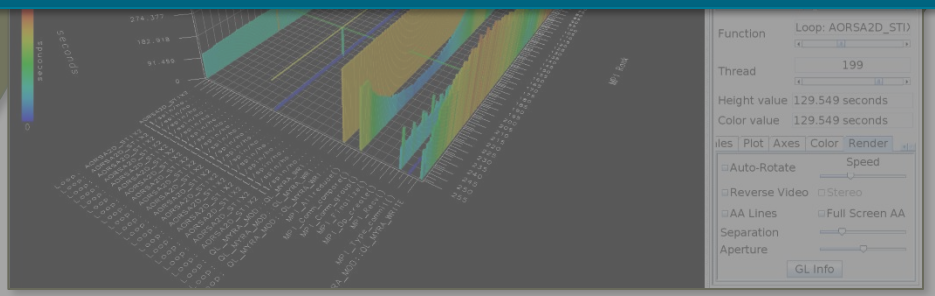
# Data Movement Optimization



Static Dataflow Analysis

Small-Scale Parametric Simulations

Dataflow IR

Performance Analysis

Reduce Data Movement

# Data Movement Optimization



**Small-Scale Parametric Simulations**

**Static Dataflow Analysis**

Dataflow IR → Performance Analysis ⇄ Reduce Data Movement →

DaCe [1]
Stateful DataFlow multiGraphs (SDFGs)

[1] Ben-Nun et al., Stateful Dataflow Multigraphs: A Data-Centric Model for Performance Portability on Heterogeneous Architectures

# Data Movement Optimization



DaCe [1]
Stateful DataFlow multiGraphs (SDFGs)

[1] Ben-Nun et al., Stateful Dataflow Multigraphs: A Data-Centric Model for Performance Portability on Heterogeneous Architectures

# Stateful DataFlow multiGraph (SDFG)

$$C = A \otimes B \qquad A \in \mathbb{R}^N, B \in \mathbb{R}^M, C \in \mathbb{R}^{N x M}$$

```
def outer_prod(A, B, C, N, M):
    for i in range(N):
        for j in range(M):
            C[i, j] = A[i] * B[j]
```



Data Containers

Parallel Region (Map)

Computations

Data Movement

State

# Static Dataflow Analysis

**Data Movement Volume**

# Static Dataflow Analysis

**Data Movement Volume**

1. Derive volume for computations
2. Propagate through graph

**Arithmetic Operations Count**

1. Count operations in AST of computations
2. Propagate through graph

**Operational Intensity**

# Static Dataflow Analysis

# Static Dataflow Analysis

$N = 8$
$M = 8$



Substitute symbols

64 Operations
Intensity: $\frac{1}{3}$

1 Operation
Intensity: $\frac{1}{3}$

64 Operations
Intensity: $\frac{1}{3}$

# Static Dataflow Analysis

$N = 8$
$M = 64$

Substitute symbols

Change symbol values to perform *scaling analysis*



**512** Operations
Intensity: $\frac{1}{3}$

**1** Operation
Intensity: $\frac{1}{3}$

**512** Operations
Intensity: $\frac{1}{3}$

# Visualization

$N = 8$
$M = 64$



Visualize data by overlaying a *heatmap*

Low volume

High volume

# Visualization

$N = 8$
$M = 64$



512 Operations
Intensity: $\frac{1}{3}$

1 Operation
Intensity: $\frac{1}{3}$

512 Operations
Intensity: $\frac{1}{3}$

Visualize data by overlaying a *heatmap*

Low operation count

High operation count

# Visualization

$N = 8$
$M = 64$

In-Situ performance data reduces context switching

Visualize data by overlaying a *heatmap*

Low operation count

High operation count



**512** Operations
Intensity: $\frac{1}{3}$

**1** Operation
Intensity: $\frac{1}{3}$

**512** Operations
Intensity: $\frac{1}{3}$

# Optimizing BERT Transformer Encoder



Data movement heatmap

# Optimizing BERT Transformer Encoder



[0:P, 0:H, 0:B, 0:SM]
[0:P, 0:H, 0:B, 0:SM]
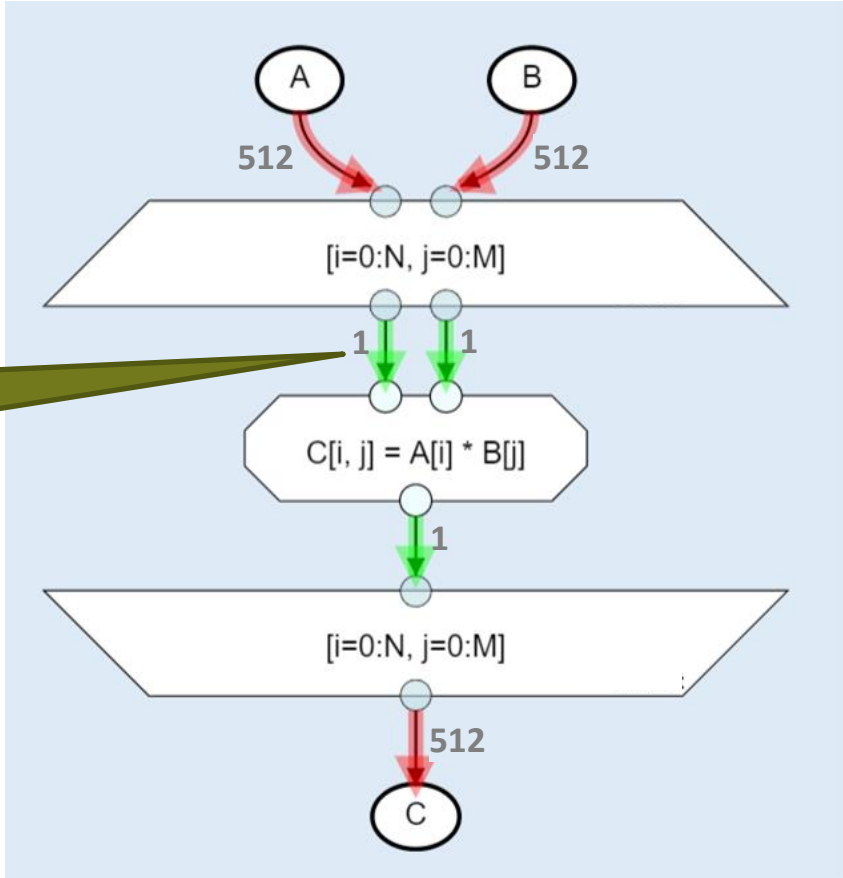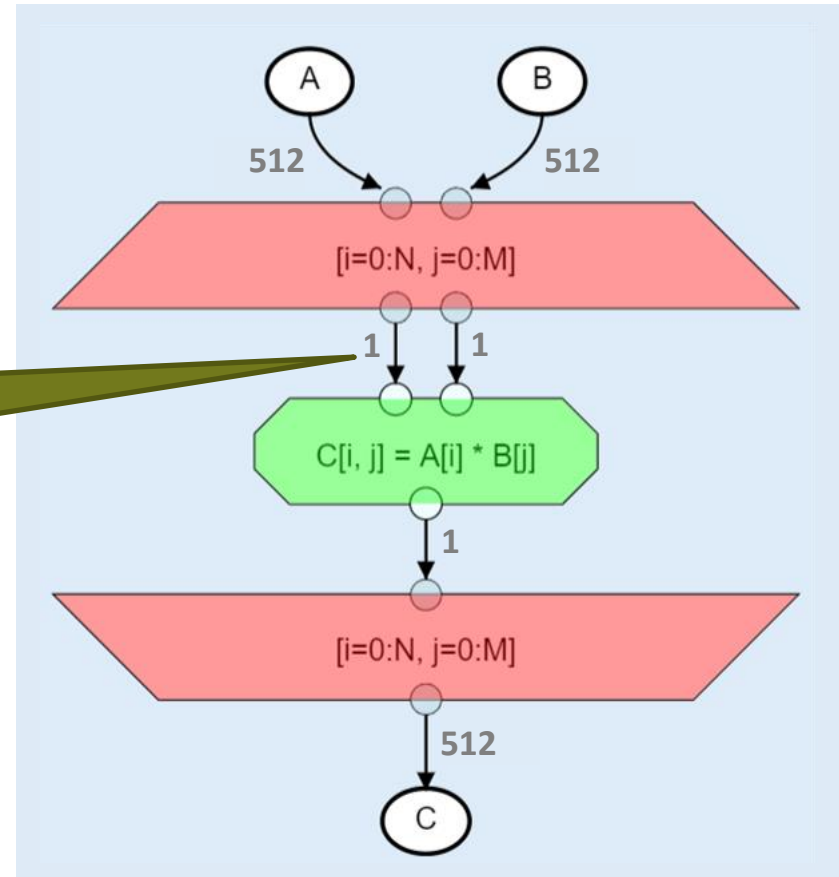[0:H, 0:B, 0:SM, 0:SM]
[0:H, 0:B, 0:SM]
[0:H, 0:B, 0:SM, 0:SM]
[0:H, 0:B, 0:SM, 0:SM]
[0:H, 0:B, 0:SM]
[0:P, 0:H, 0:B, 0:SM]
[0:H, 0:B, 0:SM, 0:SM]
[0:B, 0:SM, 0:N]
[0:B, 0:SM, 0:N]
[0:B, 0:SM, 0:N]
[0:B, 0:SM]
[0:B, 0:SM]
[0:B, 0:SM]
[0:B, 0:SM, 0:N]
[0:B, 0:SM, 0:emb]
[0:B, 0:SM, 0:emb]
[0:B, 0:SM, 0:emb]
[0:B, 0:SM, 0:N]
[0:B, 0:SM, 0:N]
[0:B, 0:SM, 0:N]
[0:B, 0:SM]
[0:B, 0:SM]
[0:B, 0:SM]
[0:B, 0:SM, 0:N]

[0:H, 0:B, 0:SM, 0:SM]
[0:H, 0:B, 0:SM]
[0:H, 0:B, 0:SM, 0:SM]
[0:H, 0:B, 0:SM, 0:SM]
[0:H, 0:B, 0:SM]
[0:H, 0:B, 0:SM, 0:SM]

einsum_gemm
[0:B, 0:SM, 0:emb]
[0:B, 0:SM, 0:emb]
[0:B, 0:SM, 0:emb]
einsum_gemm

Loops with similar bounds

22

# Optimizing BERT Transformer Encoder

# Optimizing BERT Transformer Encoder

# Optimizing BERT Transformer Encoder



Operational intensity heatmap

# Optimizing BERT Transformer Encoder

# Optimizing BERT Transformer Encoder

30.2x Speedup

# Close-Up Reuse Analysis

Simulate data reuse behavior

$$C = A \otimes B \qquad A \in \mathbb{R}^N, B \in \mathbb{R}^M, C \in \mathbb{R}^{N \times M}$$

# Close-Up Reuse Analysis

Simulate data reuse behavior

$$C = A \otimes B \qquad A \in \mathbb{R}^N, B \in \mathbb{R}^M, C \in \mathbb{R}^{N x M}$$



Specify program region

# Close-Up Reuse Analysis

Simulate data reuse behavior

$$C = A \otimes B \qquad A \in \mathbb{R}^3, B \in \mathbb{R}^4, C \in \mathbb{R}^{3x4}$$

Specify small example input parameters

$$N = 3$$
$$M = 4$$



Specify program region

# Close-Up Reuse Analysis

Simulate data reuse behavior

$$C = A \otimes B \qquad A \in \mathbb{R}^3, B \in \mathbb{R}^4, C \in \mathbb{R}^{3x4}$$

Specify small example input parameters

$$N = 3$$
$$M = 4$$

Specify program region

# Visualizing High-Dimensional Data

$$w \in \mathbb{R}^{C_{out} \times C_{in} \times K_y \times K_X}$$

4D weights of a convolution

$$C_{out} = 2$$
$$C_{in} = 3$$
$$K_Y = 4$$
$$K_X = 4$$

# Visualizing High-Dimensional Data

$$w \in \mathbb{R}^{C_{out} \times C_{in} \times K_y \times K_X}$$

$$C_{out} = 2$$
$$C_{in} = 3$$
$$K_Y = 4$$
$$K_X = 4$$

$$C_{out} = 2$$



$$K_X = 4$$

$$K_Y = 4$$

$$C_{in} = 3$$

Convolution operation

$$y[i, j, k, l] \mathrel{+}= x[i, m, k+ky, l+kx] * w[j, m, ky, kx]$$

# Access Pattern Simulation

Visually play back access pattern

$$y[i, j, k, l] \mathrel{+}= x[i, m, k+ky, l+kx] * w[j, m, ky, kx]$$

# Access Pattern Simulation



Flatten time dimension with heatmap

# Access Pattern Simulation

# Data Layout Visualization



float64 / double → element size = 8 bytes

Exposes data layout

Determine cache line using strides, line size, and element size

A

B

i=0:9   0 ● ━ 8   0
j=0:15   0 ● ━ 14   0
k=0:10   0 ● ━ 9   0

$C[i, j] \mathrel{+}= A[i, k] * B[k, j]$

C

Data layout & access pattern → **spatial locality**

# Temporal Locality

**Stack distance**, cache line granularity

Accesses to _unique_ addresses since last reference



Cache Line Size (bytes): 64

Reuse Distance Threshold:
4

**View Modes**

☐ Access Pattern / Number of Accesses
☐ Reuse Distance (Stack Distance)
◉ Median ○ Min ○ Max ○ Misses
☐ Physical Data Movement
☐ Cache Lines

A[5,6]

A[i, k] * B[k, j]

i=0:9    j=0:15    k=0:10

0 ●——— 8    0 ●——— 14    0 ●——— 9
0              0              0

$$C[i, j] \mathrel{+}= A[i, k] * B[k, j]$$

39

# Cache Misses

1. Cold miss

   Access with stack distance = ∞

2. Capacity miss

   Assuming LRU

   Access with stack distance $\geq t_d$ stack distance threshold

3. Conflict miss

   Not counted in fully-associative cache

   Calculations generalizeable [1][2]

Cache Line Size (bytes): 64

Reuse Distance Threshold:
4

$t_d = 5$

**View Modes**

☐ Access Pattern / Number of Accesses
☑ Reuse Distance (Stack Distance)
◉ Median ○ Min ○ Max ○ Misses
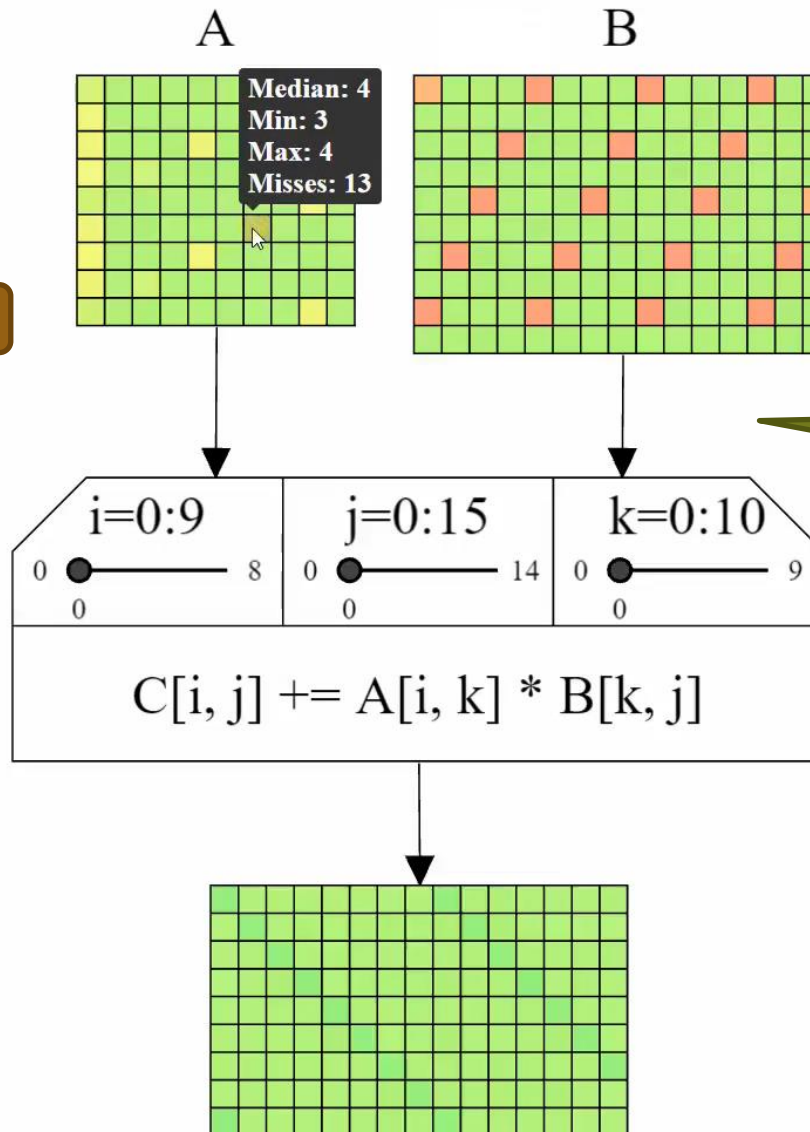☐ Physical Data Movement
☐ Cache Lines

A

Median: 4
Min: 3
Max: 4
Misses: 13

B

Physical data movement = #misses x cache line size

i=0:9          j=0:15          k=0:10
0 ●—— 8     0 ●—— 14     0 ●—— 9
  0              0              0

C[i, j] += A[i, k] * B[k, j]

C

[1] McKinley and Temam, Quantifying Loop Nest Locality Using SPEC'95 and the Perfect Benchmarks
[2] Beyls and D'Hollander, Reuse distance as a metric for cache behavior

40

# Stencil Optimization

$I = 8$
$J = 8$
$K = 5$

Accesses spread over non-contiguous dimension

Cache line shows K as contiguous dimension

Original sizes:
$I = 256$
$J = 256$
$K = 160$
Scaling Factor x32

in_field [I+4, J+4, K]

coeff [I, J, K]

i=0:8    j=0:8    k=0:5
0 ● 7    0 ● 7    0 ● 4
0        0        0

hdiff

out_field [I, J, K]

Cache Line Size (bytes): 64

Reuse Distance Threshold:
9

**View Modes**

☐ Access Pattern / Number of Accesses
☐ Reuse Distance (Stack Distance)
○ Median ○ Min ○ Max ◉ Misses
☐ Physical Data Movement
☑ Cache Lines
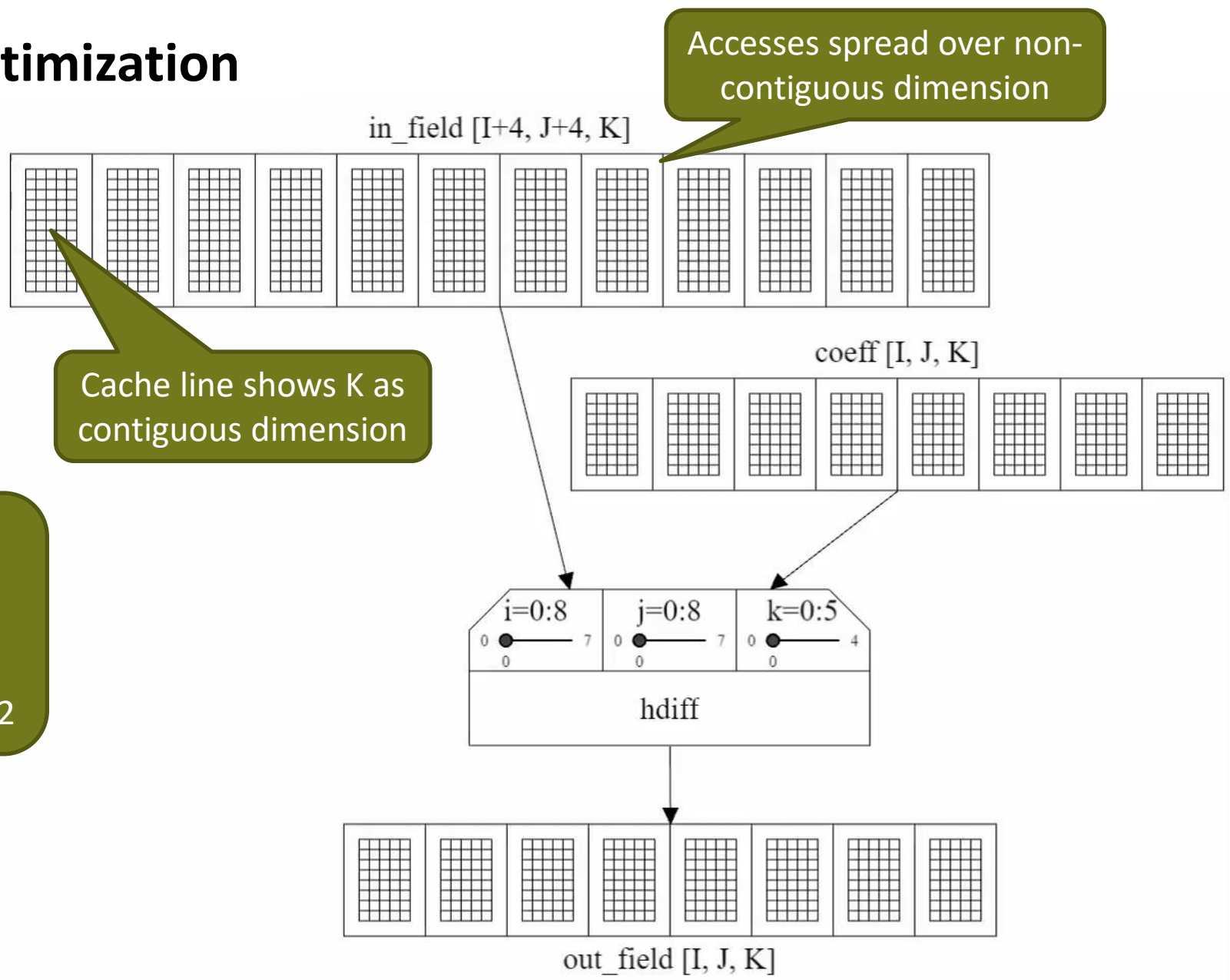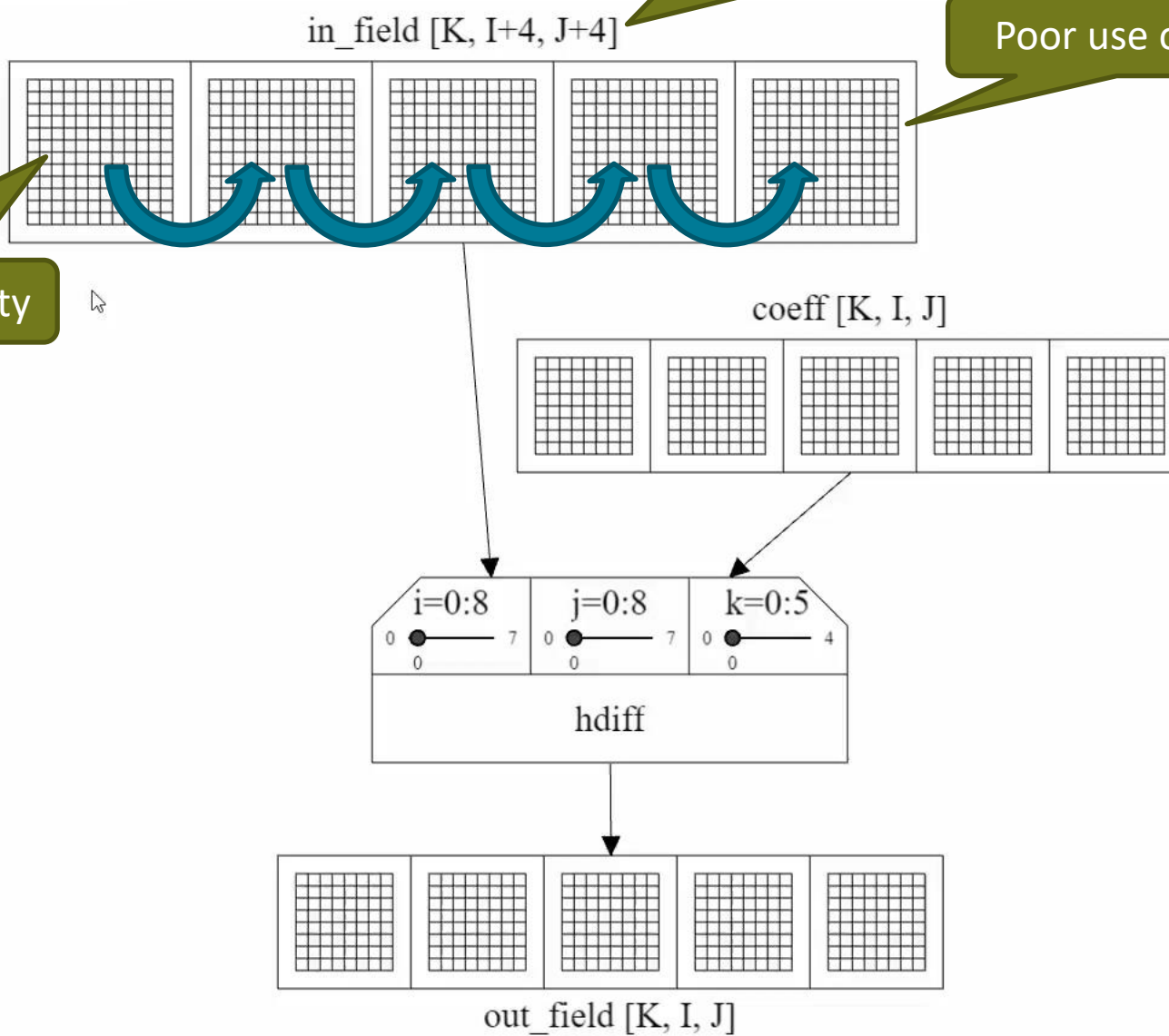
# Stencil Optimization

$I = 8$
$J = 8$
$K = 5$



Reshape data containers

Poor use of cache

Better use of spatial locality

in_field [K, I+4, J+4]

coeff [K, I, J]

i=0:8    j=0:8    k=0:5

hdiff

out_field [K, I, J]

Size (bytes): 64
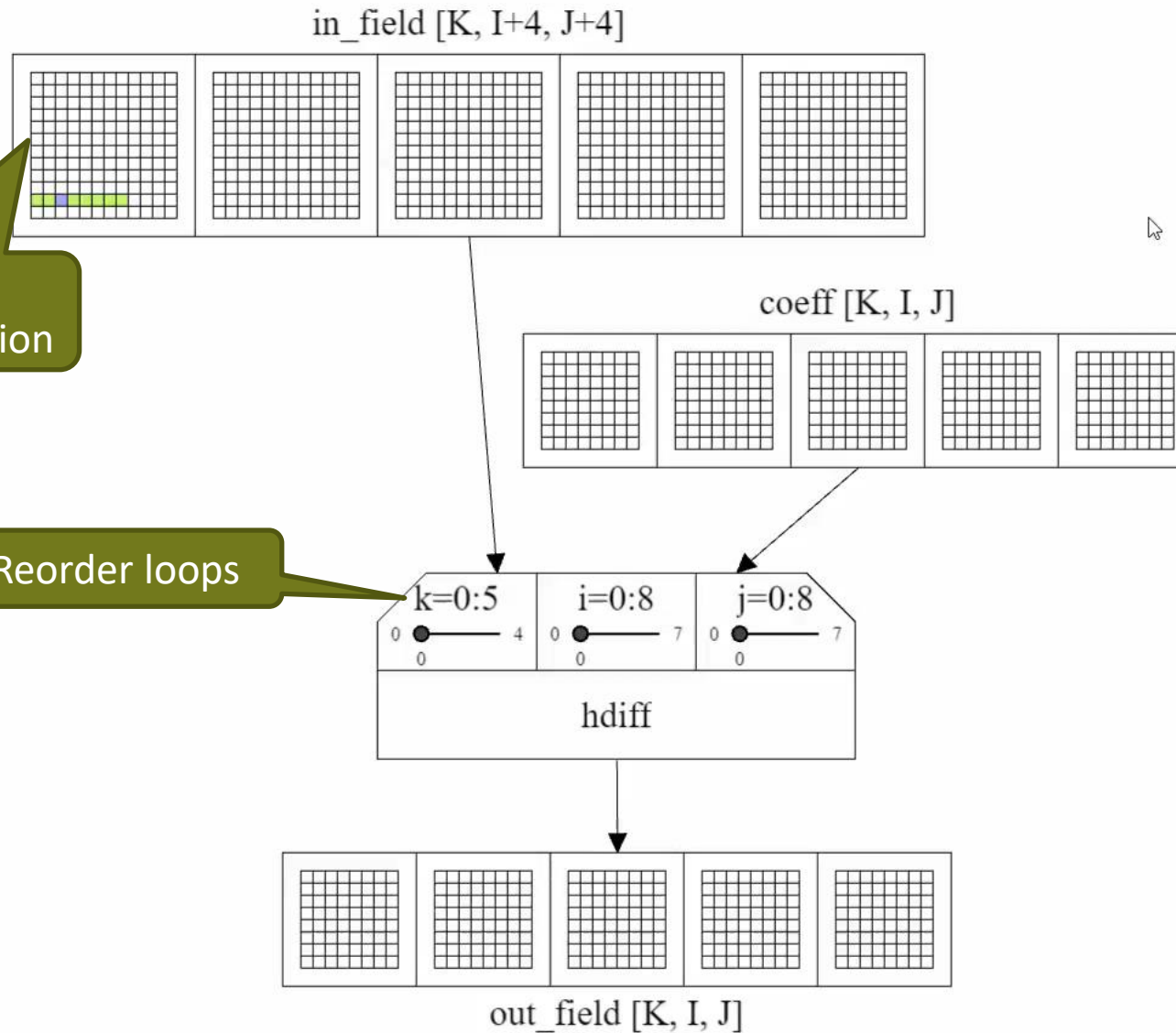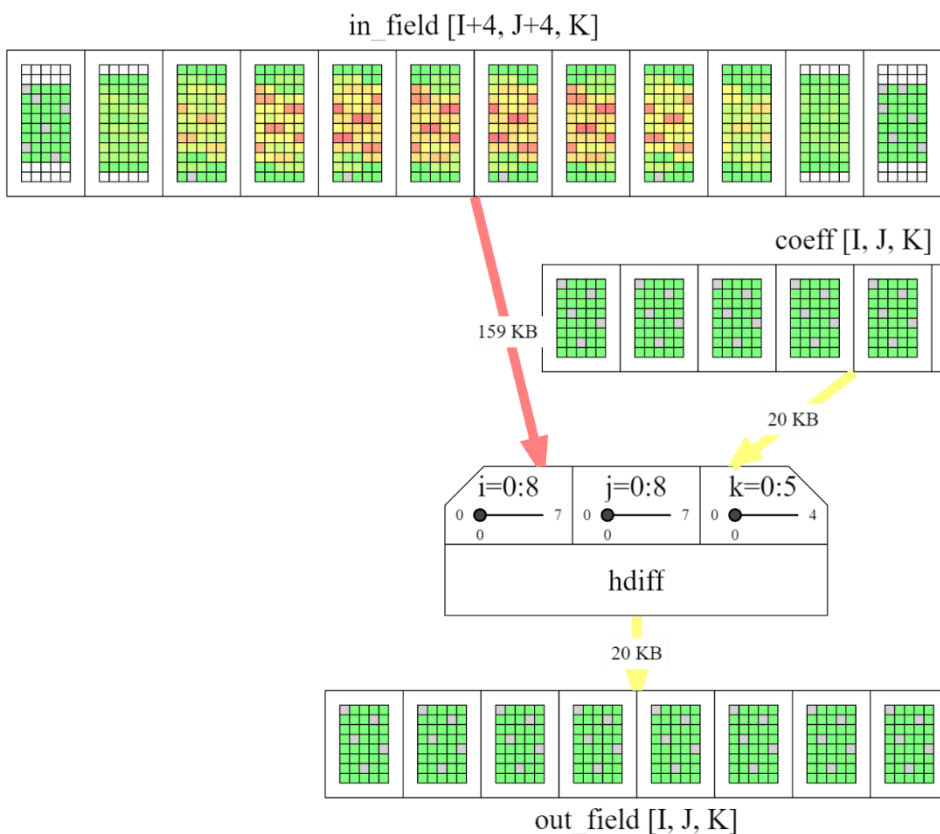
nce Threshold:
9

**View Modes**

☐ Access Pattern / Number of Accesses
☐ Reuse Distance (Stack Distance)
○ Median  ○ Min  ○ Max  ● Misses
☐ Physical Data Movement
☑ Cache Lines

# Stencil Optimization

$I = 8$
$J = 8$
$K = 5$



Iterates over contiguous dimension

Reorder loops

in_field [K, I+4, J+4]

coeff [K, I, J]

k=0:5    i=0:8    j=0:8
0 ●━━4  0 ●━━7  0 ●━━7
0        0        0

hdiff

out_field [K, I, J]

Cache Line Size (bytes): 64
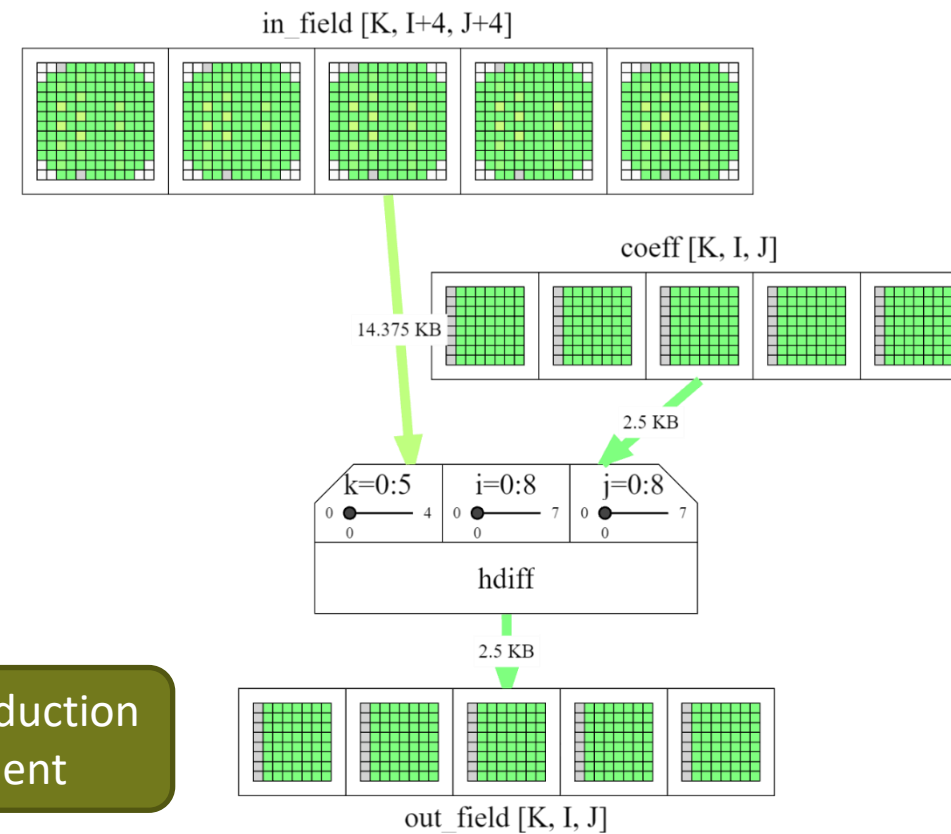
Reuse Distance Threshold:
9

**View Modes**

☐ Access Pattern / Number of Accesses
☐ Reuse Distance (Stack Distance)
○ Median  ○ Min  ○ Max  ● Misses
☐ Physical Data Movement
☑ Cache Lines

# Stencil Optimization

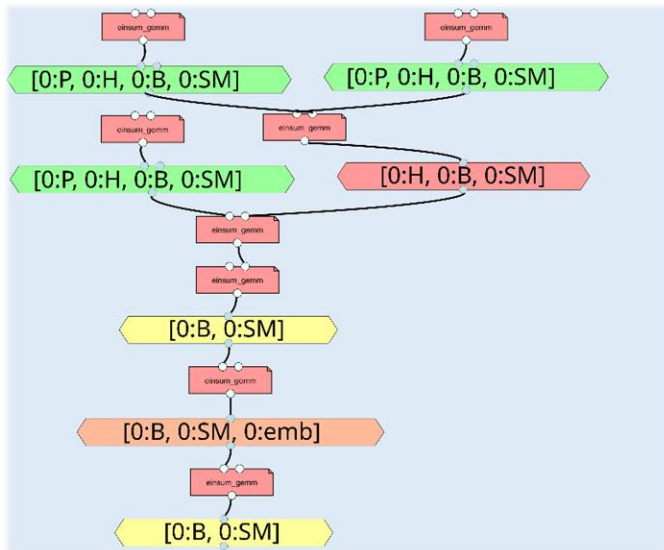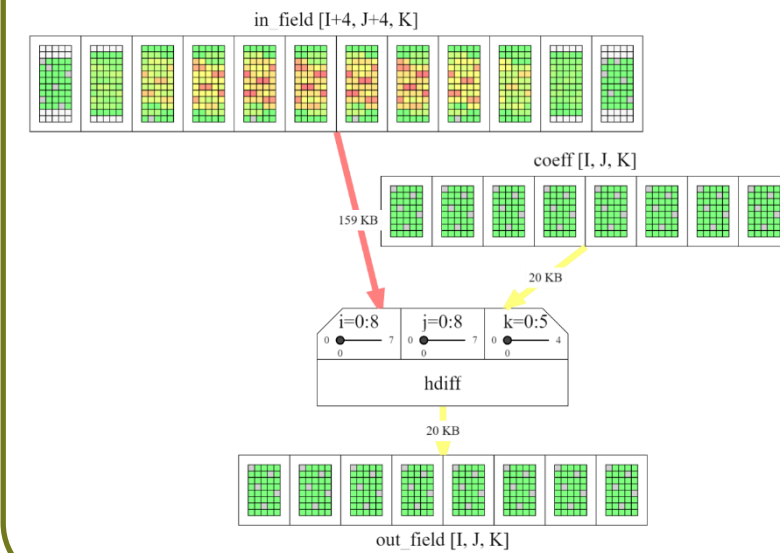Predicted 10.3x reduction in data movement

138x Speedup

9.6x Reduction in cache misses

# Conclusion



Global Data Movement



Fine-Grained Data Reuse
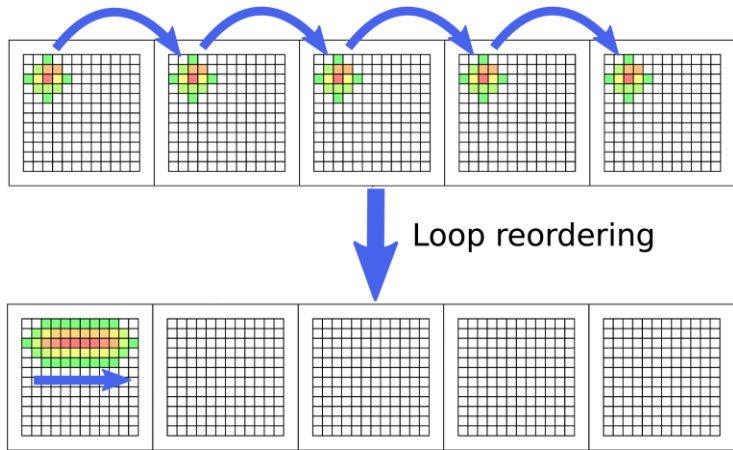
# Where Next?

## Automatic Optimization



Loop reordering

## Hardware Modelling



Reg
L1
L2
L3
Main Memory
Disk
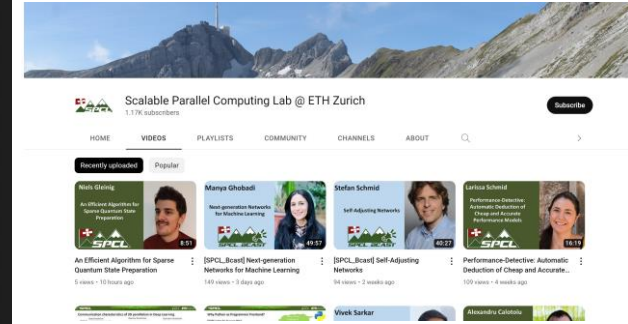Network

## Educational Tool

Thank you!

youtube.com/@spcl

twitter.com/spcl_eth

spcl.inf.ethz.ch

github.com/spcl

https://marketplace.visualstudio.com/items?itemName=phschaad.sdfv

https://github.com/spcl/dace-vscode