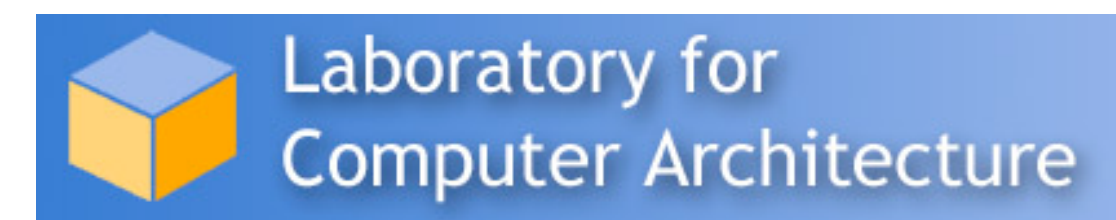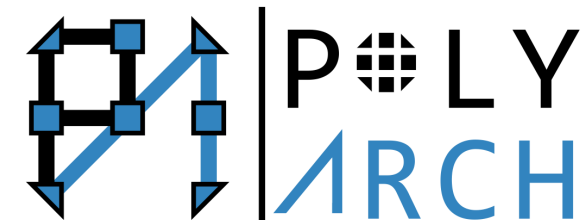# Infinity Stream

## Portable and Programmer-Friendly In-/Near-Memory Fusion
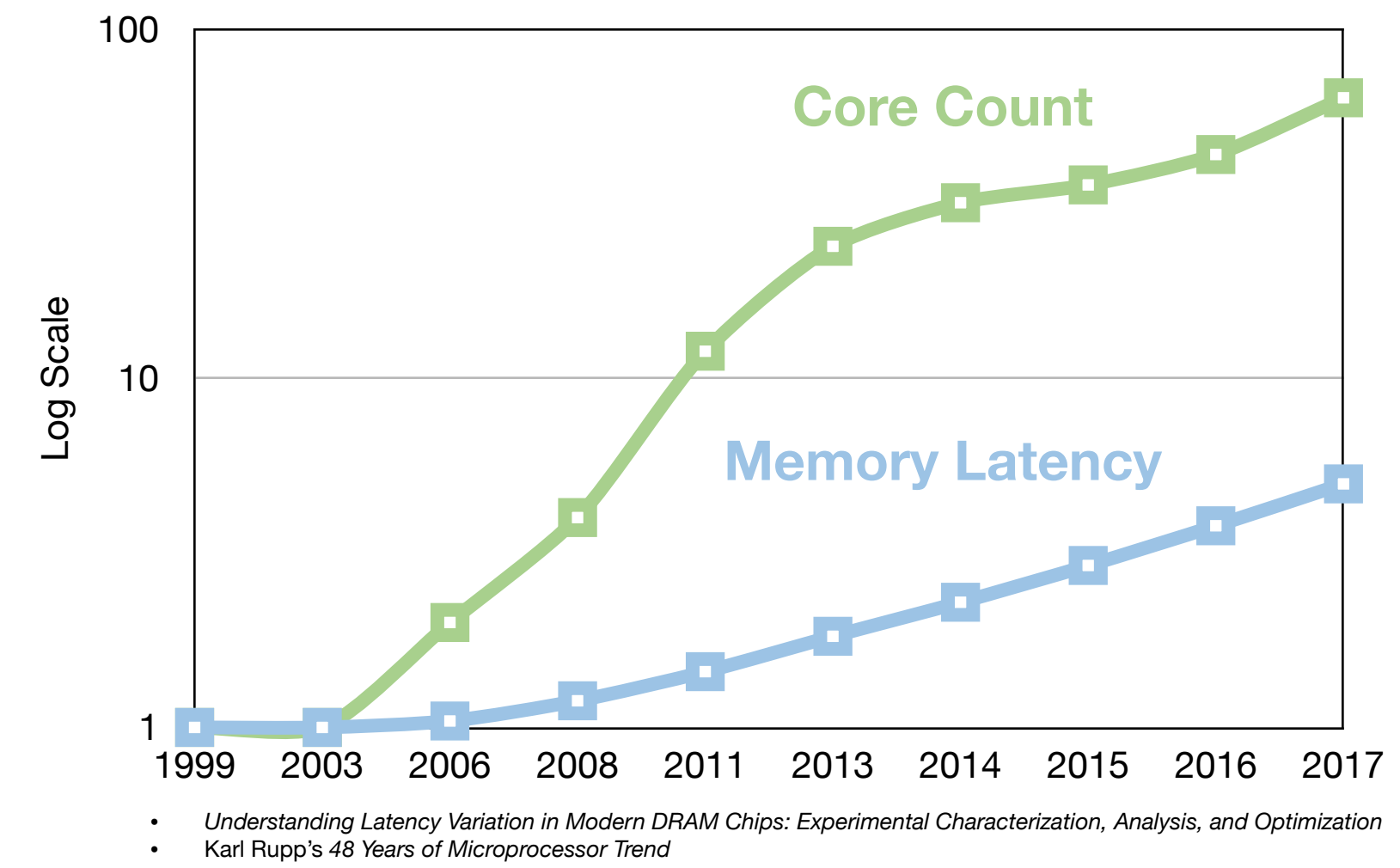
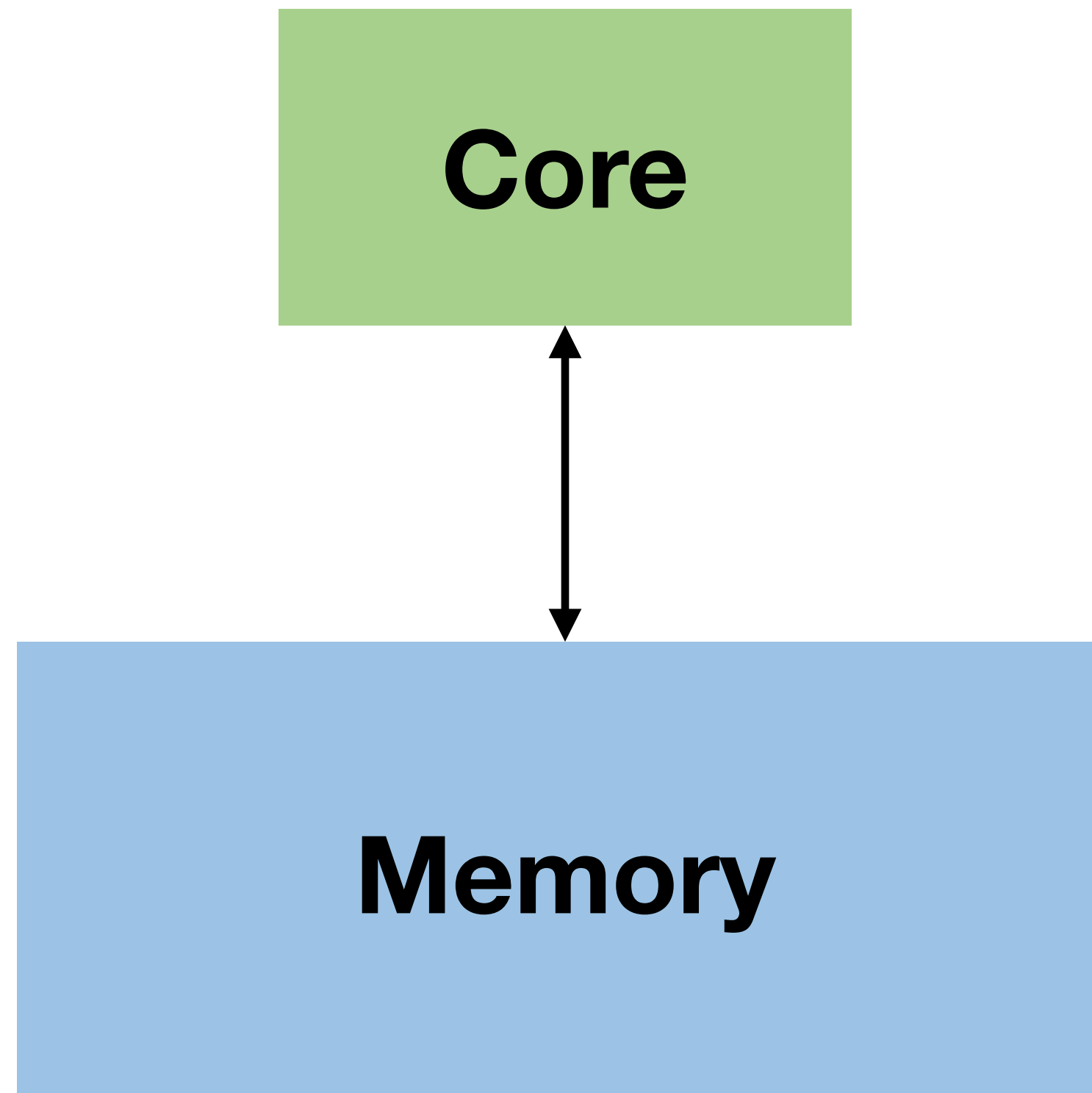Zhengrong Wang[1], Christopher Liu[1], Aman Arora[2], Lizy John[2], Tony Nowatzki[1]
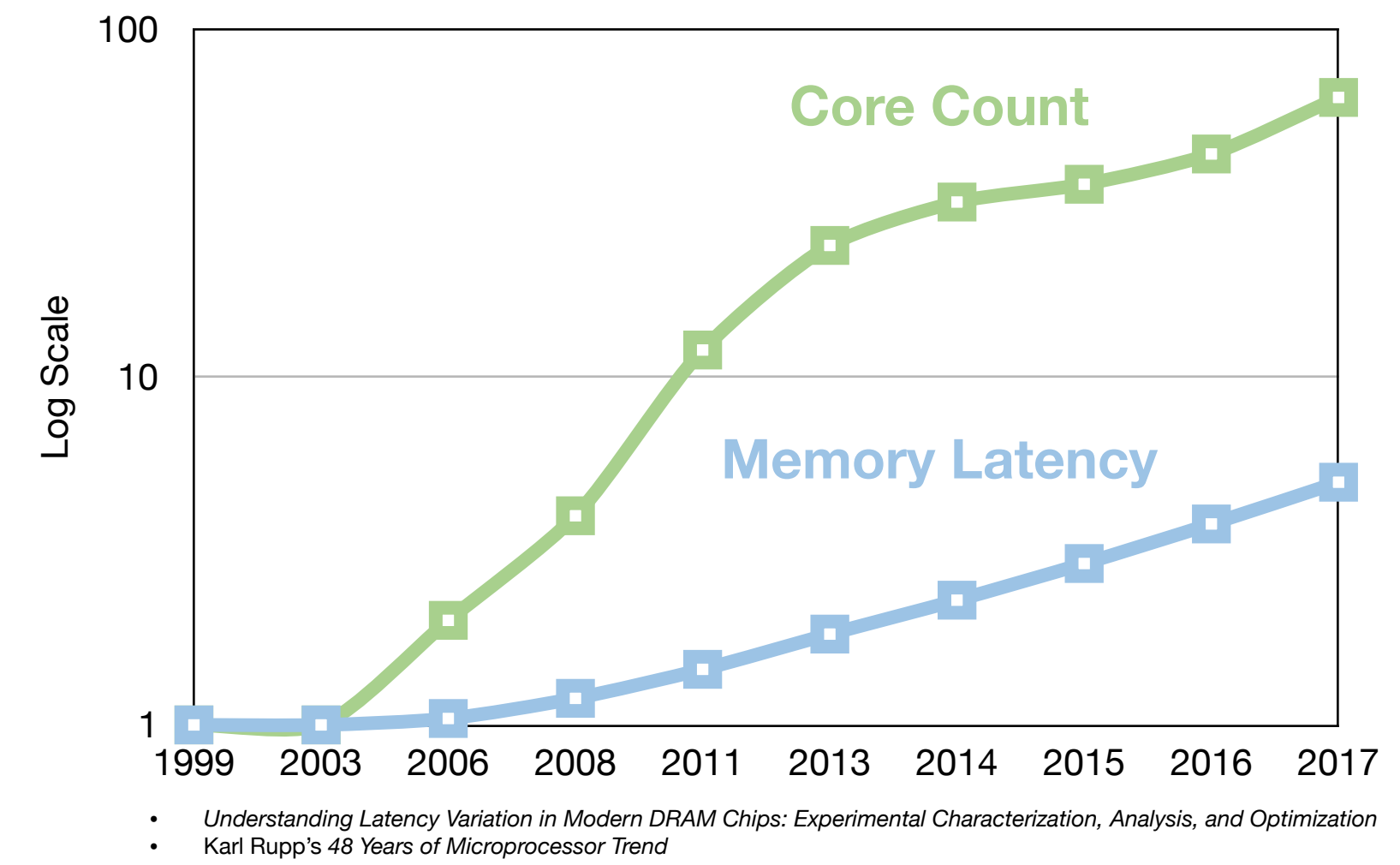
[1]UCLA, [2]UT Austin

# Von Neumann Bottleneck

*Performance limited by memory bottleneck*



Core

Memory

- *Understanding Latency Variation in Modern DRAM Chips: Experimental Characterization, Analysis, and Optimization*
- Karl Rupp's *48 Years of Microprocessor Trend*

# Von Neumann Bottleneck

*Performance limited by memory bottleneck*



**Expensive data movement**

- *Understanding Latency Variation in Modern DRAM Chips: Experimental Characterization, Analysis, and Optimization*
- Karl Rupp's *48 Years of Microprocessor Trend*

# Compute Paradigms

**More Expressive**
Less Parallelism

Less Expressive
**More Parallelism**

## In-Core

*Von Neumann Model*

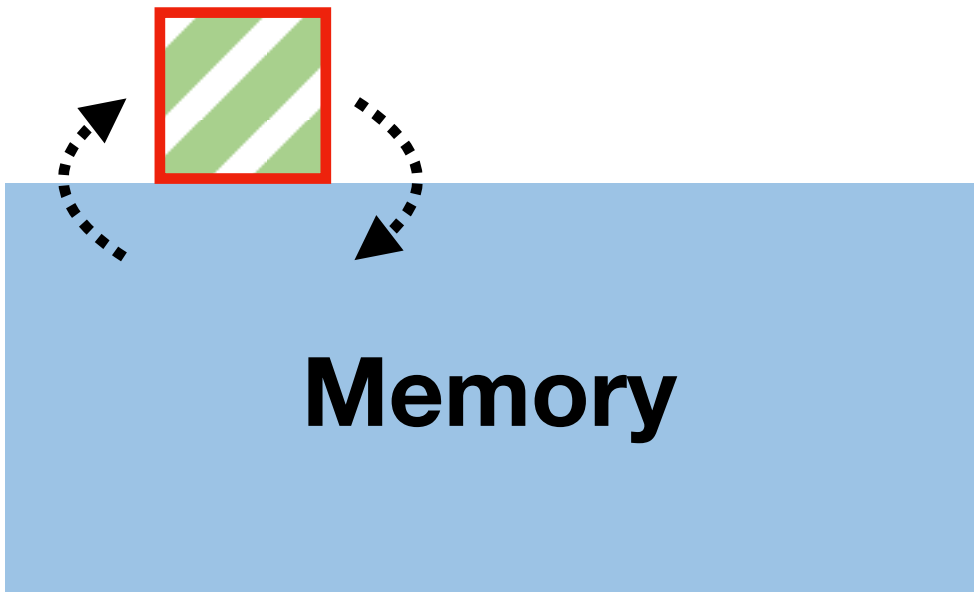## Near-Memory

*Spatially Local Simple Cores*

## In-Memory

*Massive Vector Processing*

👍 Supports complex control flow

👎 Von Neumann bottleneck

👍 Programs similar to in-core without control flow

👎 Limited parallelism
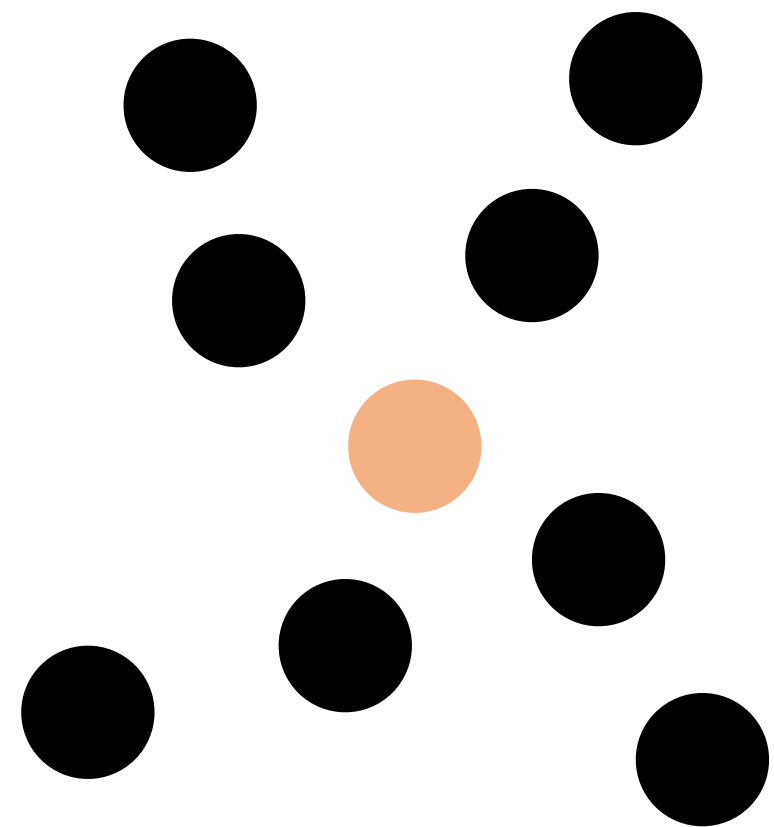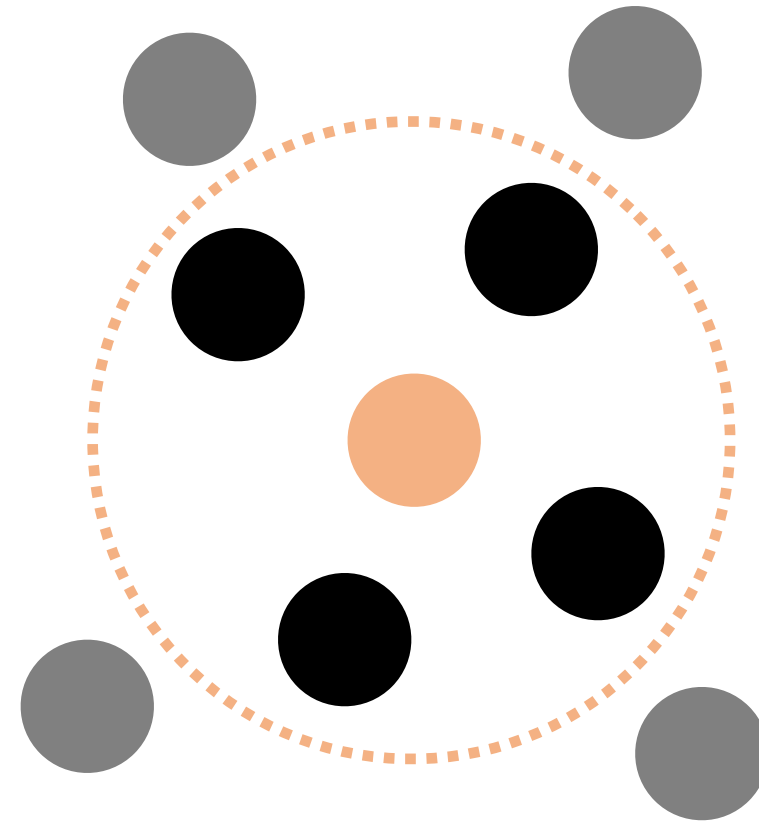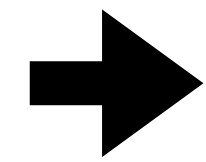
👍 Massive amounts of parallelism
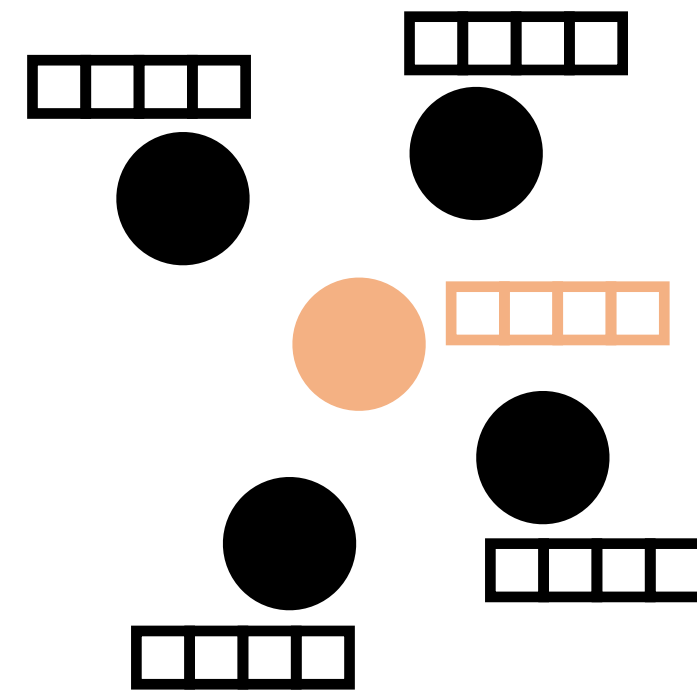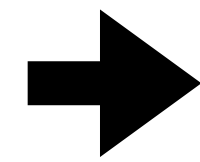
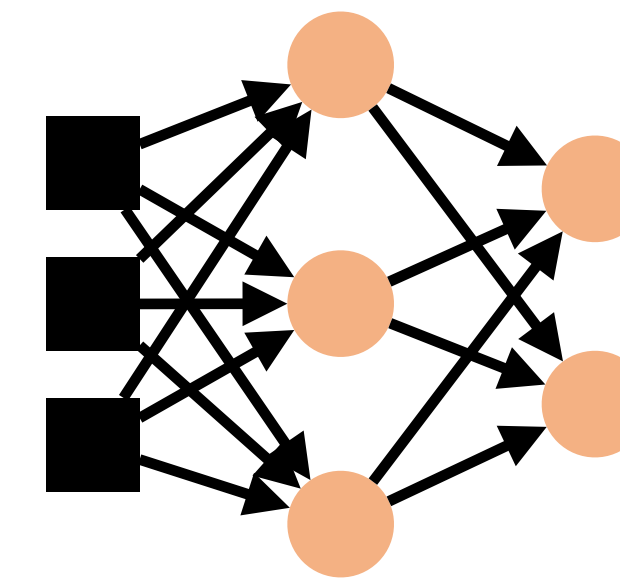👎 Difficult to program due to many restrictions
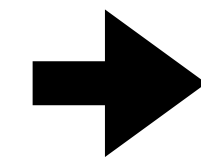
5

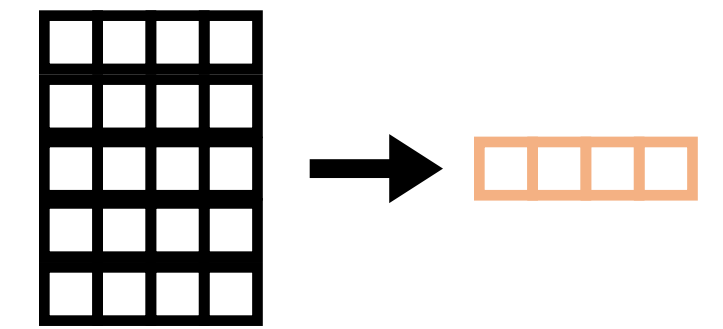# Example: PointNet++



**Furthest Sample**
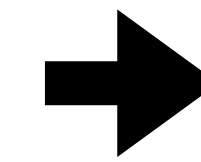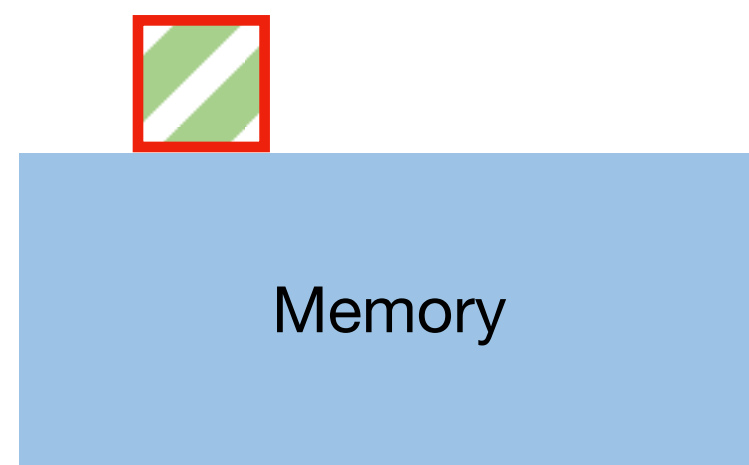Irregular memory access

**Ball Query**
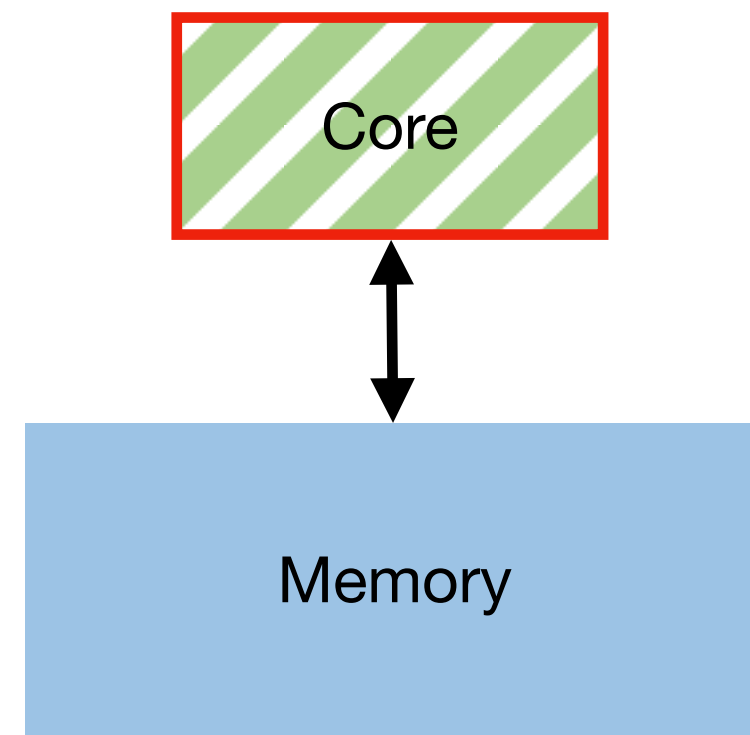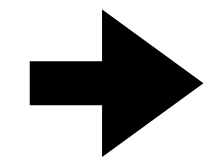Complex control

**Gather**
Irregular memory access

**Multi-layer Perceptron**
Matrix computations

**Aggregate**
Element-wise *max()*

Core

Memory

**Near-Memory**

Memory

**In-Core**

Core

Memory

**Near-Memory**

Memory

**In-Memory**

Core

Memory

**In-Core**

# There is an *Adoption Problem*

# Programming is *Difficult*



```
Mat matmul(Mat core_A, Mat core_B, Mat core_C):
    InMat in_A = inMemcpy(core_A, M * N, fromCore)
    InMat in_B = inMemcpy(core_B, M * N, fromCore)
    InMat in_C = inMalloc(M * N)
    for m in [0, M):
      for n in [0, N):
        InVec in_V = in_A[m][:] * in_B[:][c]
        in_V = inPartialReduce(+, in_V, rounds=3)
        NearVec near_V = nearMalloc(K)
        near_V = nearMemcpy(in_V, K / 2 / 2 / 2,
                            fromIn)

        NearScalar near_dotpdt = nearReduce(near_V)
        core_C[m][n] = coreMemcpy(near_dotpdt,
                            fromNear)
```

# *Hide* Hardware Details

**Core**

**Near-Memory Compute**

**In-Memory Compute**

**Memory**

```
Mat matmul(Mat A, Mat B, Mat C):
  for m in [0, M):
    for n in [0, N):
      for k in [0, K):
        C[m][n] += A[m][k] * B[k][n]
```

**Easier to program**

# Outline

**Plain C Code**

**2** Unified Abstraction

Agnostically represents
1. In-core compute
2. In-memory compute
3. Near-memory compute

**4** Partitioning
Static vs Dynamic

Portable binary without
exposing µArch details to
compiler

**3** µArch-specific
Optimizations

Specialize binary to take
advantage of µArch details

µArch Extensions

Hardware support for
compute paradigms

**1** **Compute Paradigms**

*In-Core*

Core

Memory

*In-Memory*

Memory

*Near-Memory*

Memory

# Outline

**2**

Unified Abstraction

Agnostically represents
1. In-core compute
2. In-memory compute
3. Near-memory compute

**4**

Partitioning
Static vs Dynamic

Portable binary without
exposing µArch details to
compiler

**3**

µArch-specific
Optimizations

Specialize binary to take
advantage of µArch details

µArch Extensions

Hardware support for
compute paradigms

**1**

**Compute Paradigms**

Plain C Code

*In-Core*

Core

Memory

*In-Memory*

Memory

*Near-Memory*
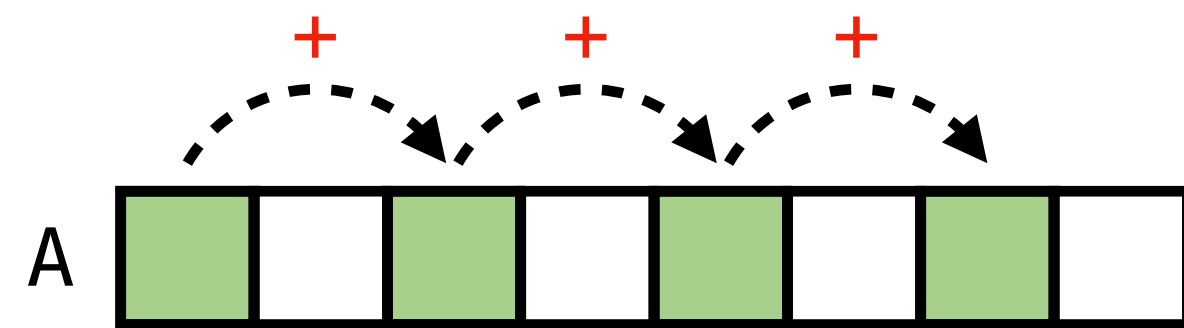
Memory

# Near-Memory Compute

*Offload computation* closer to the data

$C[i] = A[i] \& B[i]$

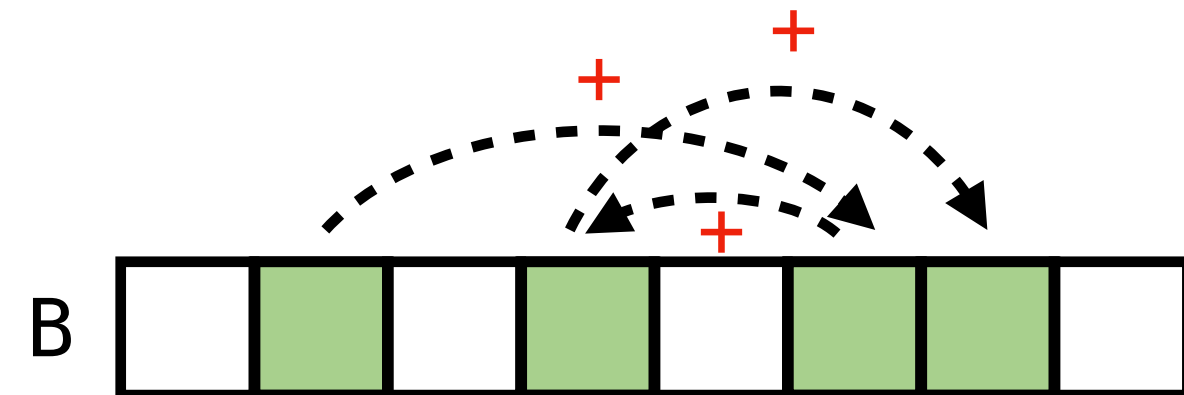Supports complex **memory access patterns**
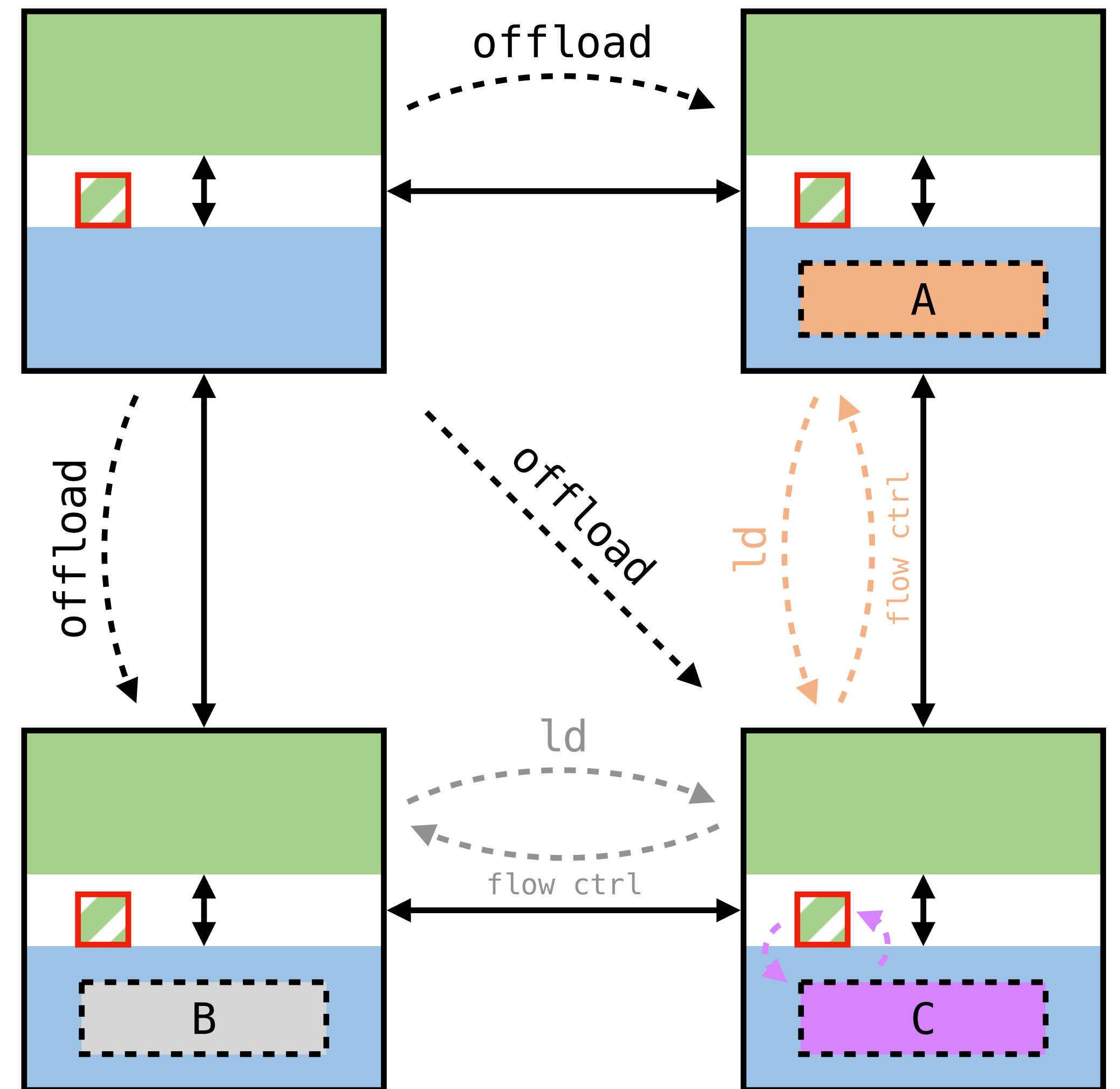
***Affine Patterns***
e.g., `A[2*i]`

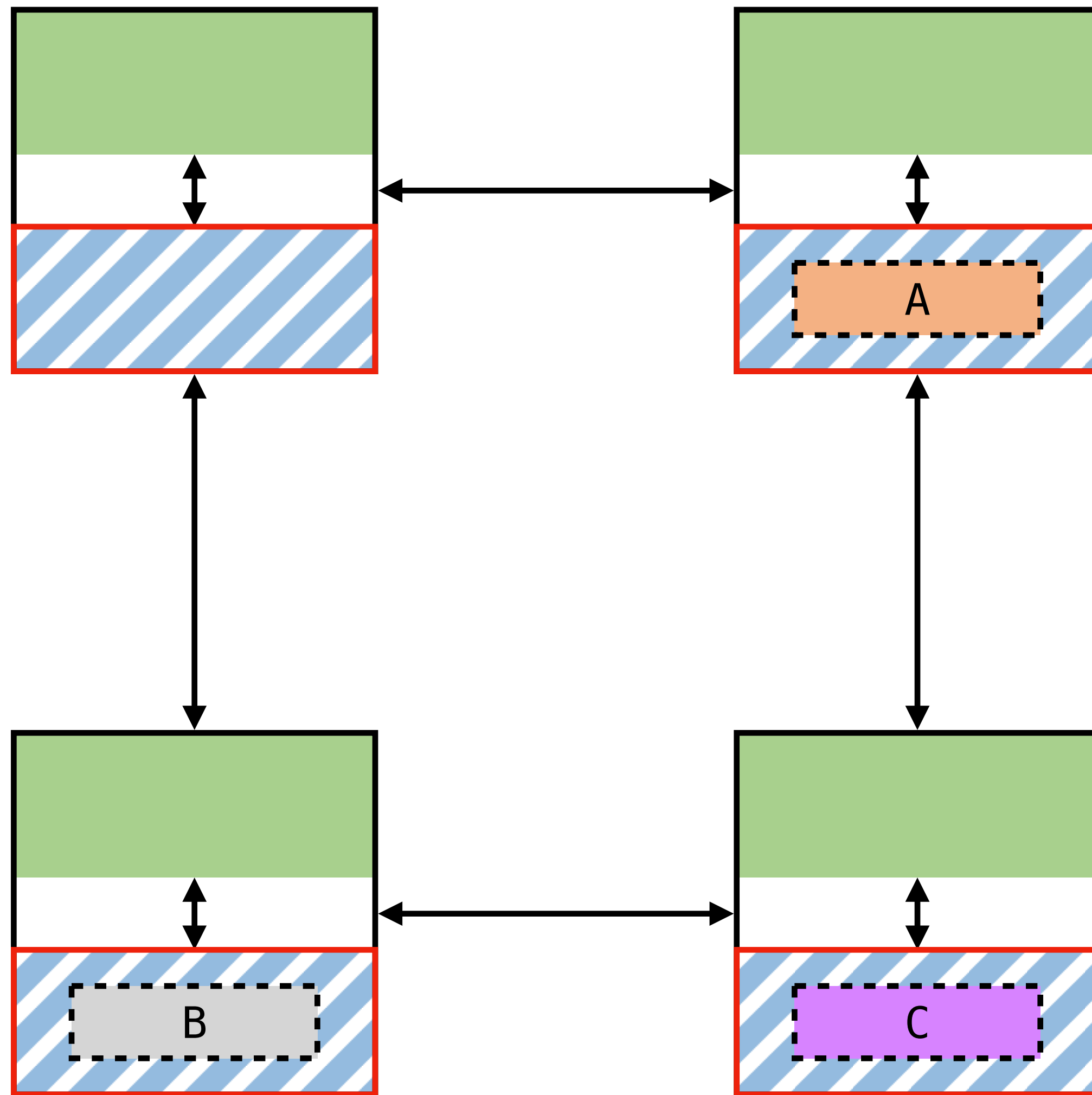***Indirect Patterns***
e.g., `B[A[i]]`

with **reduction** capabilities
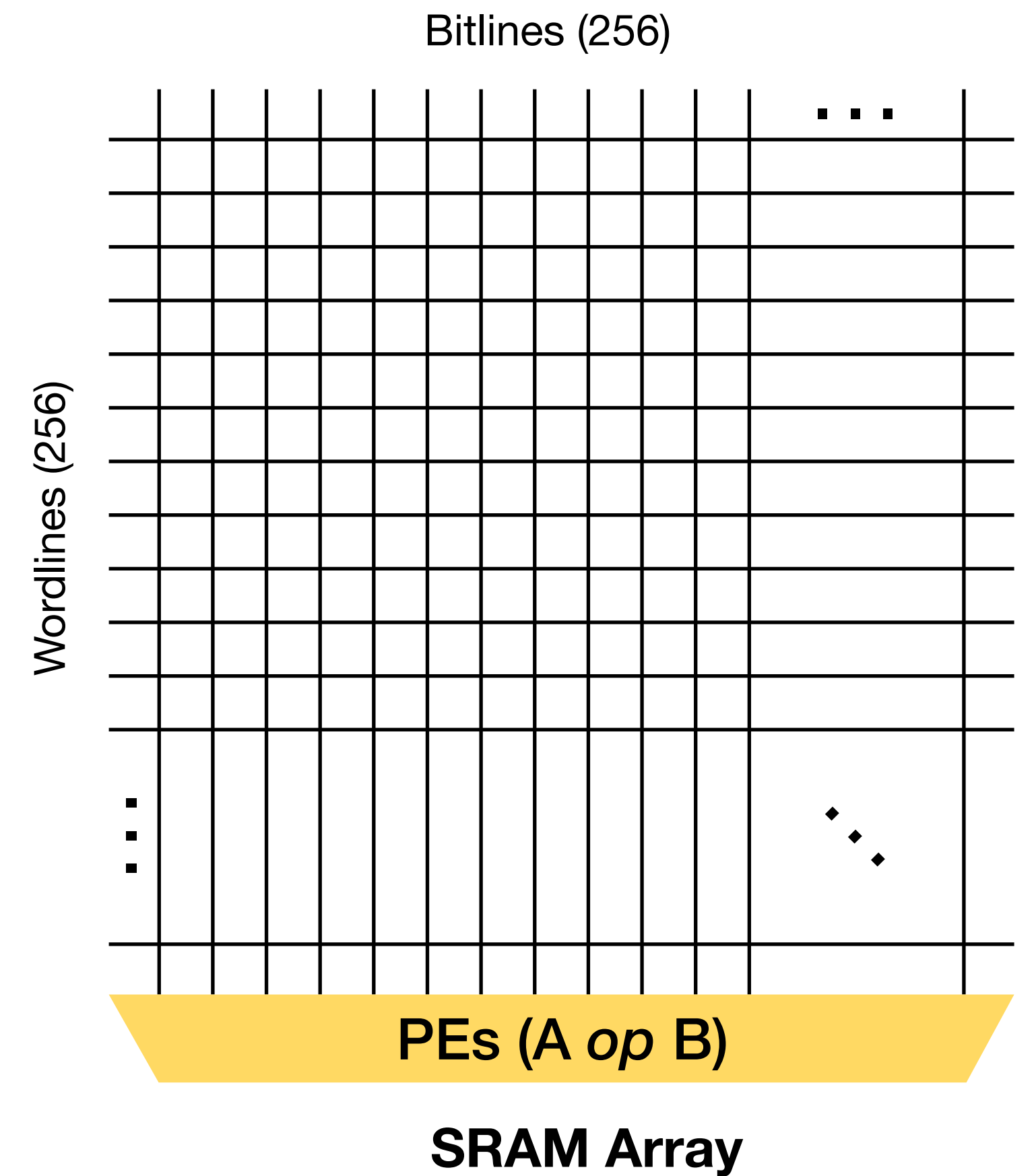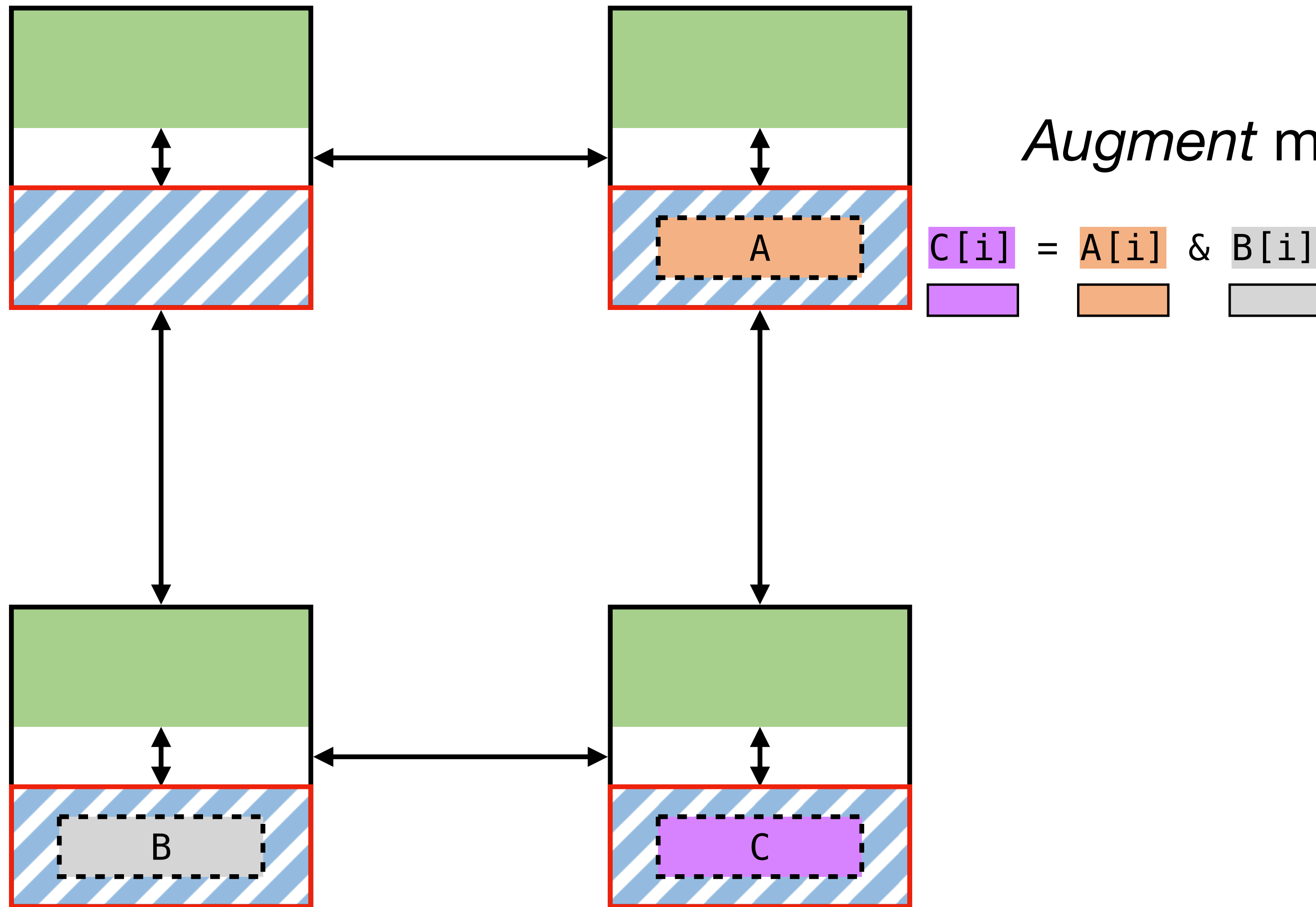
**Limitation:** Lower compute width

**In-Memory Compute**

*Augment* memory technology with compute
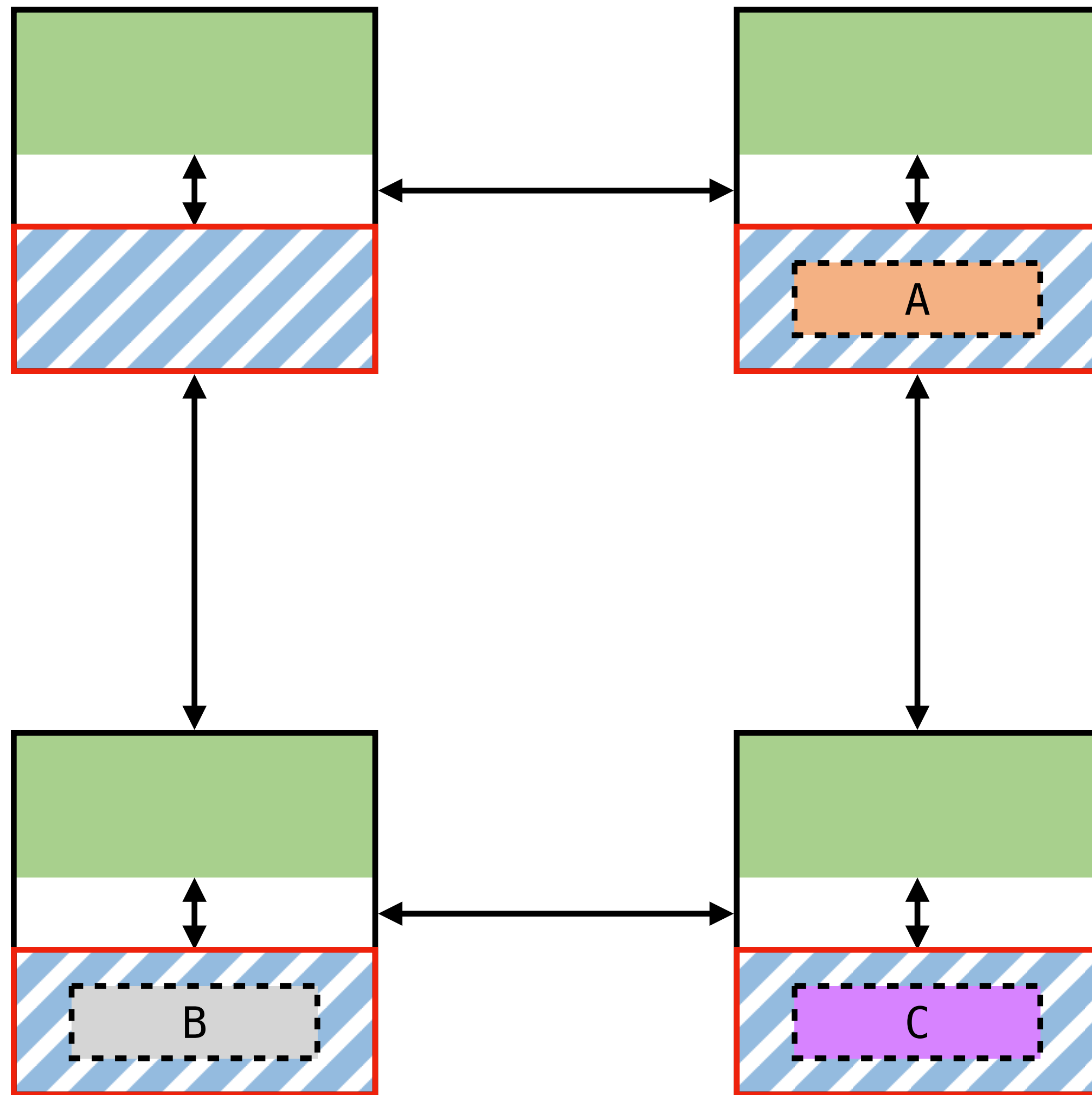
`C[i]` = `A[i]` & `B[i]`

# In-Memory Compute

*Augment* memory technology with compute

C[i] = A[i] & B[i]
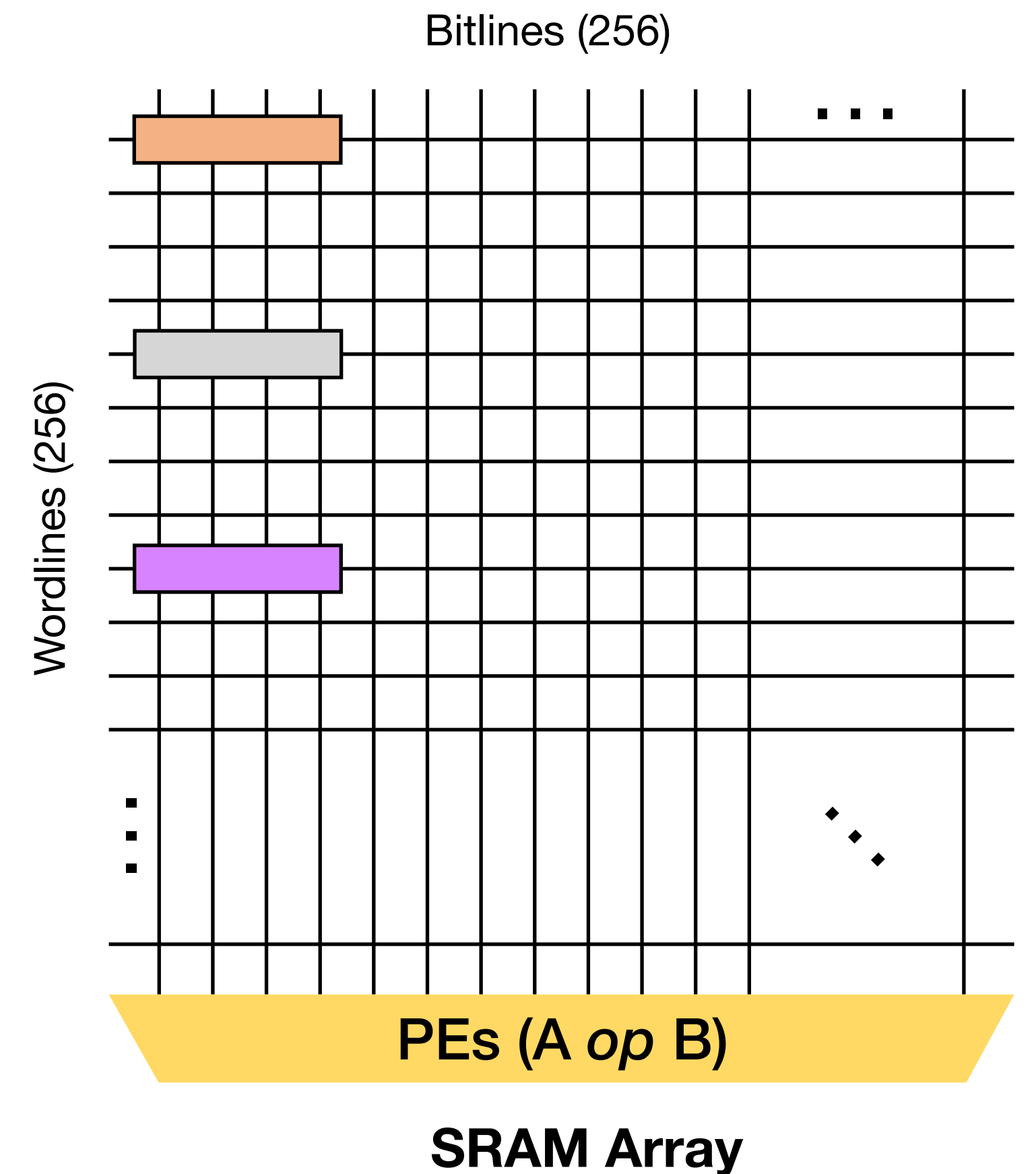
Bitlines (256)

Wordlines (256)

PEs (A *op* B)

**SRAM Array**

16

**In-Memory Compute**

*Augment* memory technology with compute

$C[i] = A[i] \ \& \ B[i]$

***Standard Data Layout***

Bitlines (256)

Wordlines (256)

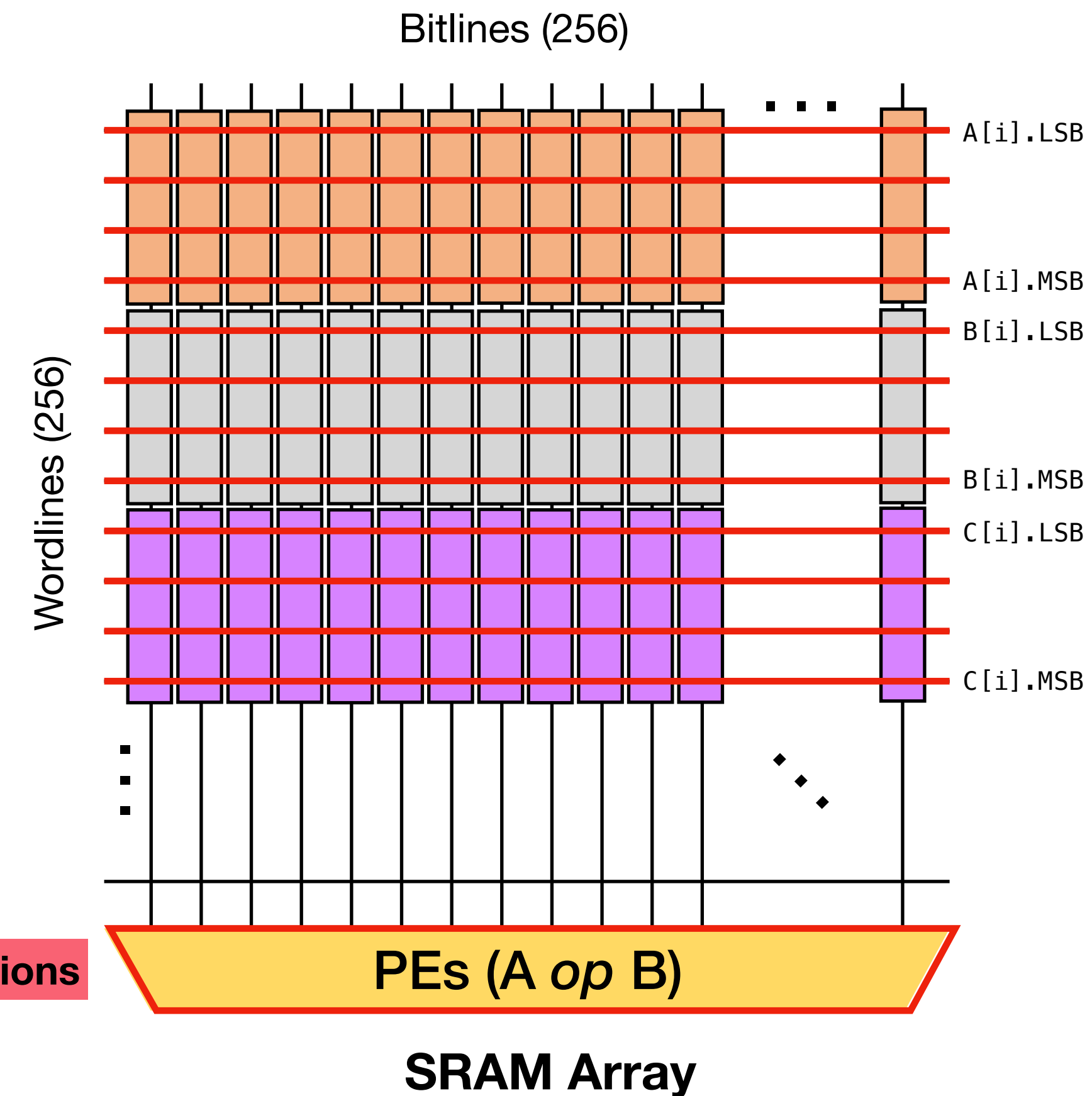PEs (A *op* B)

**SRAM Array**

A

B

C

*Augment* memory technology with compute

A

C[i] = A[i] & B[i]

Supported Operations

+ − * / % & | ^ << >>
sin cos exp log sqrt

Offers **massive** vector parallelism

B

C

*Bit-Serial Data Layout*

Bitlines (256)

. . .

A[i].LSB

A[i].MSB
B[i].LSB

B[i].MSB
C[i].LSB

C[i].MSB

Wordlines (256)

Bitwise Operations

PEs (A *op* B)

**SRAM Array**

*Augment* memory technology with compute

**Bit-Serial Data Layout**

Bitlines (256)

Supported Operations
```
+ - * / % & | ^ << >>
sin cos exp log sqrt
```

Offers **massive** vector parallelism

$C[i]$ = $A[i]$ & $B[i]$

A
B
C

A[i].LSB

A[i].MSB
B[i].LSB

B[i].MSB
C[i].LSB

C[i].MSB

Wordlines (256)

PEs (A *op* B)

**Challenge:** Requires data alignment & transpose

**SRAM Array**

20

# Compute Paradigms

## In-Memory Compute

Offers **massive** vector parallelism

**Challenge:** Requires data alignment & transpose

## Near-Memory Compute

Supports complex **memory access patterns** with **reduction** capabilities

**Limitation:** Lower compute width

# Programming Considerations

## 1. Orchestration

Meet the requirements for each compute paradigm

- Data alignment
- Data layout & tiling
- Bit-serial transpose
- Managing on-chip space

## 2. Fused In-/Near-Memory

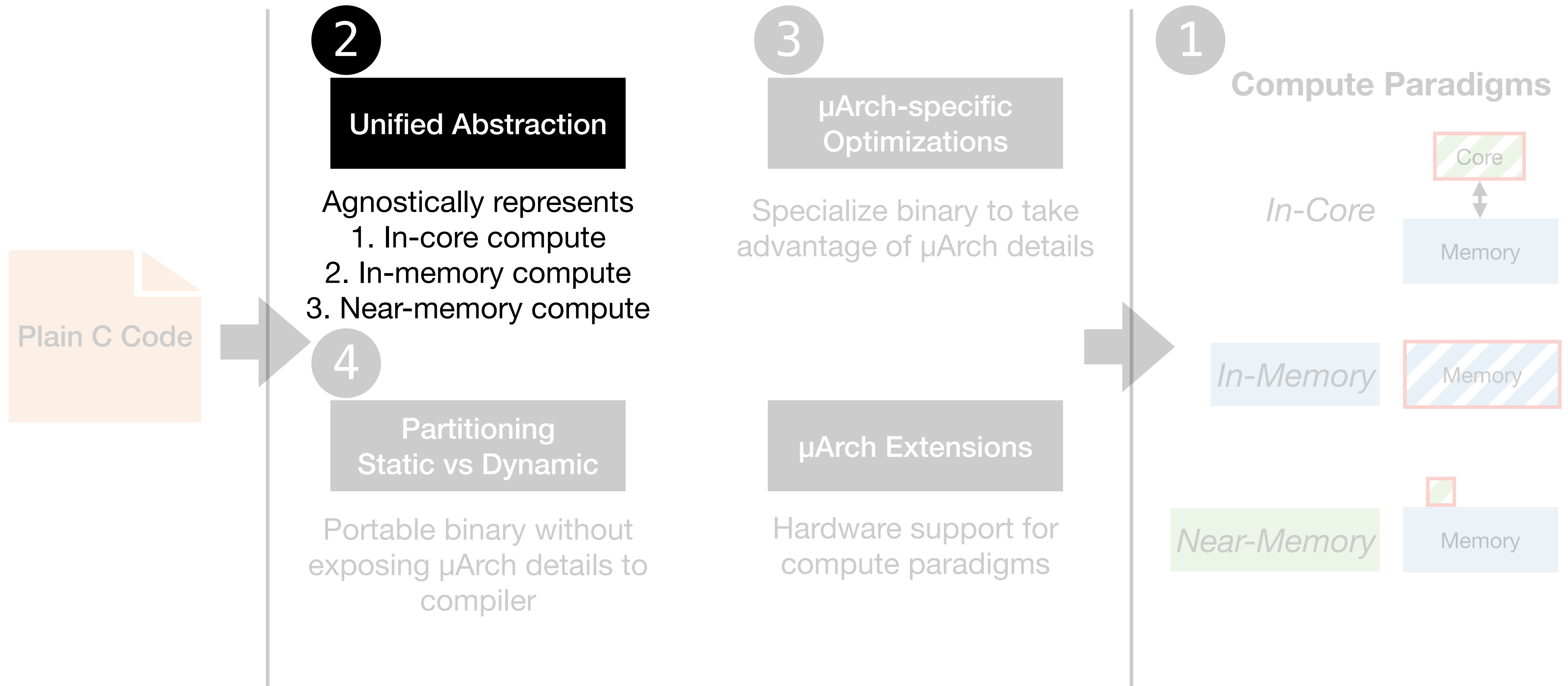Statically & dynamically take advantage of each compute paradigm

- **In-Memory:** Large input size & element-wise compute
- **Near-Memory:** Small input size or irregular memory patterns
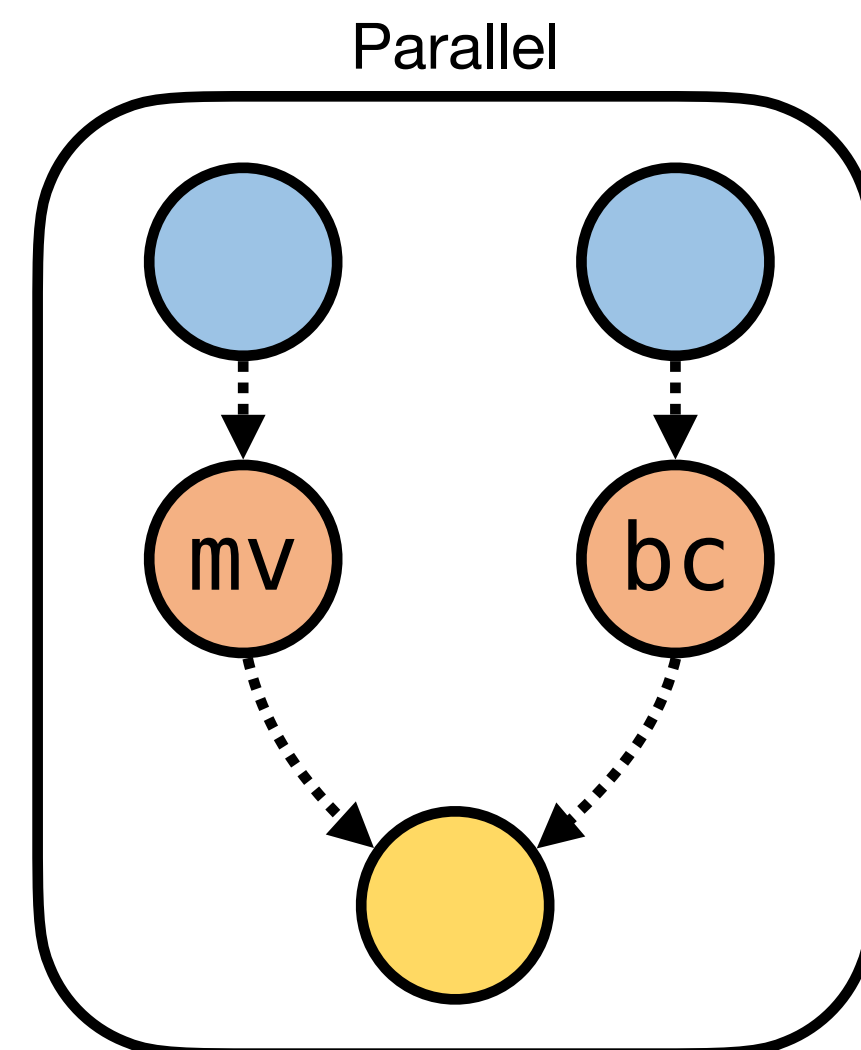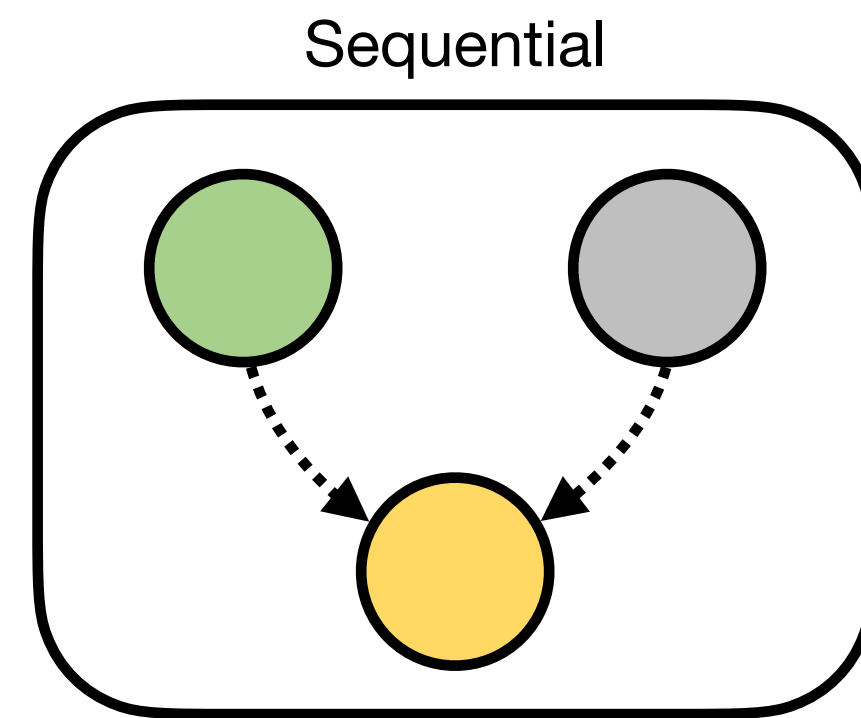- **Fusion** of in-/near-memory

## 3. Portability

Target a large variety of microarchitectures with a single binary

# Programming Considerations

## 1. Orchestration

Meet the requirements for each compute paradigm

- Data alignment
- Data layout & tiling
- Bit-serial transpose
- Managing on-chip space

## 2. Fused In-/Near-Memory

Statically & dynamically take advantage of each compute paradigm

- **In-Memory:** Large input size & element-wise compute
- **Near-Memory:** Small input size or irregular memory patterns
- **Fusion** of in-/near-memory

## 3. Portability

Target a large variety of microarchitectures with a single binary

# Transparency
minimizes programmer burden

# Outline

**2**

**Unified Abstraction**

Agnostically represents
1. In-core compute
2. In-memory compute
3. Near-memory compute

**4**

**Partitioning
Static vs Dynamic**

Portable binary without
exposing µArch details to
compiler

**3**

**µArch-specific
Optimizations**

Specialize binary to take
advantage of µArch details

**µArch Extensions**

Hardware support for
compute paradigms

**1**

**Compute Paradigms**

Core

*In-Core*

Memory

*In-Memory*

Memory

Memory

*Near-Memory*

Memory

Plain C Code

# Dataflow Representation



**Near-Memory Compute**

- Memory access pattern
- Computation

**In-Memory Compute**

- N-D Tensors
- Data alignment
- Spatial reuse

**Fused In-/Near-Memory**

- Reduction operations
- Fusion

# Stream Dataflow Graph

○ Memory access pattern

**1D Filter**

```
for i in [0, N–1):
    B[i] = F[0] x A[i  ]
         + F[1] x A[i+1]
```

**Memory Access Pattern**

| a | b | c | d | e | f | g | h | encoded by **Stream Node** $A_0$

| a | b | c | d | e | f | g | h | encoded by **Stream Node** $A_1$

# Stream Dataflow Graph

○ Computation

**1D Filter**

```
for i in [0, N-1):
  B[i] = F[0] x A[i  ]

       + F[1] x A[i+1]
```

# Tensor Dataflow Graph

**Observation:** contiguous memory access may be tensorized

○      N-D Tensors

**1D Filter**

```
for i in [0, N–1):
  B[i] = F[0] x A[i  ]

       + F[1] x A[i+1]
```



**But** we need to map *data* to *hardware*

30

# Tensor Dataflow Graph

**1D Filter**

```
for i in [0, N–1):
  B[i] = F[0] x A[i  ]

       + F[1] x A[i+1]
```

**Observation:** contiguous memory access may be tensorized



Maps
*data → virtual vector lane*

**Global Lattice Space**

Each cell represents a virtual vector lane

Describes N-dimensional alignment

Maps
*virtual vector lane →
physical vector lane*

lane 1 lane 2 lane 3 lane 4 lane 5 lane 6 lane 7 lane 8

**Hardware Vector Processor**

# Tensor Dataflow Graph

## 1D Filter

```
for i in [0, N-1):
  B[i] = F[0] x A[i  ]

       + F[1] x A[i+1]
```

## Global Lattice Space

Each cell represents a virtual vector lane

Describes N-dimensional alignment

Memory access is represented as a hyperrectangle

# Tensor Dataflow Graph

**1D Filter**

```
for i in [0, N-1):
    B[i] = F[0] x A[i  ]

         + F[1] x A[i+1]
```

Memory access is represented as a hyperrectangle

○ Data alignment

**Move Node** $mv$ realigns hyperrectangle

36

# Matrix-Vector Multiplication

```
for m in [0, M):
    for k in [0, K):
        C[m] += A[m][k] * B[k]
```

## Example



A          B          C

# Matrix-Vector Multiplication

```
for m in [0, M):
  C[m] = sum(A[m][:] * B[k])
```
Vectorize

○ Reduction operations

**Example**



A       B       C

40

# Matrix-Vector Multiplication

```
for m in [0, M):
    C[m] = sum(A[m][:] * B[k])
```

Spatial

**Example**



A          B          C

# Matrix-Vector Multiplication

```
for m in [0, M):
    C[m] = sum(A[m][:] * B[k])
```

Spatial



○ Reduction operations

**Example**



A · B = C

43

# Matrix-Vector Multiplication

```
for m in [0, M):
    C[m] = sum(A[m][:] * B[k])
```

Spatial



X

+ reduce dim=0

○ Reduction operations

**Example**



A            B            C

44

# Matrix-Vector Multiplication

```
for m in [0, M):          Spatial
    C[m] = sum(A[m][:] * B[k])
```

○ Reduction operations

## Reduction Tree          **Step**

$A_m$ x B

0

1

+

**Wasted Parallelism**          2

+ + +

$A_m$          B

X

+  reduce dim=0

V

+  reduce

Insufficient parallelism

Use near-memory to perform reduction

47

# Matrix-Vector Multiplication

Spatial
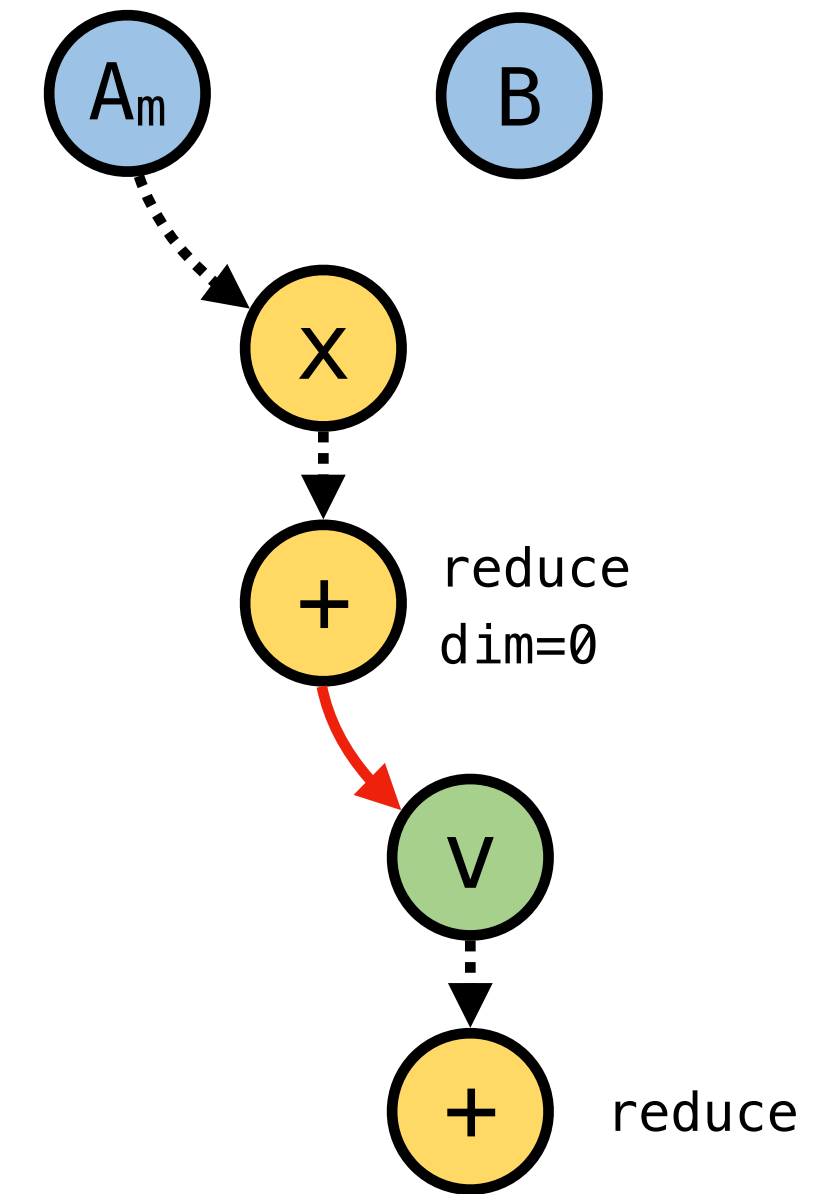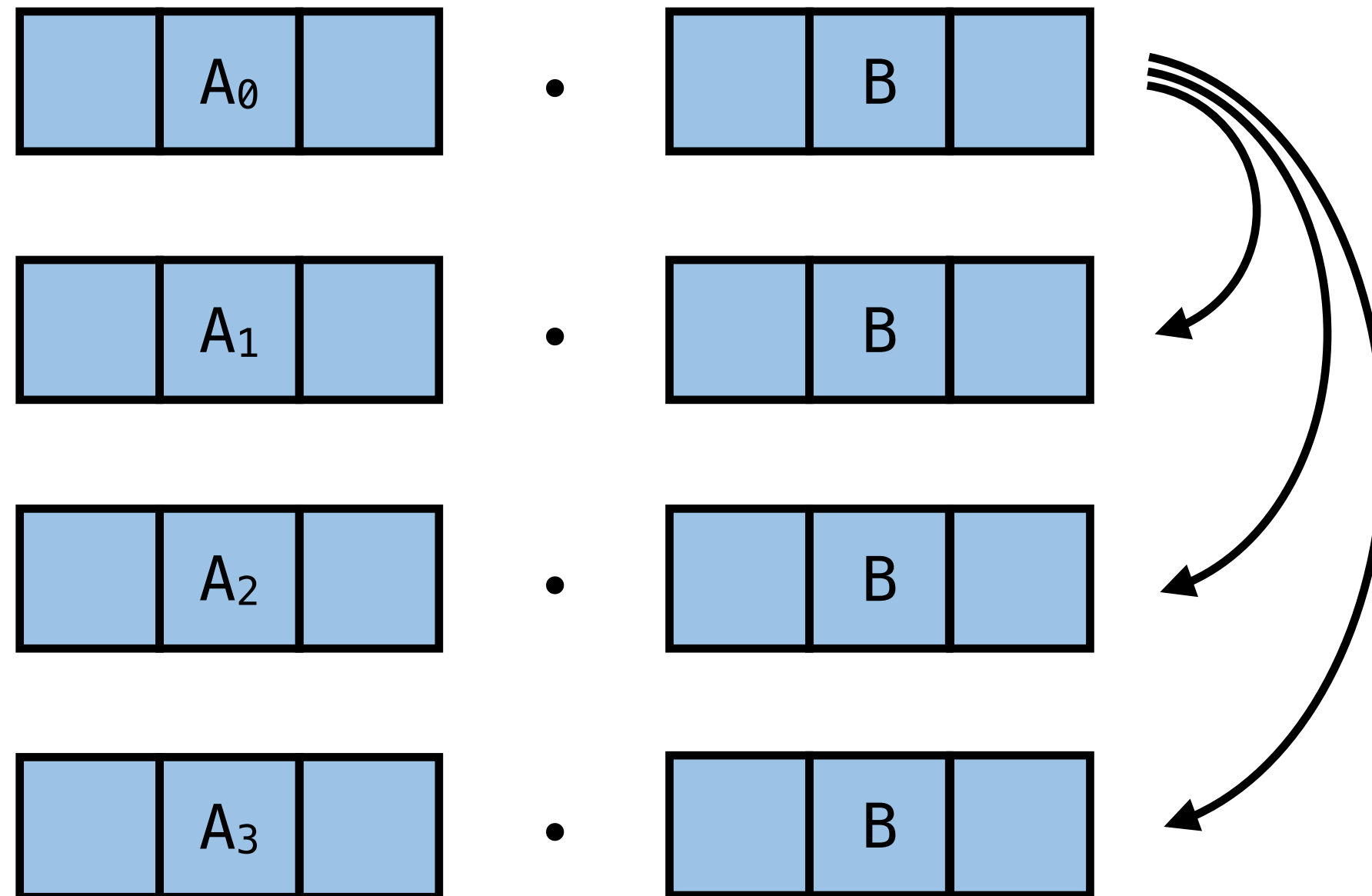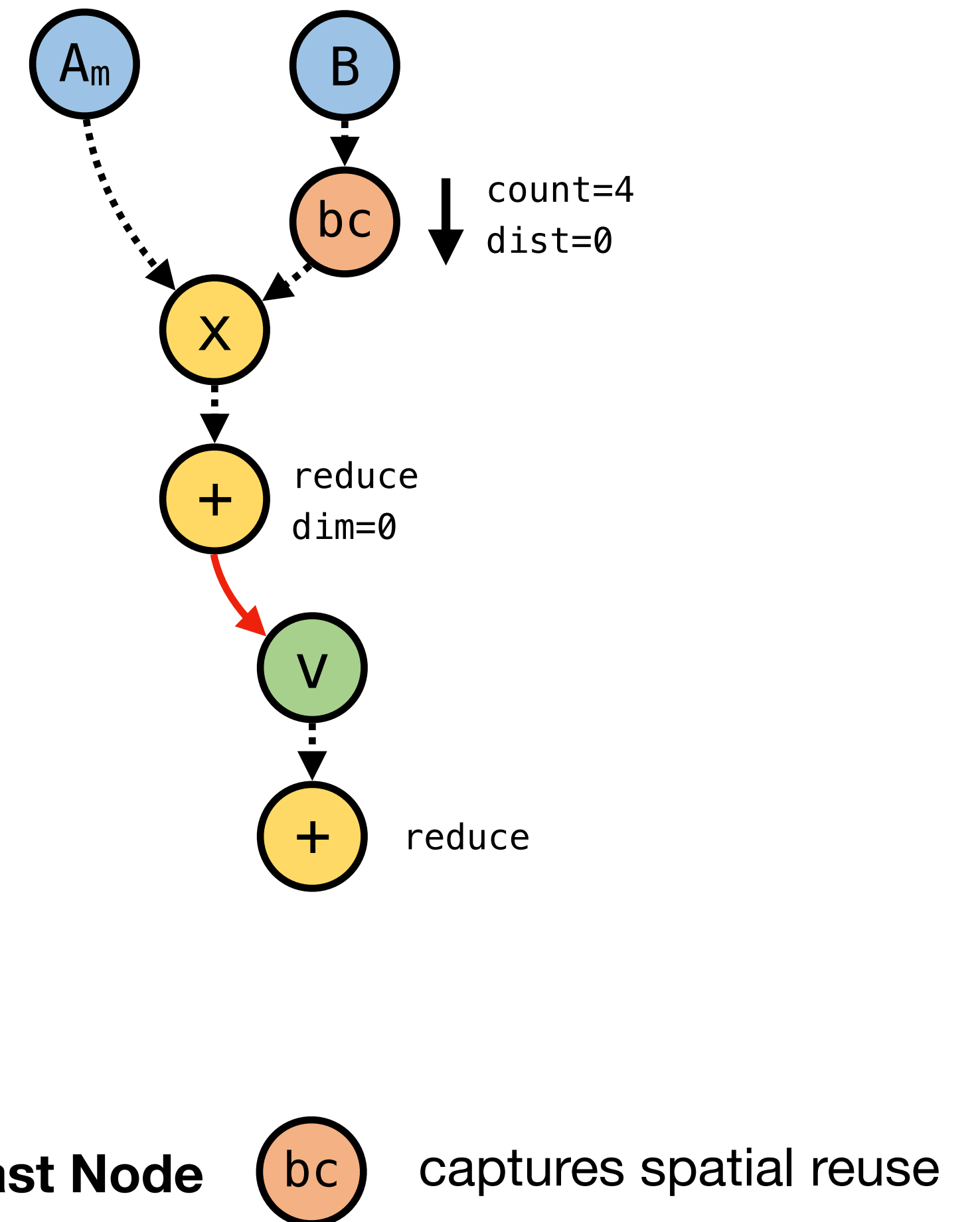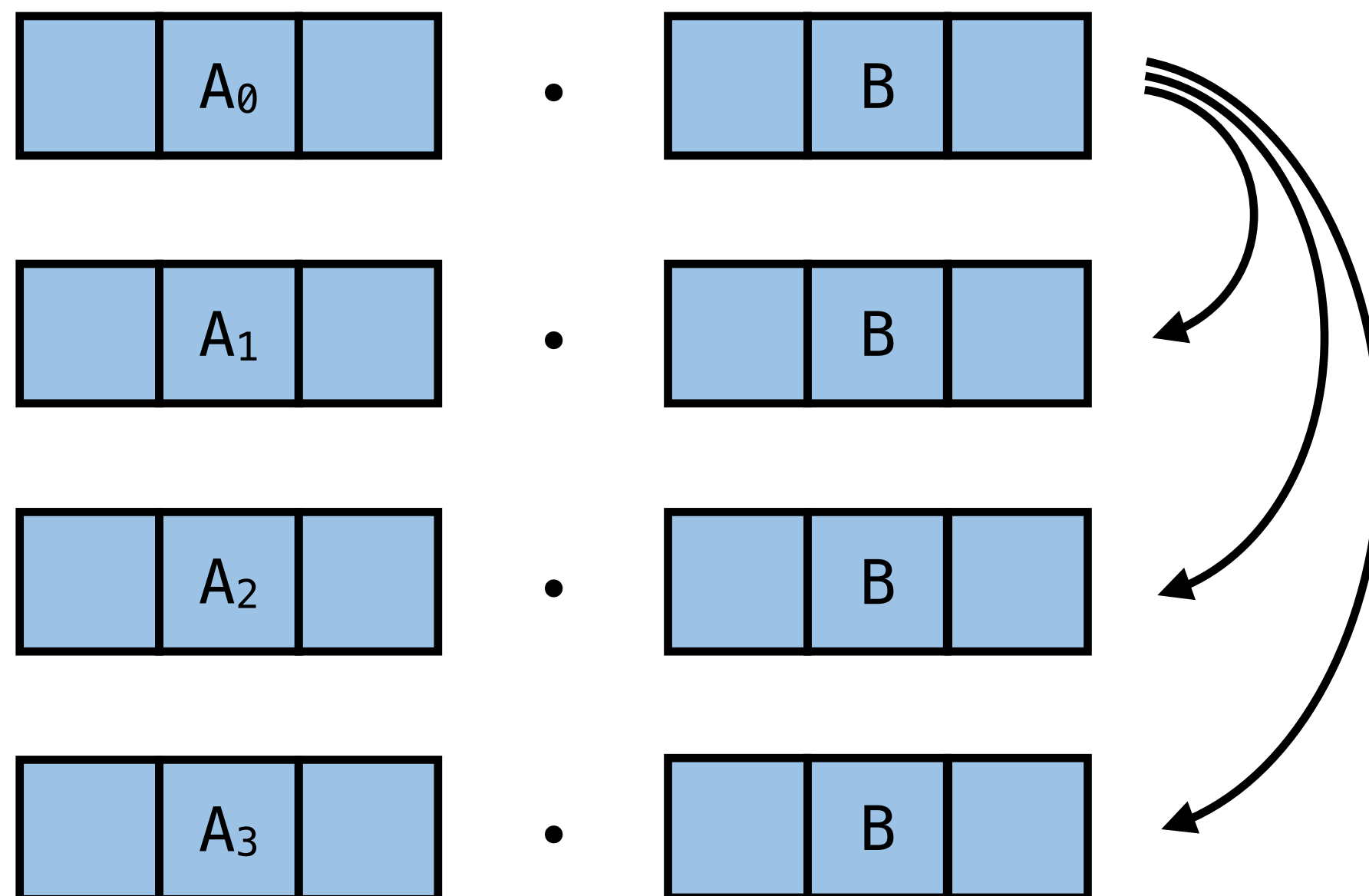
```
for m in [0, M):
    C[m] = sum(A[m][:] * B[k])
```

**Idea:** Expose more parallelism through spatial reuse

○ Spatial reuse

**Example**



$A_m$     B

X

+  reduce dim=0

V

+  reduce

50

# Matrix-Vector Multiplication

Spatial

```
for m in [0, M):
    C[m] = sum(A[m][:] * B[k])
```

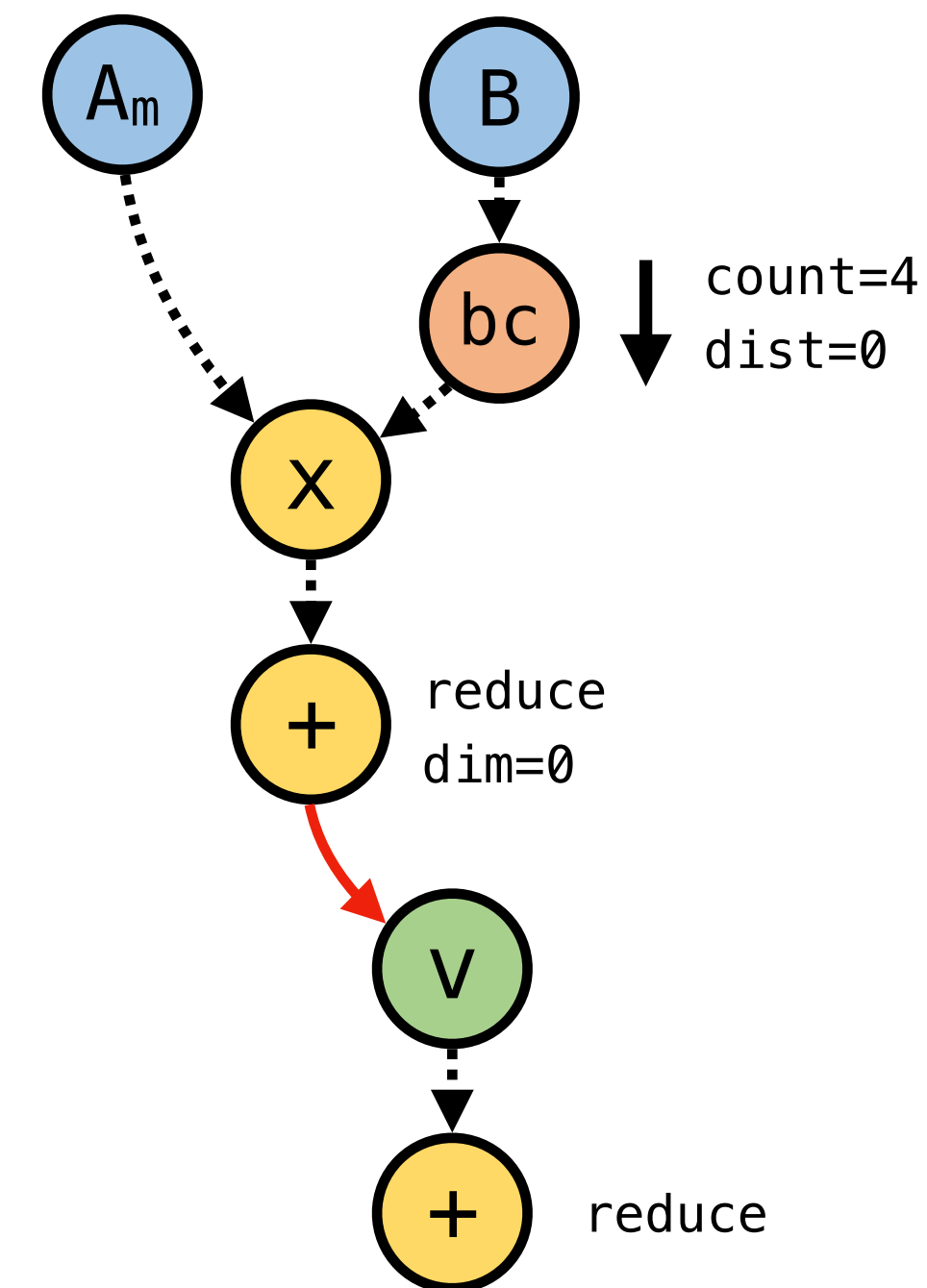**Idea:** Expose more parallelism through spatial reuse
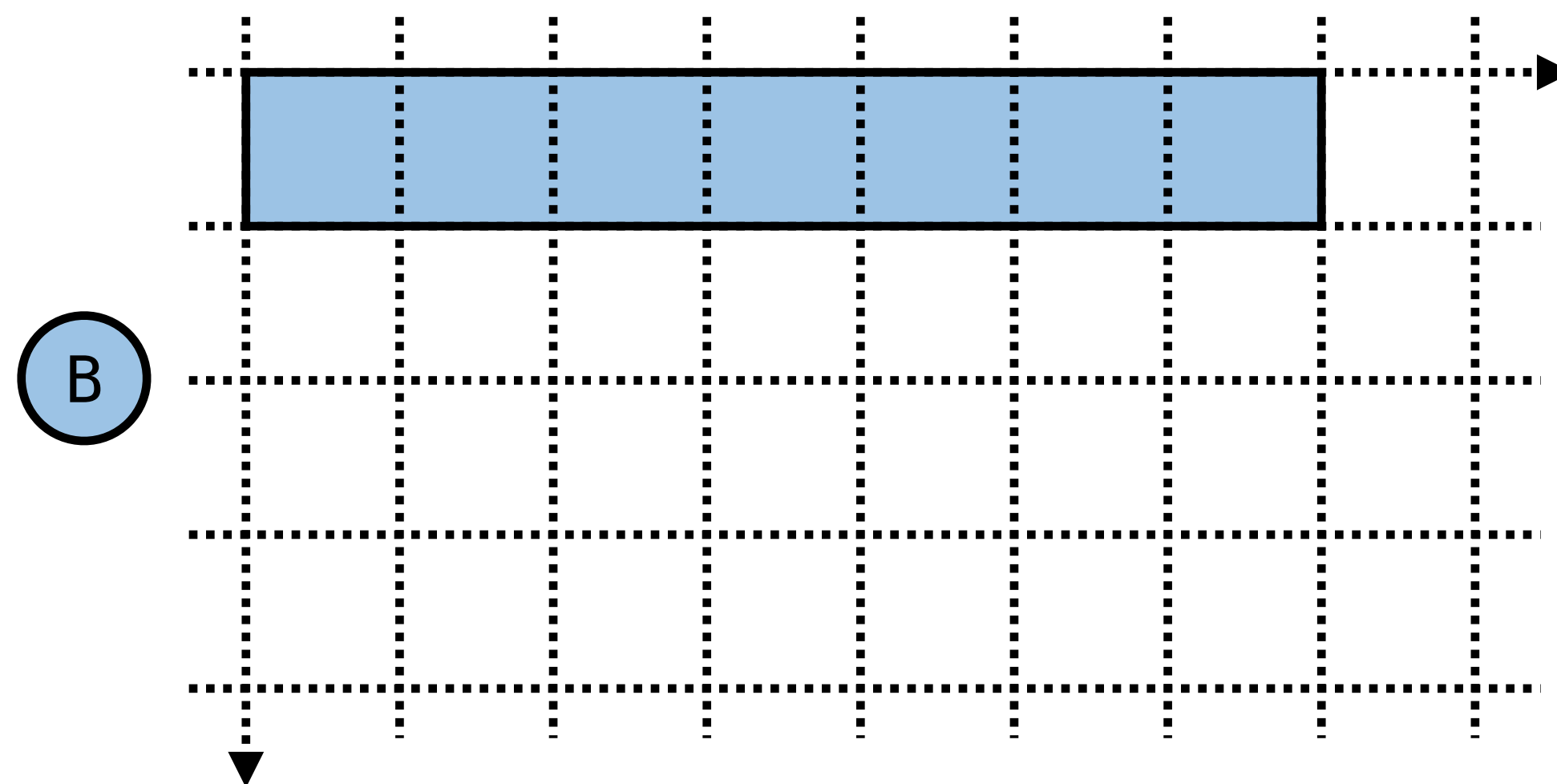
**Example**

Duplicate

52

# Matrix-Vector Multiplication

Spatial

```
for m in [0, M):
    C[m] = sum(A[m][:] * B[k])
```

**Idea:** Expose more parallelism through spatial reuse

**Example**



○ Spatial reuse

count=4
dist=0

reduce
dim=0

reduce

Duplicate

**Broadcast Node**  bc  captures spatial reuse

55

# Matrix-Vector Multiplication

Spatial

```
for m in [0, M):
    C[m] = sum(A[m][:] * B[k])
```

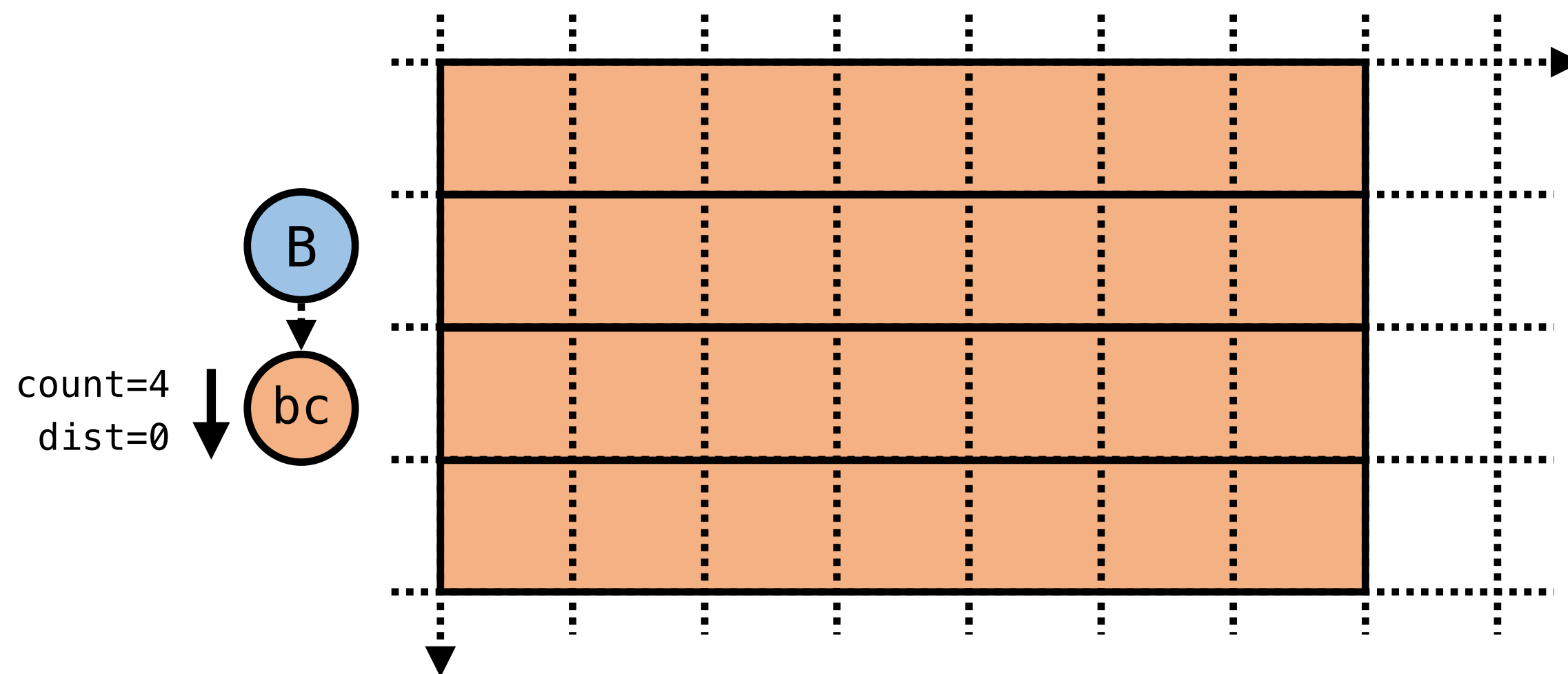**Idea:** Expose more parallelism through spatial reuse

○ Spatial reuse

A_m    B

bc    count=4
      dist=0

x

+    reduce
     dim=0

v

+    reduce

**Broadcast Node**    bc    captures spatial reuse

56

# Matrix-Vector Multiplication

Spatial

```
for m in [0, M):
    C[m] = sum(A[m][:] * B[k])
```

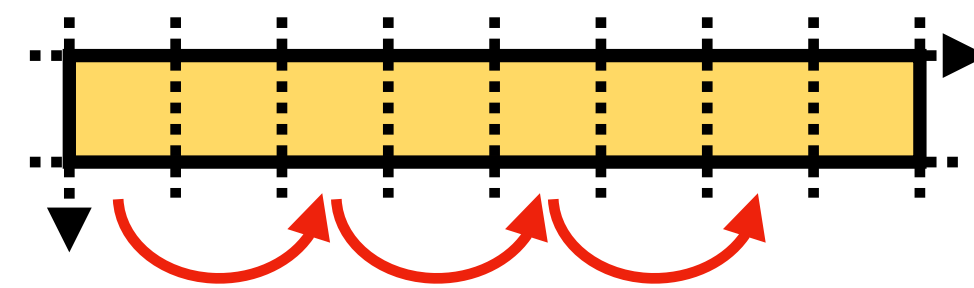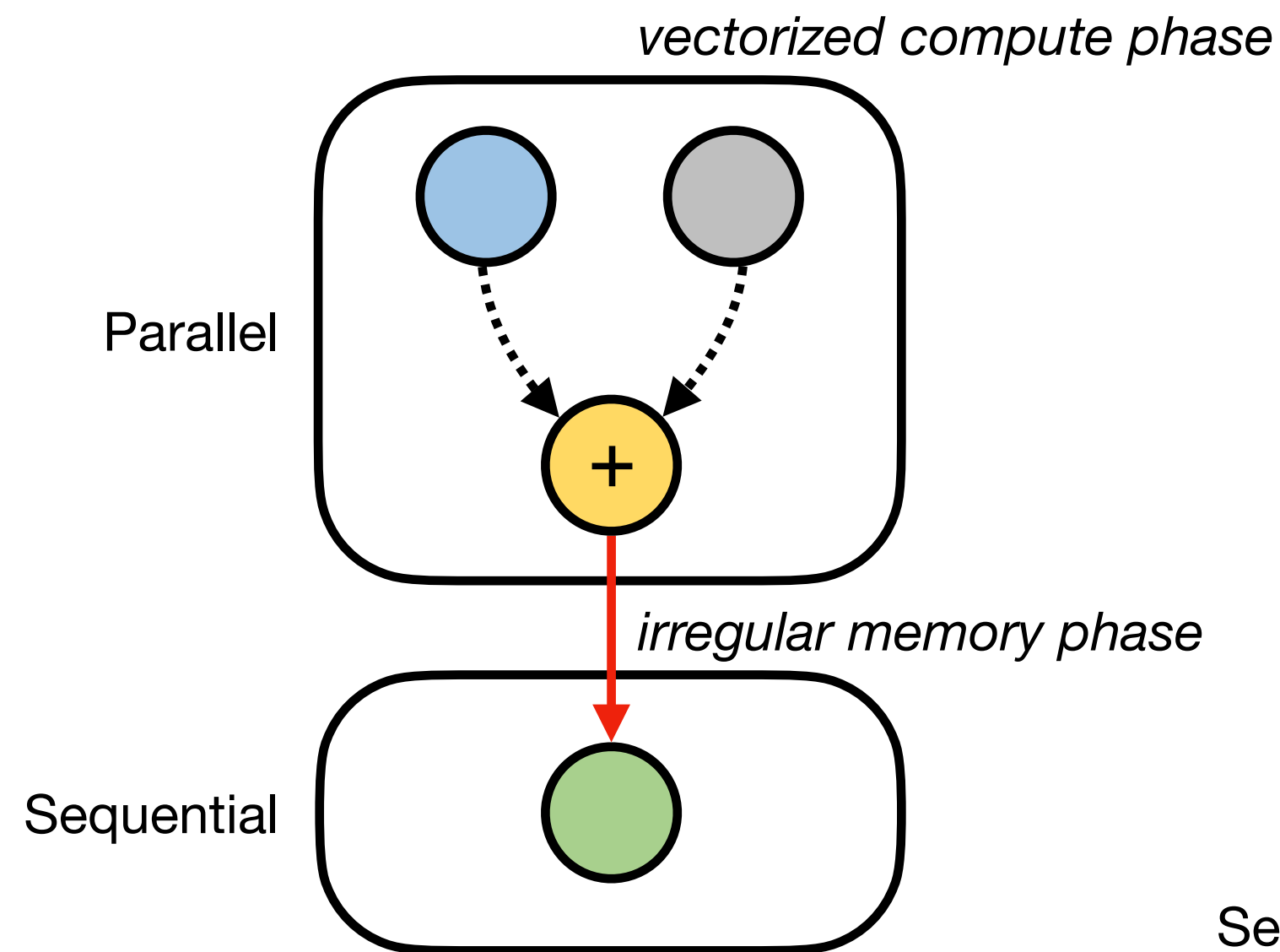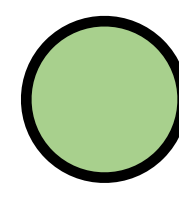**Idea:** Expose more parallelism through spatial reuse

○ Spatial reuse



**Broadcast Node** (bc) captures spatial reuse

# Stream ⬅➡ Tensor

## Load as Stream

*vectorized compute phase*

Parallel

*irregular memory phase*

Fusion

Sequential

Accesses *tensor data* with any memory access pattern supported by streams
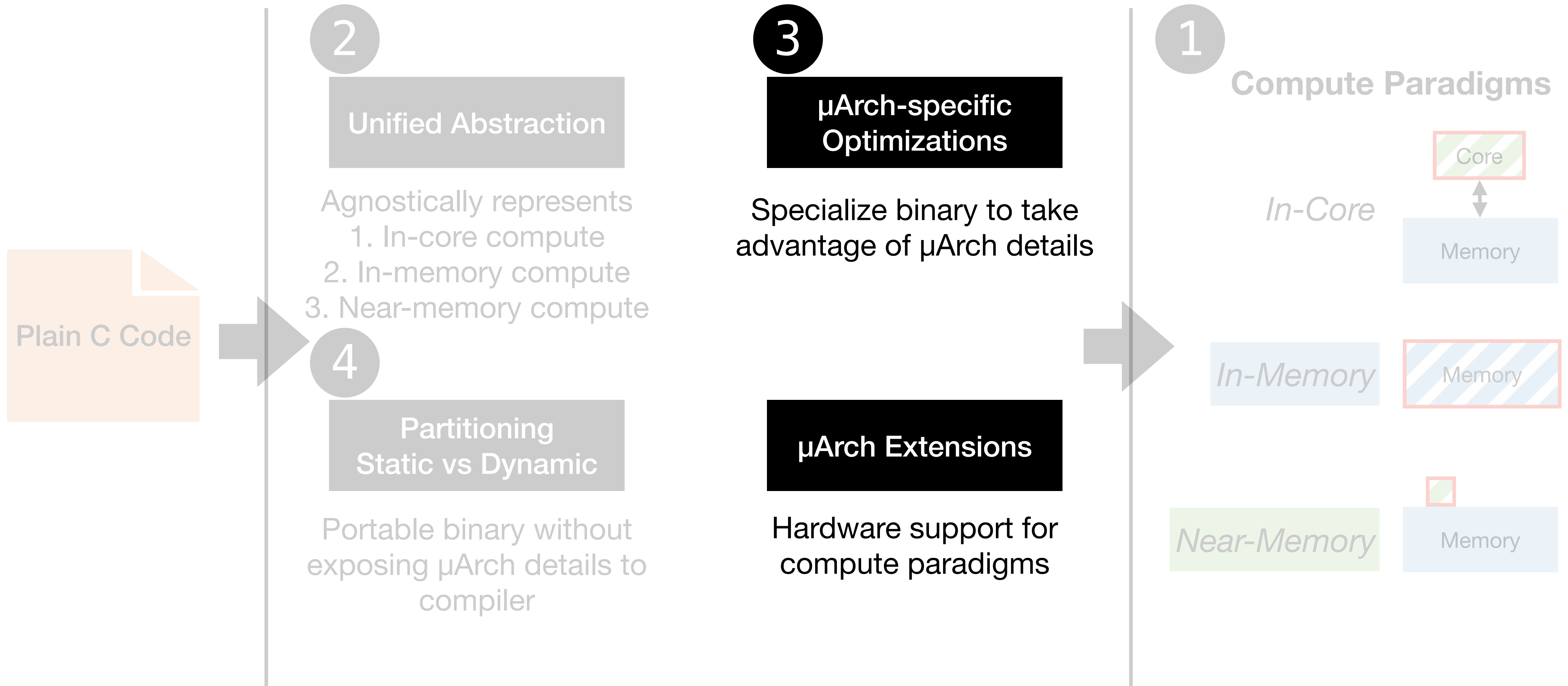
## Reduction Stream

Parallel

reduce

Sequential

reduce

Special case of
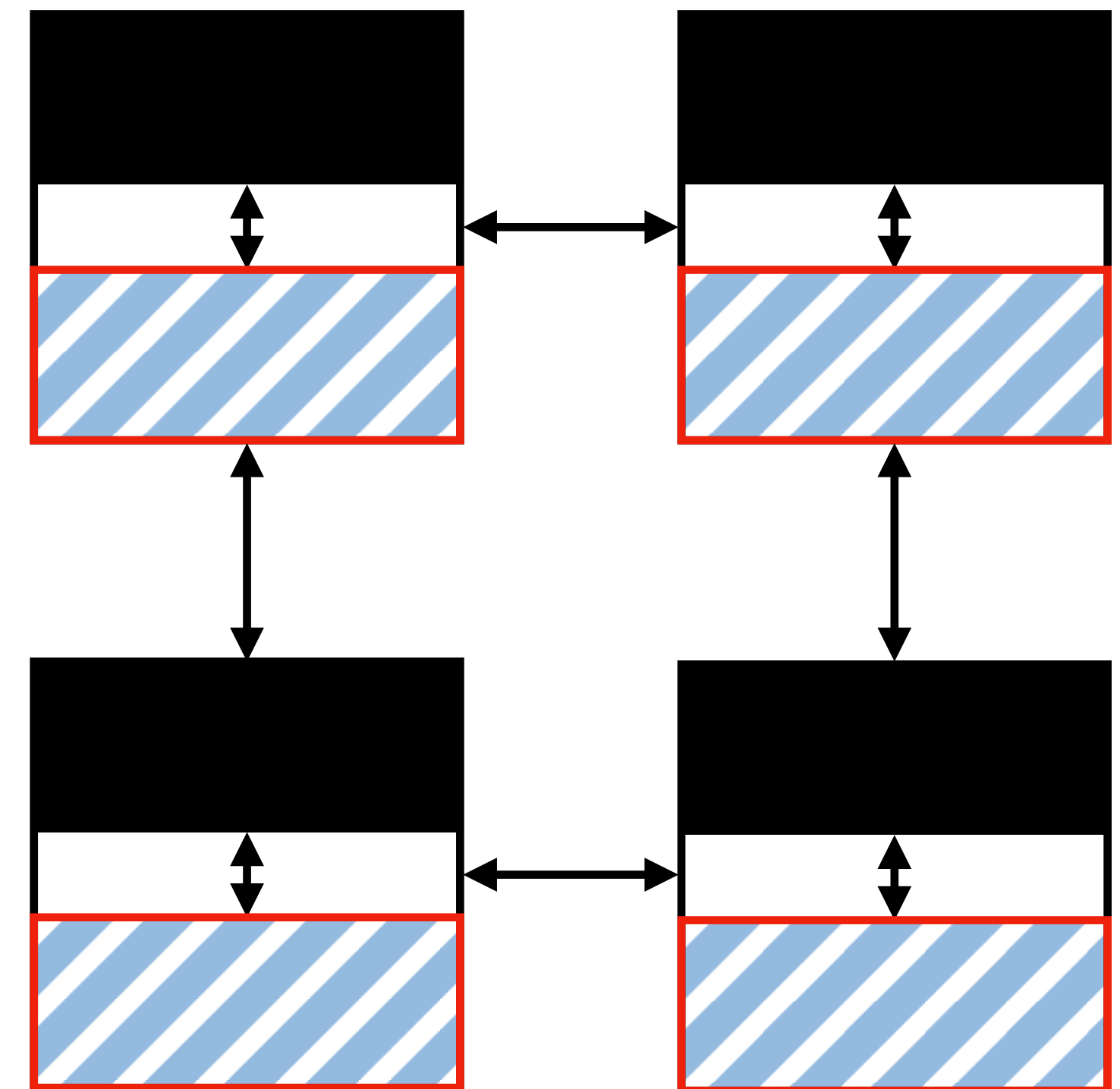*Load as Stream*

## Store as Tensor

*irregular memory phase*

Sequential

*vectorized compute phase*

Parallel

Irregular memory access

*Set up*

Contiguous memory

58

# Outline

**2**

Unified Abstraction

Agnostically represents
1. In-core compute
2. In-memory compute
3. Near-memory compute

**4**

Partitioning
Static vs Dynamic

Portable binary without
exposing µArch details to
compiler

**3**

µArch-specific
Optimizations

Specialize binary to take
advantage of µArch details

µArch Extensions

Hardware support for
compute paradigms

**1**

Compute Paradigms

Core

*In-Core*

Memory

*In-Memory*     Memory

*Near-Memory*     Memory

Plain C Code

# Mapping to Hardware
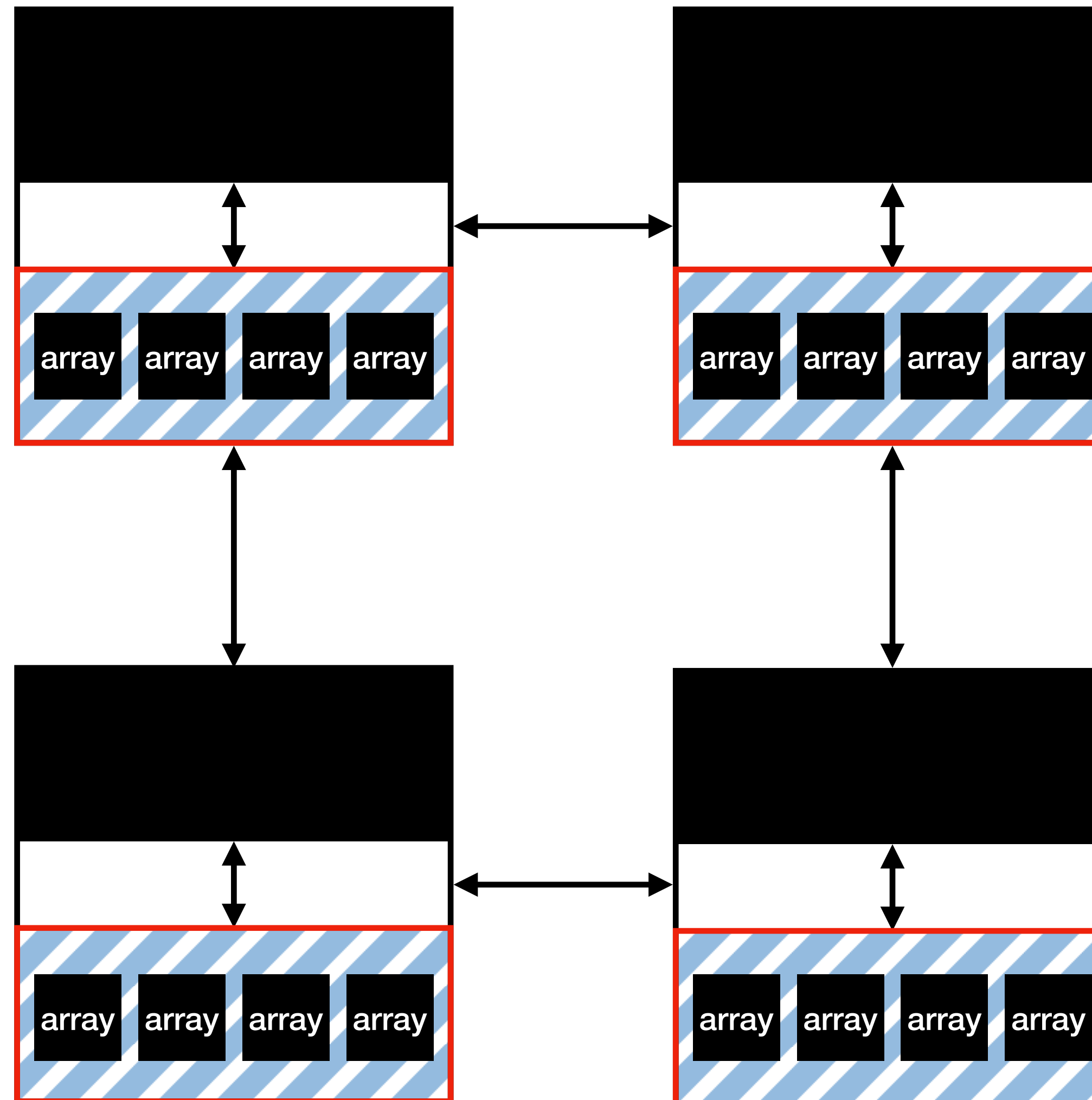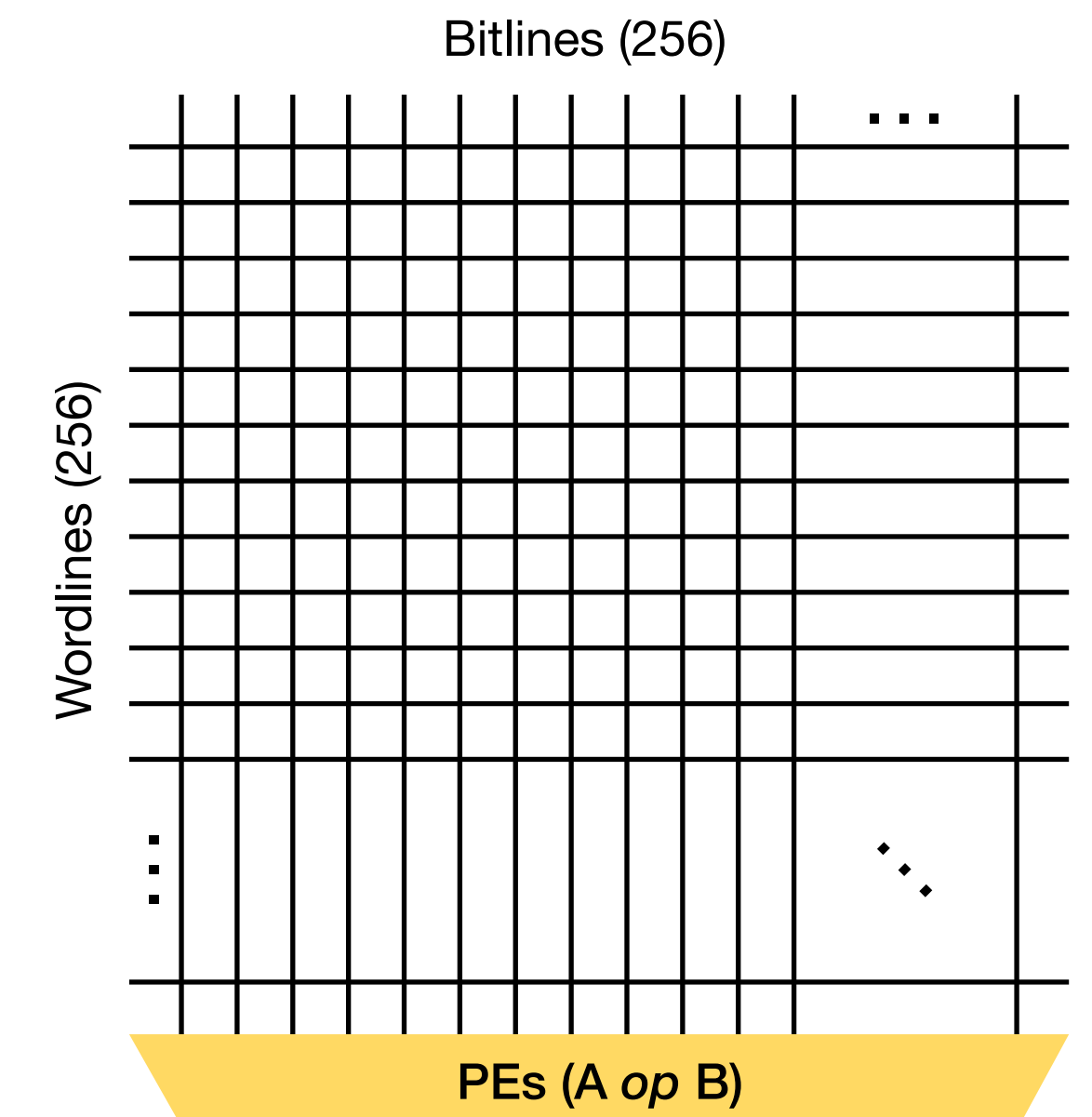
**Global Lattice**

**SRAM Banks**

# In-Memory Organization

**Question**
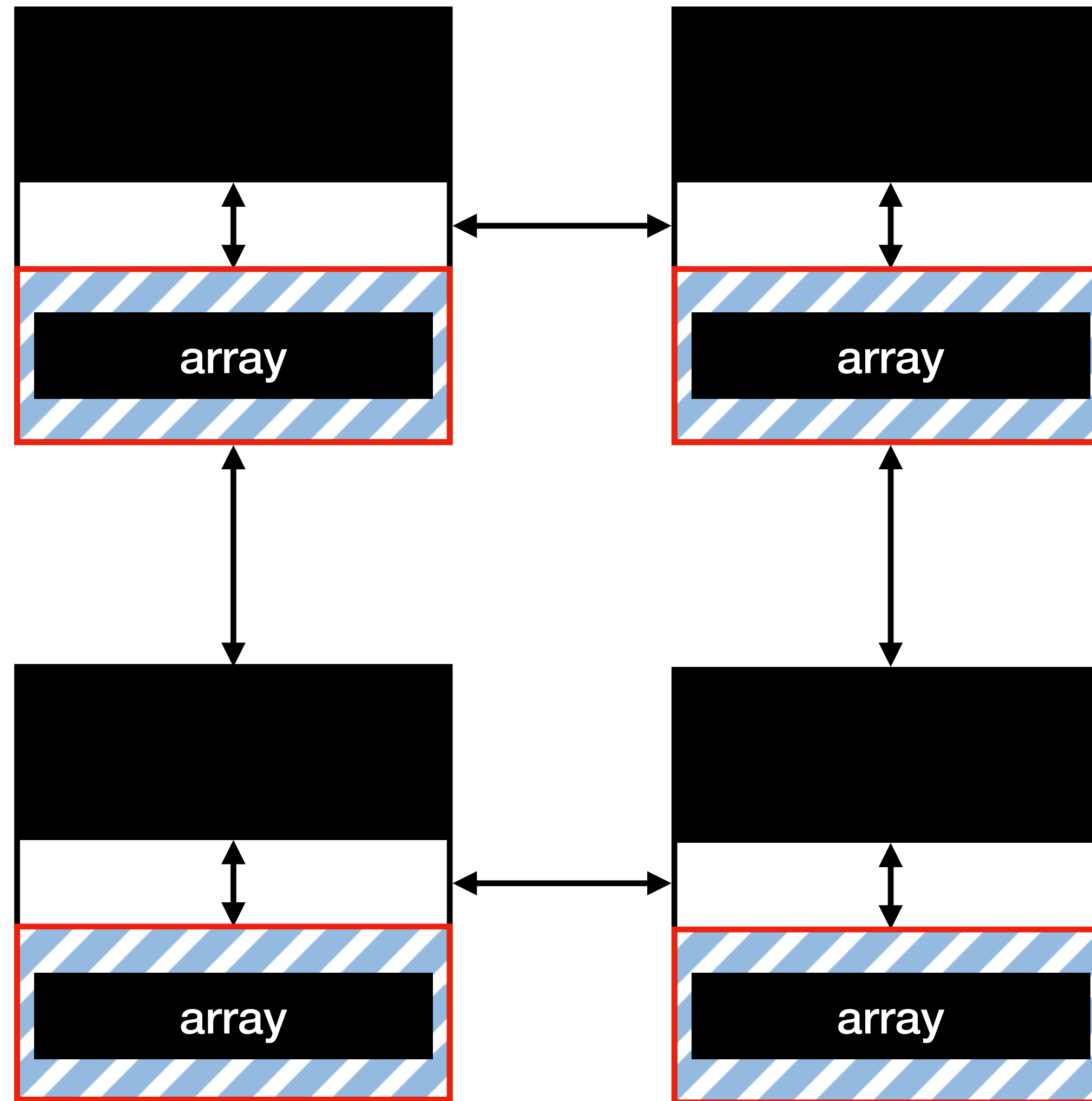What are the effects of tiling on on inter-bank traffic?

Each last-level cache bank consists of many in-memory SRAM arrays
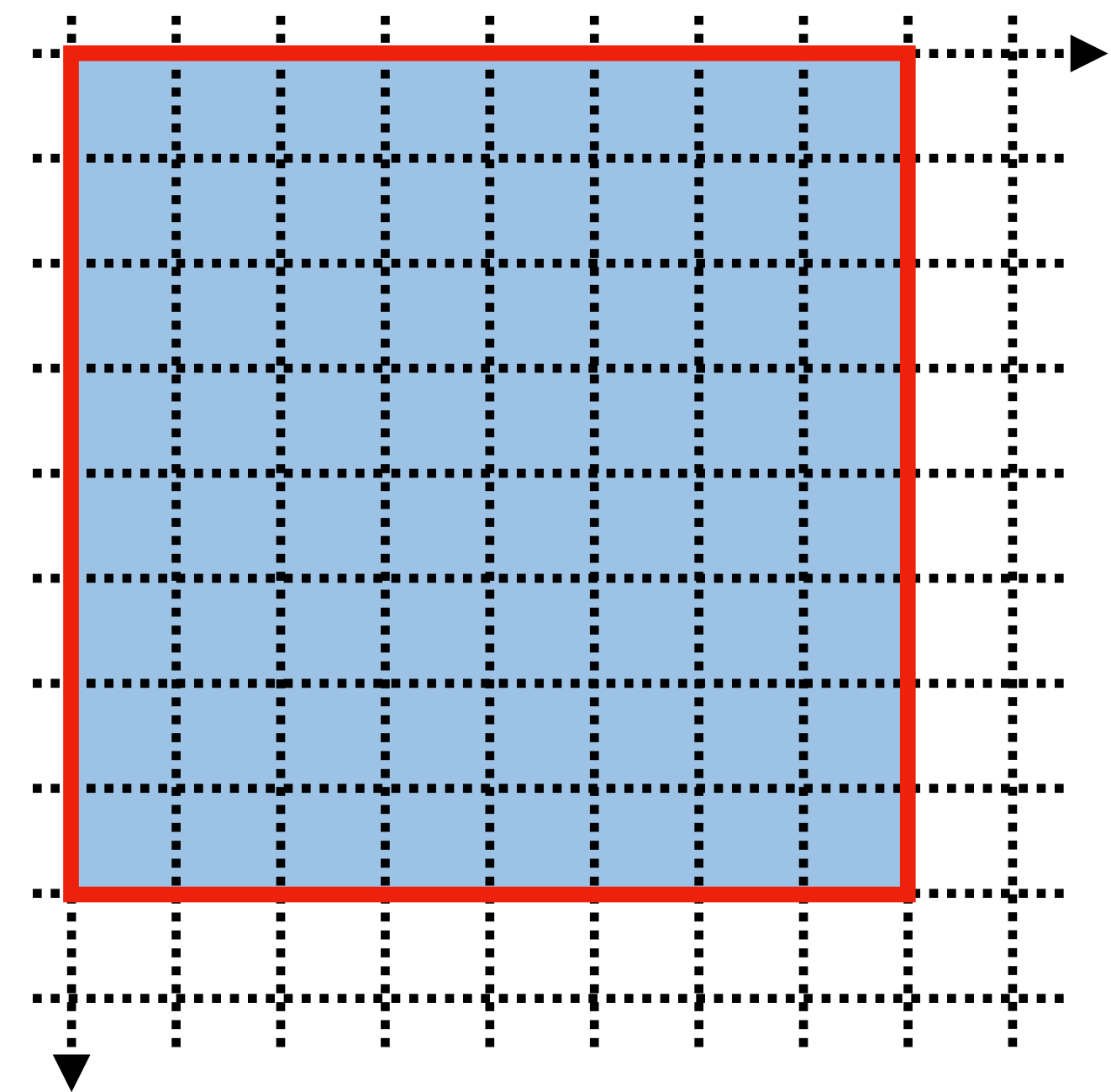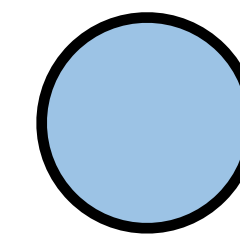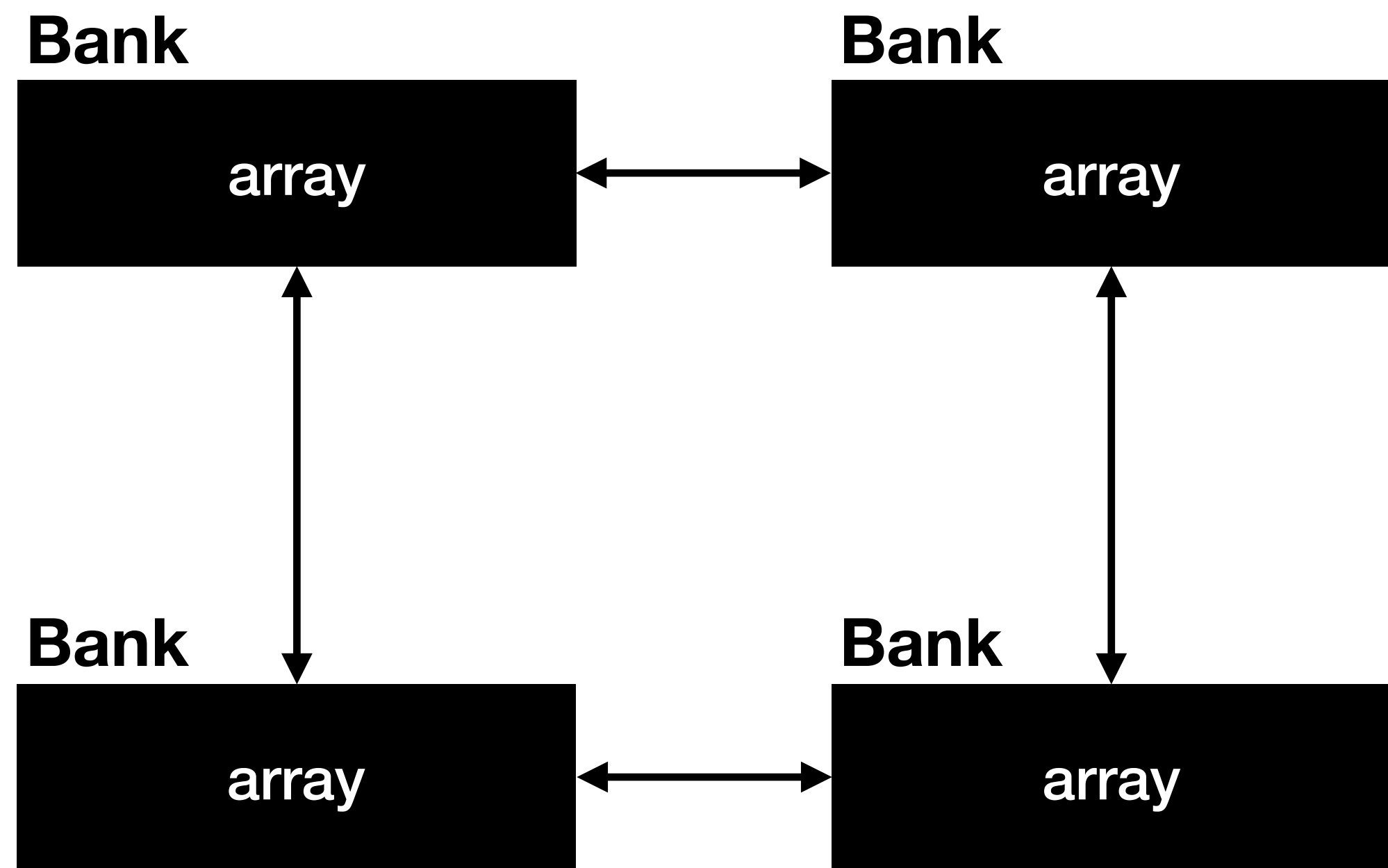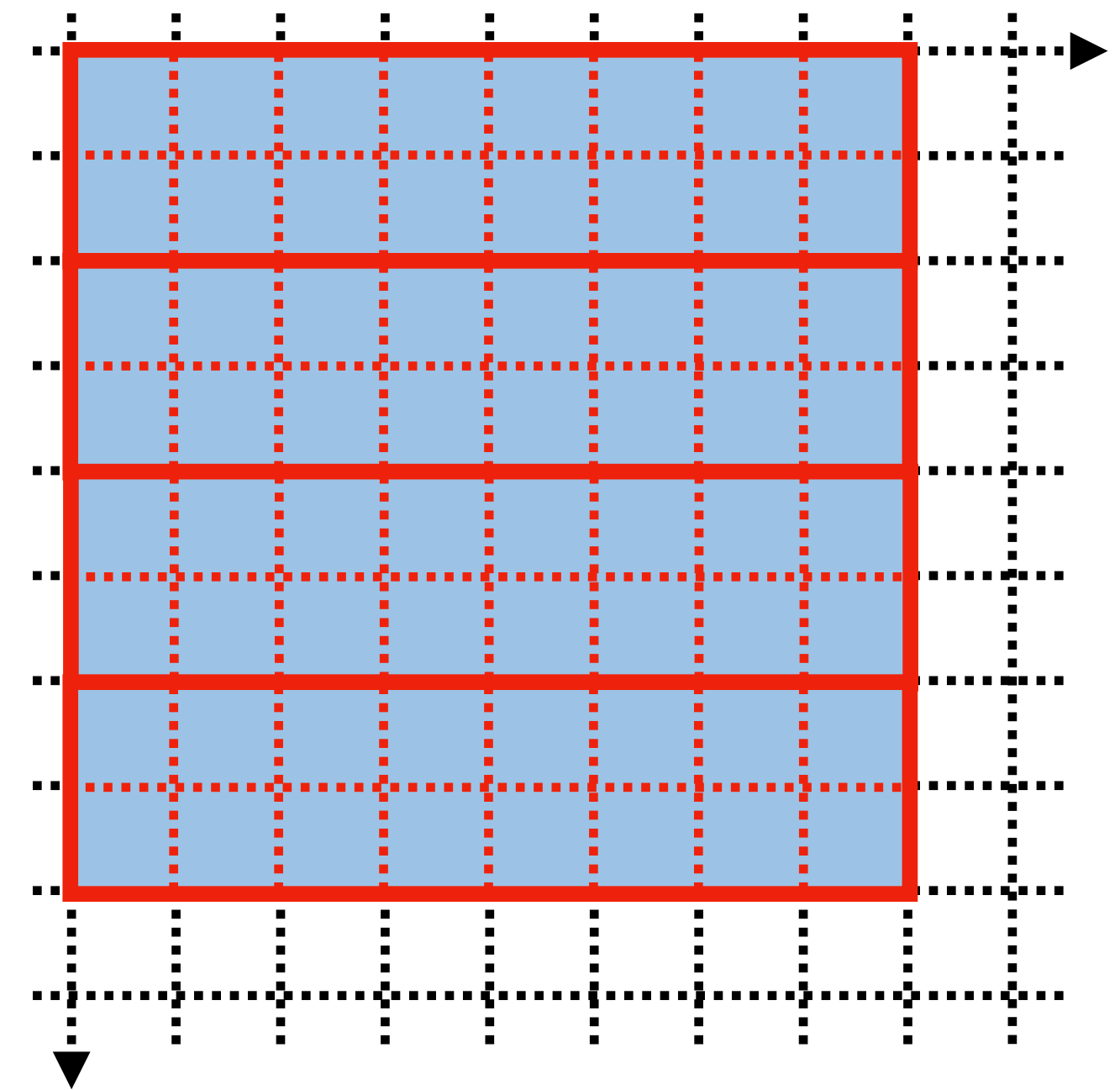
# In-Memory Organization

**Question**
What are the effects of tiling on
on inter-bank traffic?

Abstract many SRAM arrays
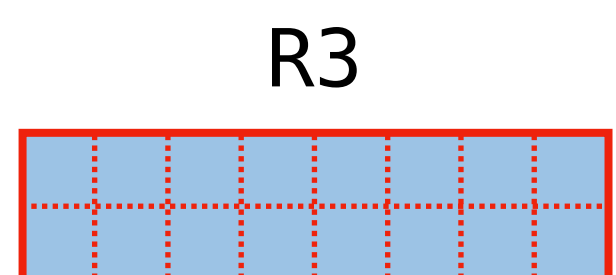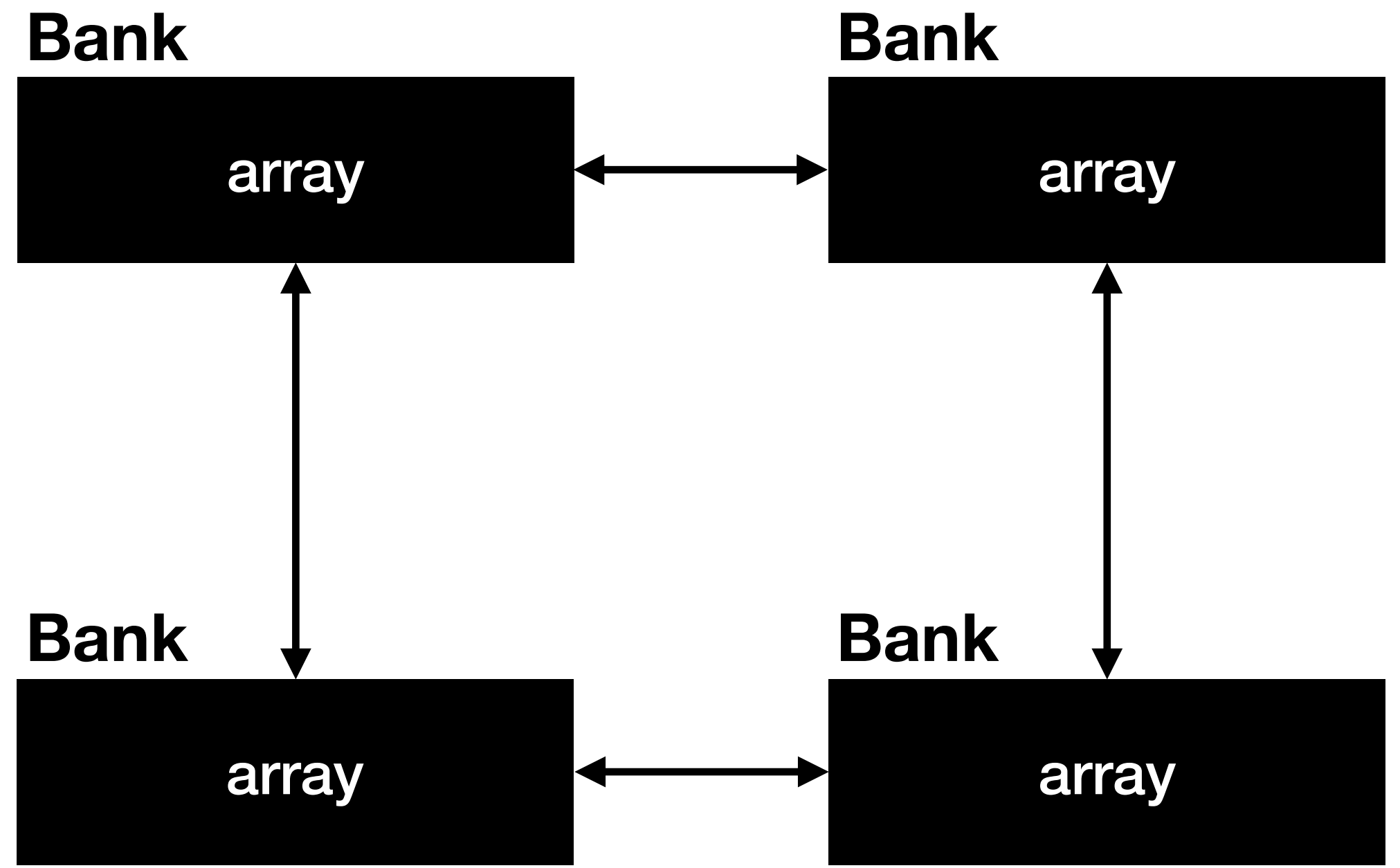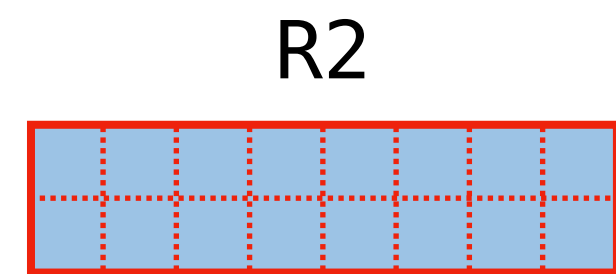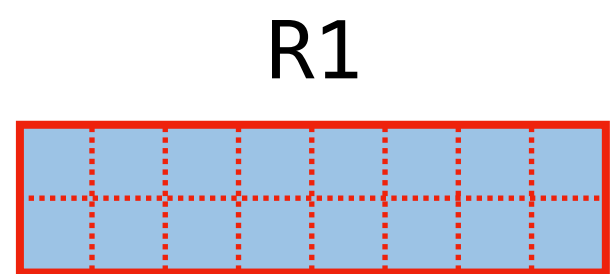into a single large SRAM array



array

array

array

array

**Bank**

array

**Bank**

array

**Bank**

array

**Bank**

array

Lattice cells are a sub-matrix of some larger matrix

R1

R2

**Bank**

array

**Bank**

array

**Bank**

array

**Bank**

array

R3

R4

R1

Q1

R2

Q2

**Bank**
array

**Bank**
array

**Bank**
array

**Bank**
array

R3

Q3

R4

Q4

1

mv

74

R1

Q1

R2

Q2

**Bank**
array

**Bank**
array

**Bank**
array

**Bank**
array

R3

Q3

R4

Q4

1

mv

Unmapped data

R1

Q1

R2

Q2

**Bank**
array

**Bank**
array

**Bank**
array

**Bank**
array

R3

Q3

R4

Q4

1

mv

78

Unmapped data

R1  Q1  R2  Q2

**Bank** array  **Bank** array

**Bank** array  **Bank** array

R3  Q3  R4  Q4

81

1 → mv

**Inter-Bank Data Movement**

Row     Square

<

R1

Q1

R2

Q2

**Bank**
array

**Bank**
array

**Bank**
array

**Bank**
array

1 mv

R3

Q3

R4

Q4

84

Unmapped data

R1    Q1    R2    Q2

Bank
array

Bank
array

Bank
array

Bank
array

R3    Q3    R4    Q4

1  mv

Inter-Bank Data Movement

Row    Square

>

87

# Inter-bank Data Movement
*depends on*
# Tiling Size + Move Direction

# Inter-bank Data Movement
## *depends on*
# Tiling Size + Move Direction

**Static**

Compiler

→

**Dynamic**

Runtime &
JIT Compiler

## Layout Hints

**1. Inter-Bank Traffic:**

*Best tiling size that minimizes inter-bank data movement*

Tensor (T) is moved along dimensions $d_1$, $d_2$, $d_3$, …

**2. Data Alignment:**

*Tensors A & B must have matching tiling size*

Tensors (A) & (B) are used together

# Microarchitecture Extensions

**Layout Override Table**

Tracks which arrays are currently transposed

**Tensor Transpose Unit**

Transposes array elements to and from bit-serial format

Core

L1 $

*Tags*

L2 $

Tensor Controller

LOT

Stream Engine

Stream Engine

Tensor Controller

LOT

L3 Way $i$

Tags

Router

SRAM Array  SRAM Array   SRAM Array  SRAM Array

SRAM Array  SRAM Array   SRAM Array  SRAM Array

TTU

More details in paper

# Outline

**②**

**Unified Abstraction**

Agnostically represents
1. In-core compute
2. In-memory compute
3. Near-memory compute

**③**

**μArch-specific Optimizations**

Specialize binary to take advantage of μArch details

**①**

**Compute Paradigms**

*In-Core*

Core

Memory

**④**

**Partitioning
Static vs Dynamic**

Portable binary without exposing μArch details to compiler

**μArch Extensions**

Hardware support for compute paradigms

Plain C Code

*In-Memory*

Memory

*Near-Memory*

Memory

# Outline

**Static**

**Dynamic**

**Compute Paradigms**

**Plain C Code**

**Unified Abstraction**

Agnostically represents
1. In-core compute
2. In-memory compute
3. Near-memory compute

**μArch-specific Optimizations**

Specialize binary to take advantage of μArch details

**μArch Support**

Hardware support for compute paradigms

*In-Core*

Core

Memory

*In-Memory*

Memory

*Near-Memory*

Memory

# Key Challenges



**Static**

Plain C Code → **Unified Abstraction**

Agnostically represents
1. In-core compute
2. In-memory compute
3. Near-memory compute

**Dynamic**

**μArch-specific Optimizations**

Specialize binary to take advantage of μArch details

**μArch Extensions**

Hardware support for compute paradigms

**Compute Paradigms**

*In-Core* — Core ⟷ Memory

*In-Memory* — Memory

*Near-Memory* — Memory

# Outline

**Static**

**Dynamic**

**Compute Paradigms**

Plain C Code

Compiler

Runtime

JIT Compiler

µArch Extensions

*In-Core*

Core

Memory

*In-Memory*

Memory

*Near-Memory*

Memory

94

**Static**

**Dynamic**

**Compute Paradigms**



Plain C Code

Compiler

Runtime

JIT Compiler

μArch Extensions

*In-Core*

Core

Memory

*In-Memory*

Memory

*Near-Memory*

Memory

Requirements

- Data ***alignment & layout***
- ***Portable*** binary
- ***Decide between*** in-core, in-/near-memory, and fusion
- Managing ***on-chip space***
- ***Tiling***
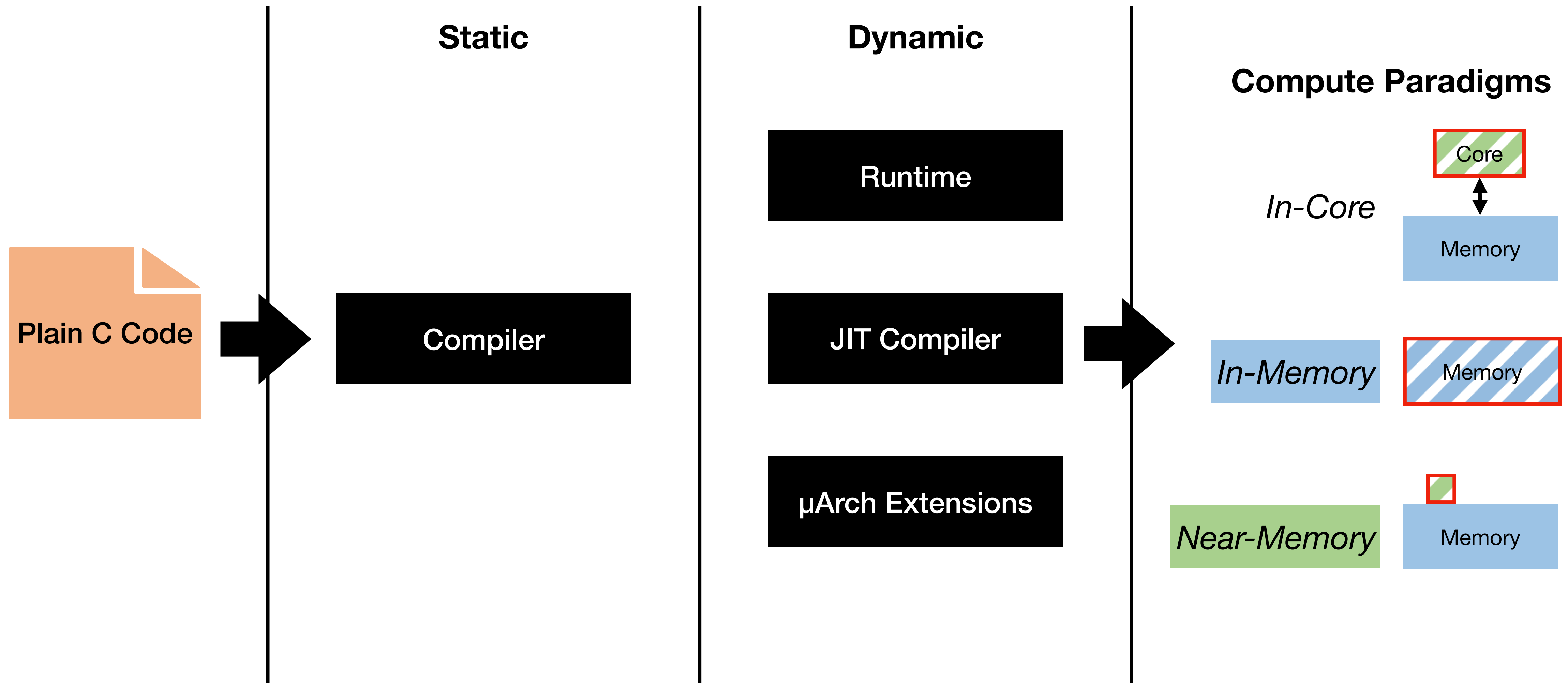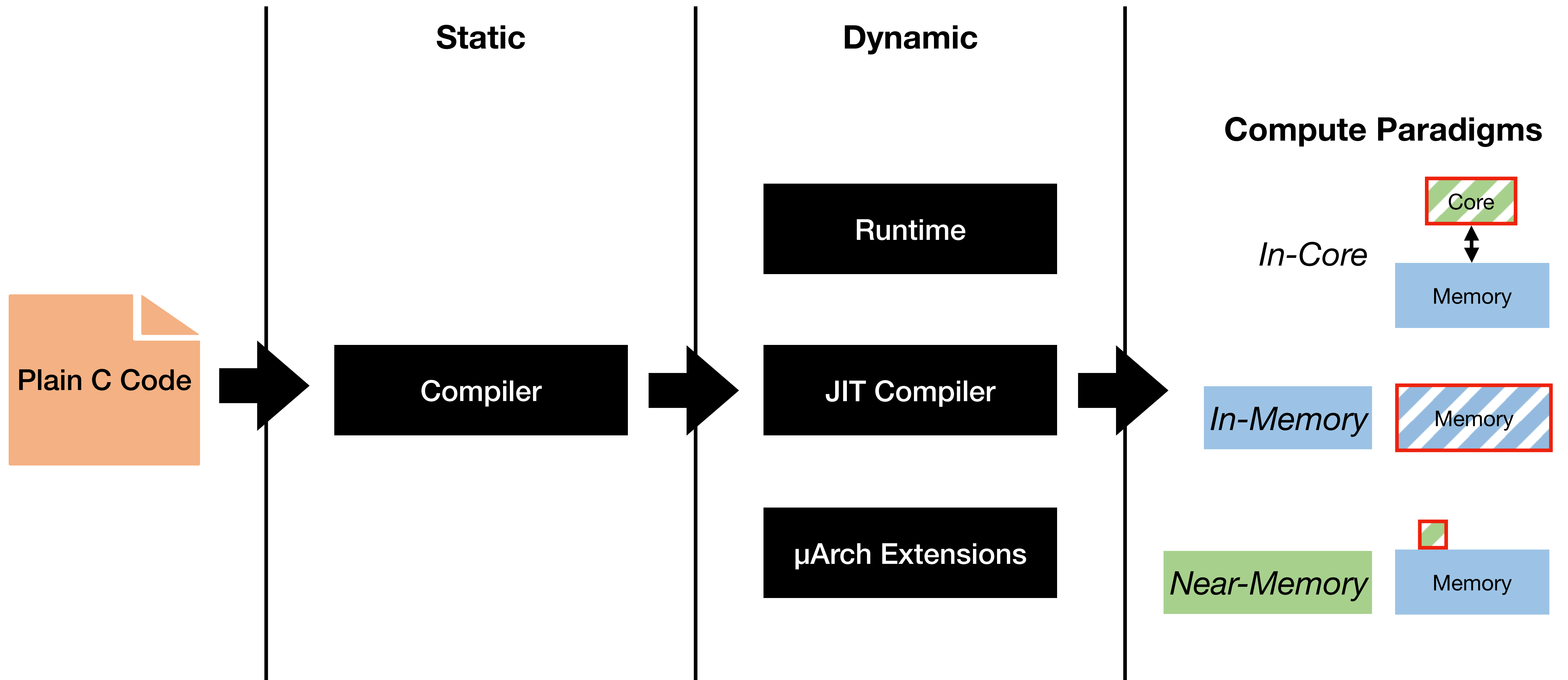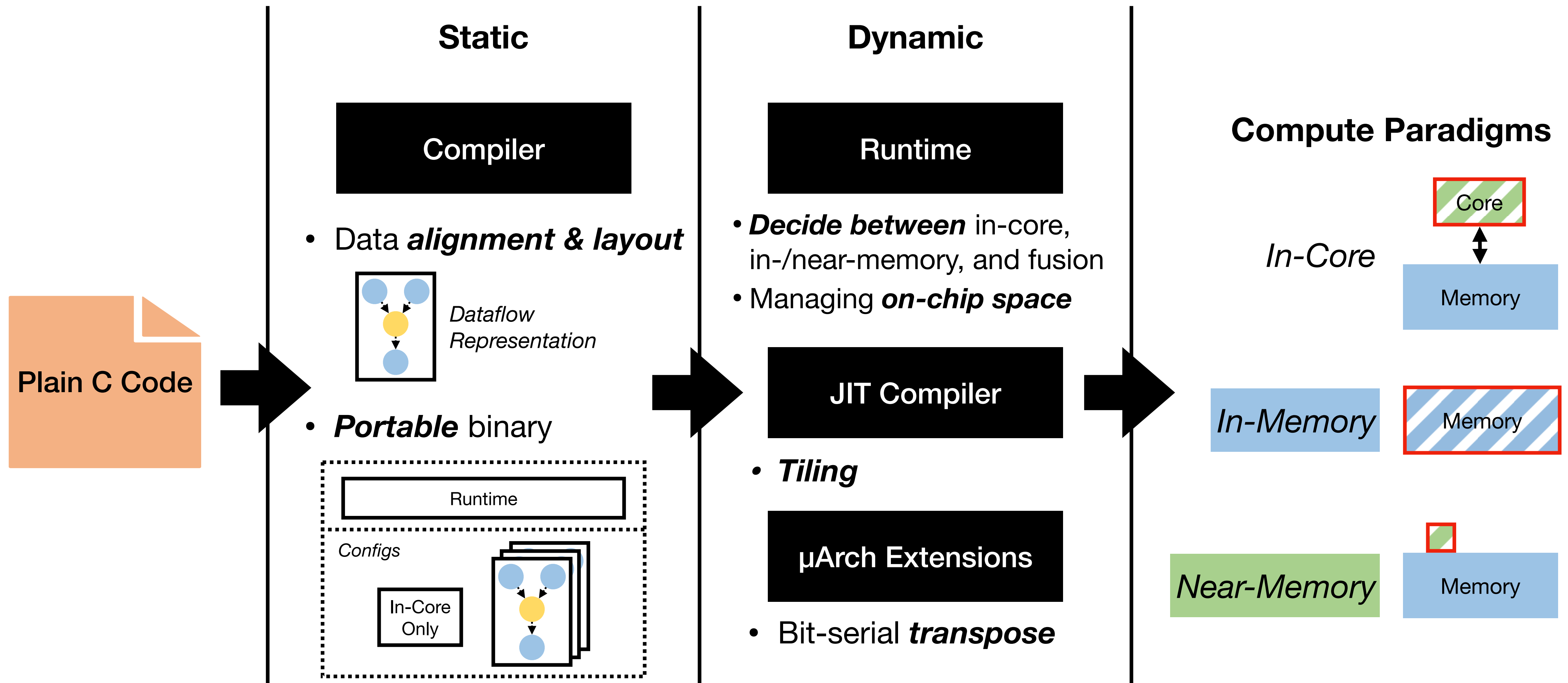- Bit-serial ***transpose***

# The Infinity Stream Approach

**Static**

**Dynamic**

**Compute Paradigms**

Plain C Code

**Compiler**

- Data *alignment & layout*

*Dataflow Representation*

- *Portable* binary

Runtime

*Configs*

In-Core Only

**Runtime**

- *Decide between* in-core, in-/near-memory, and fusion
- Managing *on-chip space*

**JIT Compiler**

- *Tiling*

**µArch Extensions**

- Bit-serial *transpose*

*In-Core*

Core

Memory

*In-Memory*

Memory

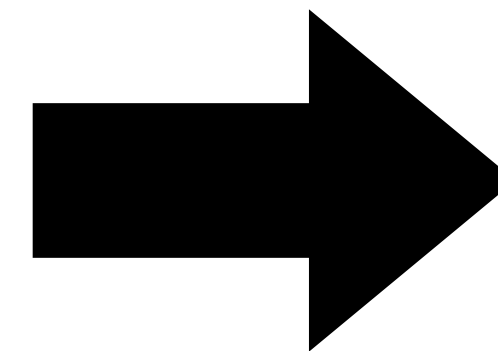*Near-Memory*

Memory

101

# Simulation Methodology
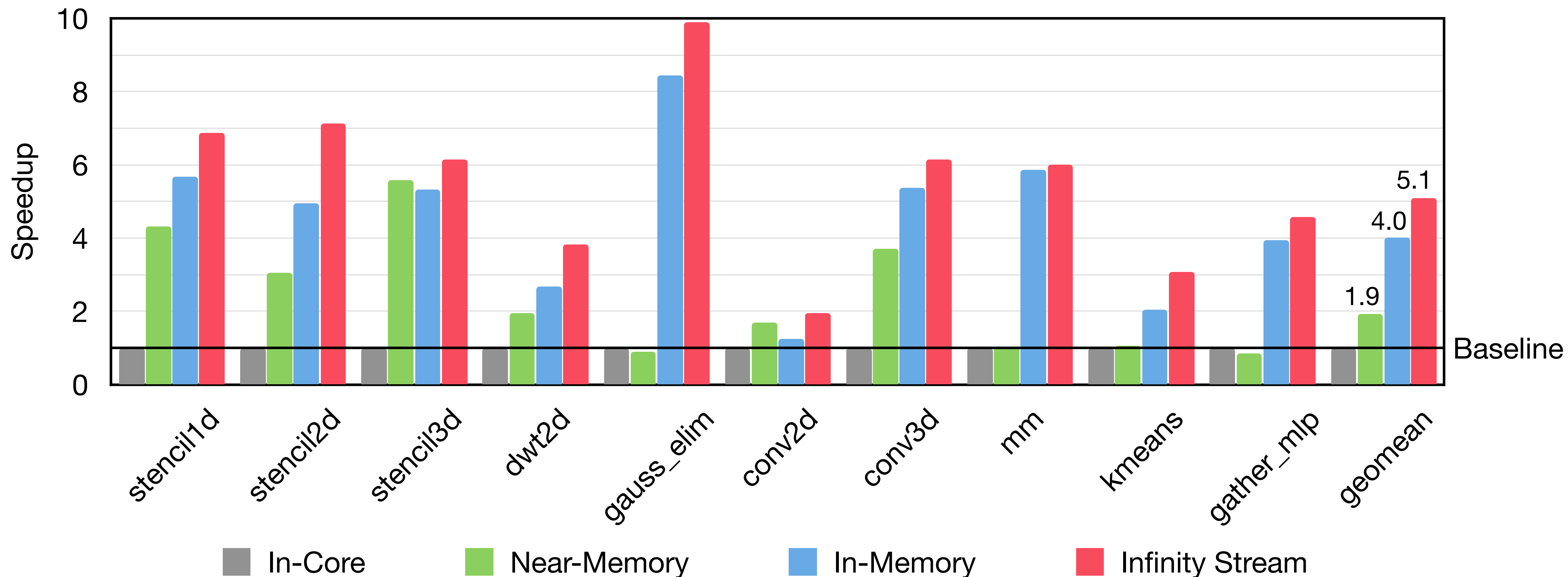
## Simulator

gem5 with partial AVX-512 support

## Hardware Configuration

- 64 Cores (8x8)
- 18 ways (16 SRAM arrays/way)
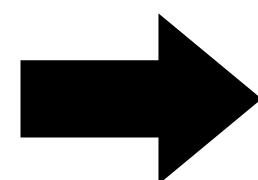- 8KB SRAM arrays (256x256)

L3 Total: 144MB

# Fused Execution Speedup



Legend: In-Core, Near-Memory, In-Memory, Infinity Stream

- 64 Cores (8x8)
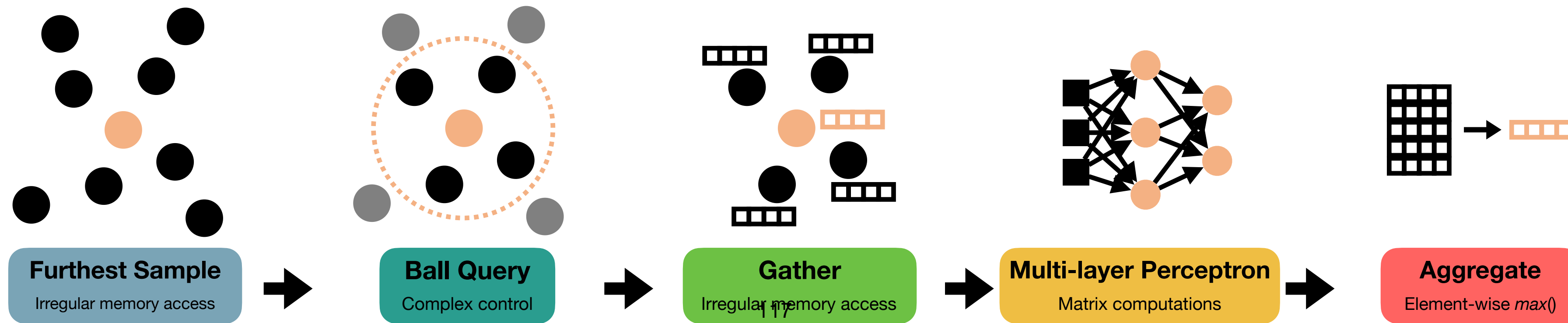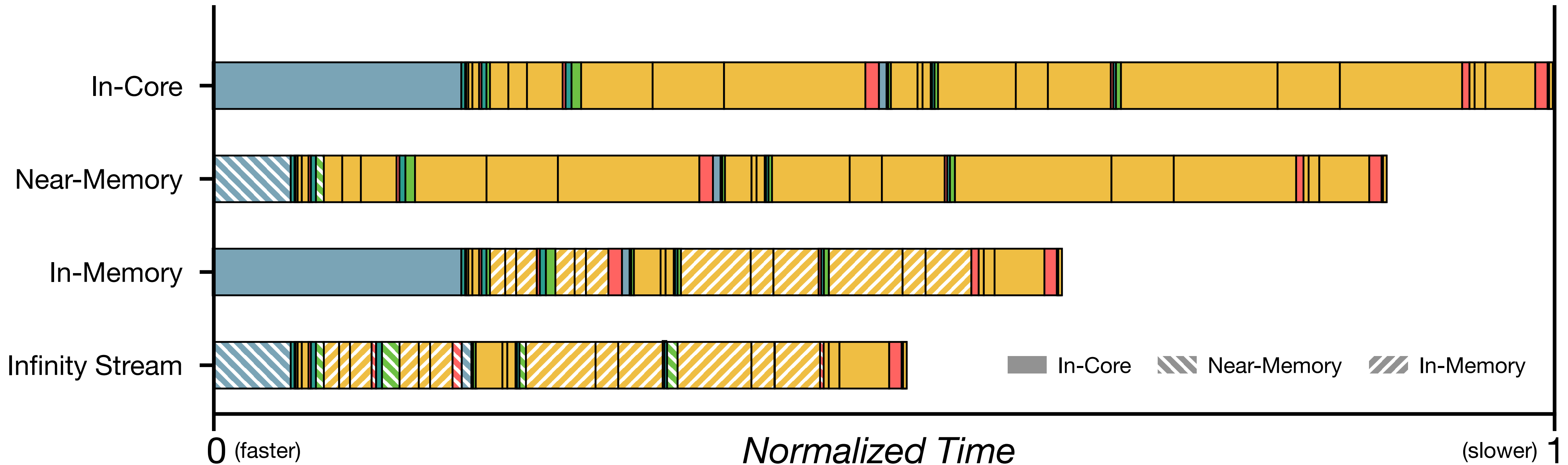- 18 ways (16 SRAM arrays/way)
- 8KB SRAM arrays (256x256)

L3 Total: 144MB

**Area** 6.52%
overhead

104

**Energy Efficiency**
- 5.6x over In-Core
- 2.4x over Near-Memory
- 1.6x over In-Memory

# Case Study: PointNet++

**Stream**
Memory Access Pattern
S

**Tensor**
Vectorized Memory Access
T

**Move**
Virtual Vector Alignment
mv

**Broadcast**
Spatial Reuse
bc

**Global Lattice Space**
Virtual Vector Lanes

**In-Core**
👍 Complex Control Flow

Core

Memory

👎 Von Neumann Bottleneck

**In-Memory**
👍 Massive Vector Processing

Memory

👎 Strict Alignment Req.

**Near-Memory**
👍 Complex Memory Patterns

Memory

👎 Low Compute Width

# Infinity Stream

**In-/Near-Memory Fusion**

Sequential

Parallel

Supports memory irregularity

Extremely high compute width

💡 Hints

JIT Compiler

**Runtime Tiling**
Minimizes Network Traffic

**Area**
6.52% overhead

**Energy Efficiency**
5.62x ⚡

Tensor Transpose Unit

Layout Override Table

**uArch Extensions**

**Fusion Speedup**

In-Core
Near-Memory 1.9x
In-Memory 4.0x
Infinity Stream 5.1x

118