

# Intentional Programming for HPC SoC Computing

**Wen-mei Hwu**

University of Illinois, Urbana-Champaign

MulticoreWare, Inc.

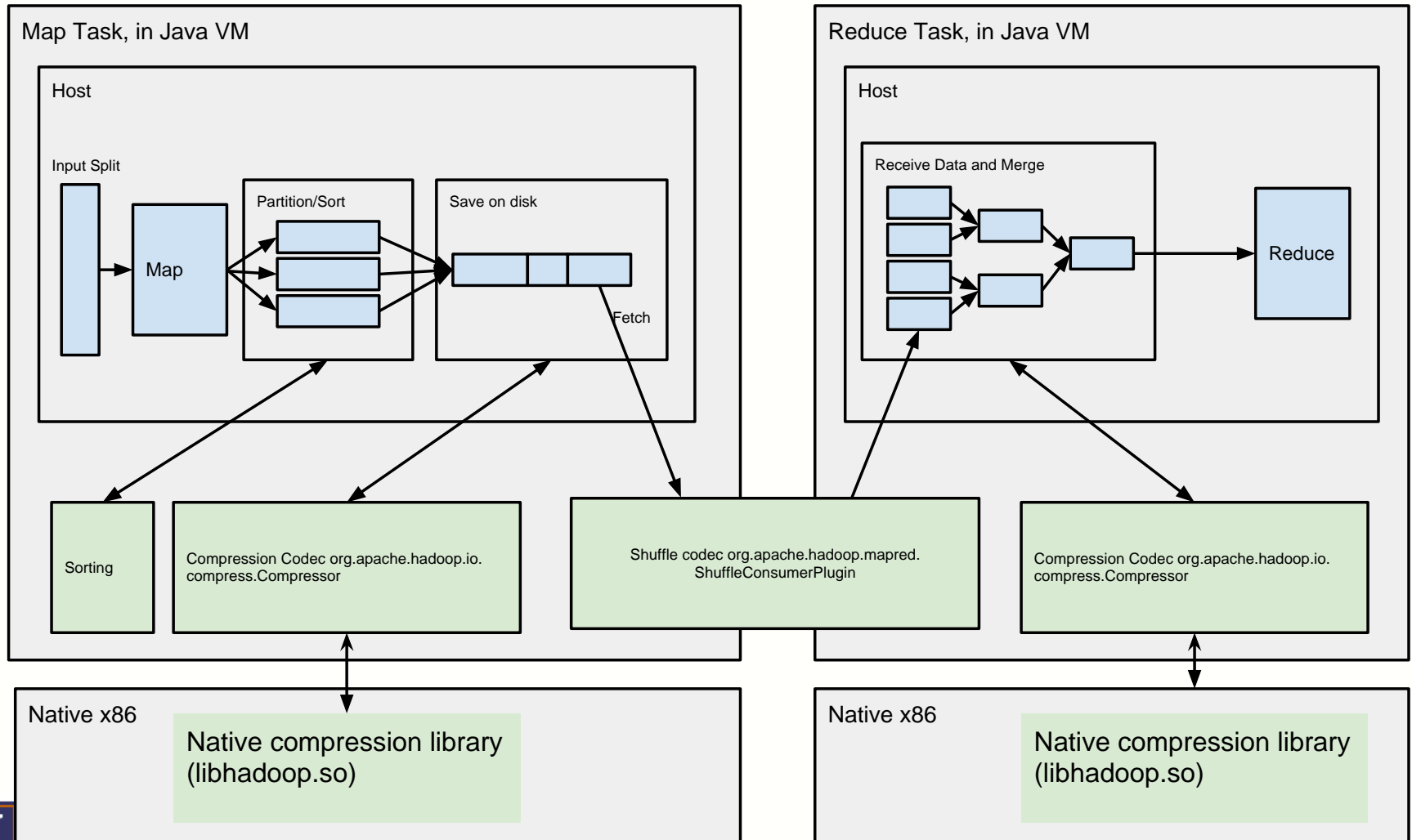


# SoC HPC

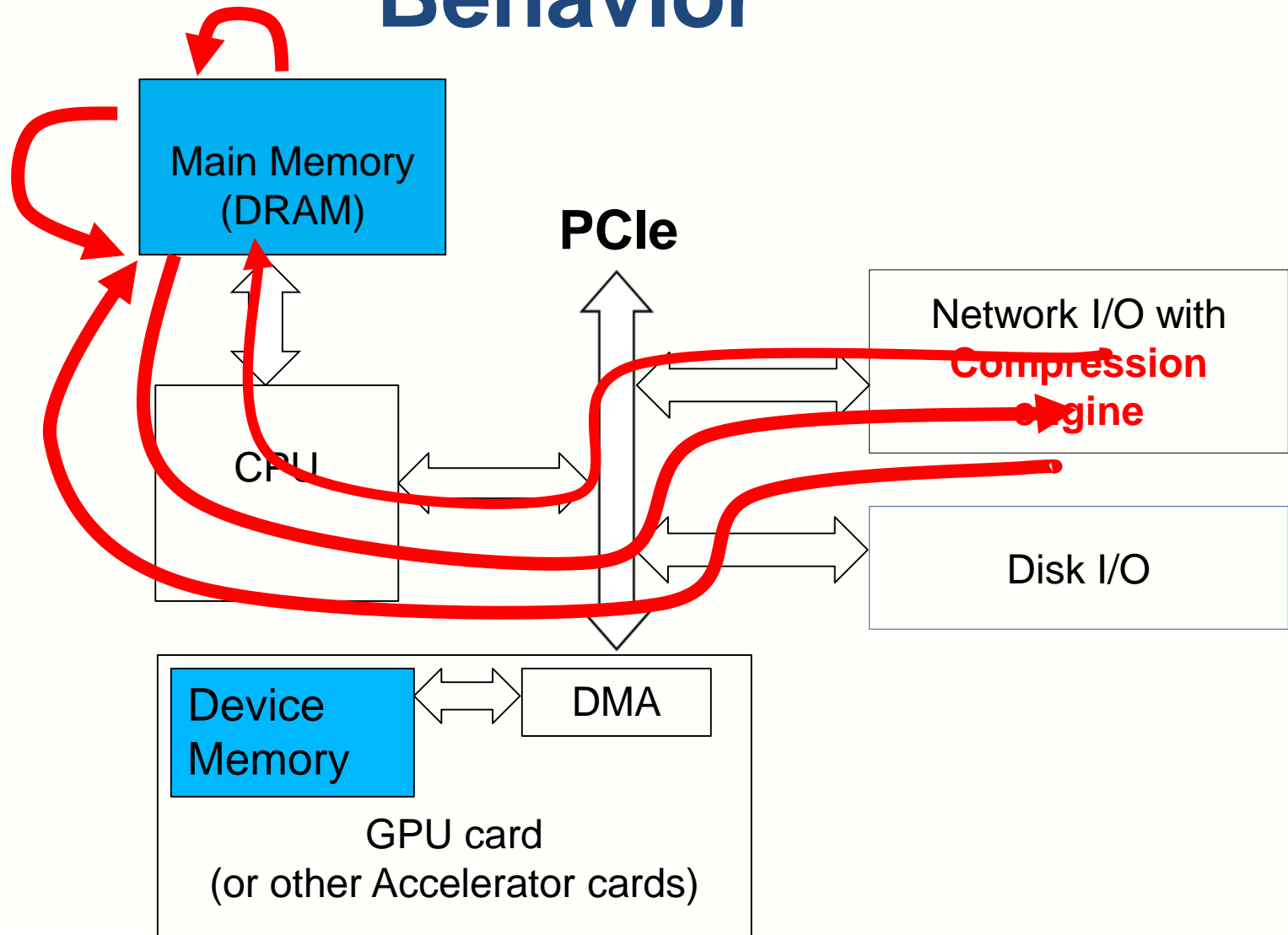
## Opportunities and Challenges

- Packaging and architecture innovations to
  - Reduce system-level data movement
    - Reduced latency and increased throughput
  - Enable distributed, heterogeneous computing
- Major impediment is software inertia.

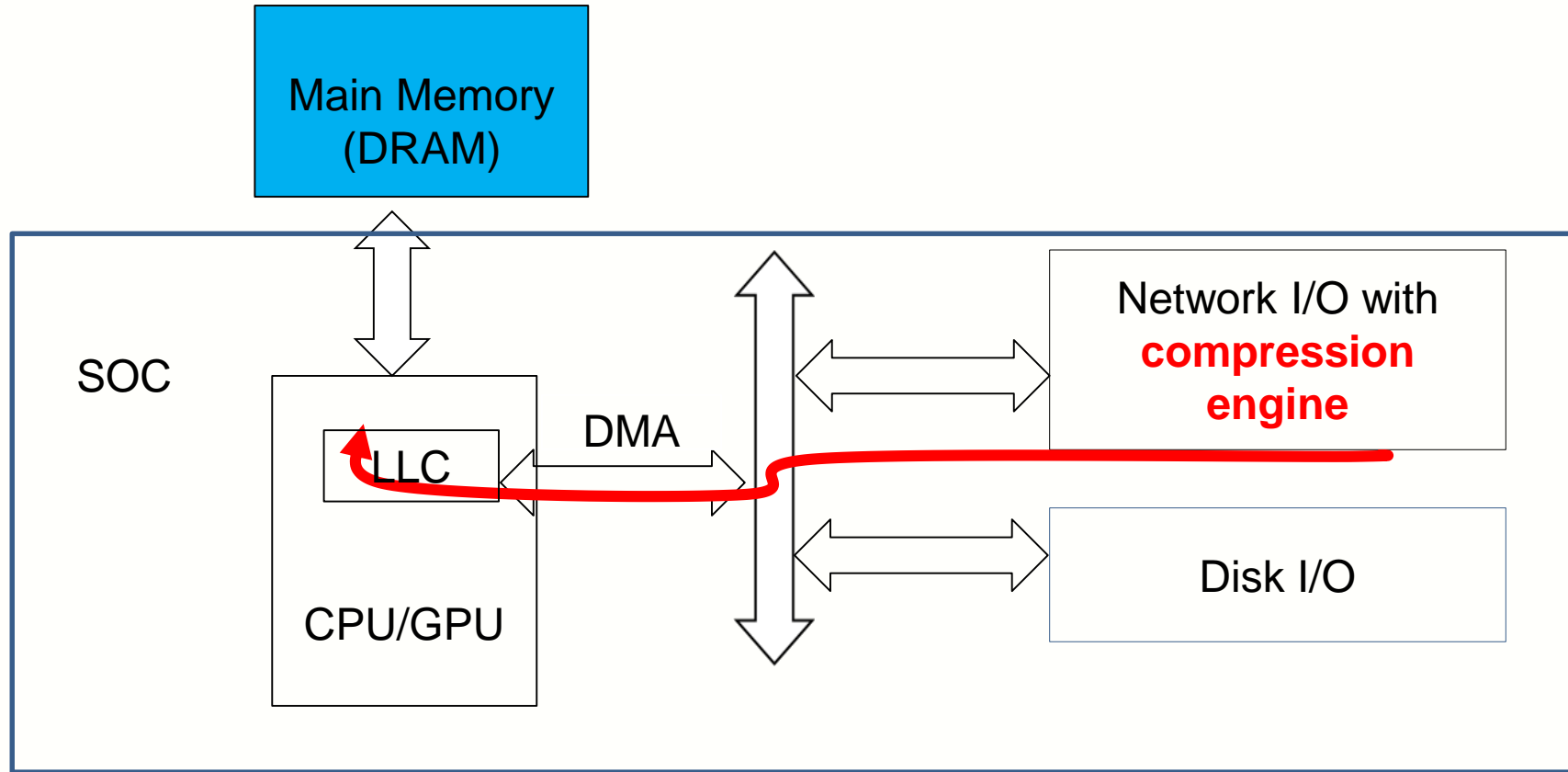
# Acceleration Challenges in Hadoop



# State of the Art Data Transfer Behavior



# Desired SoC Data Transfer Behavior

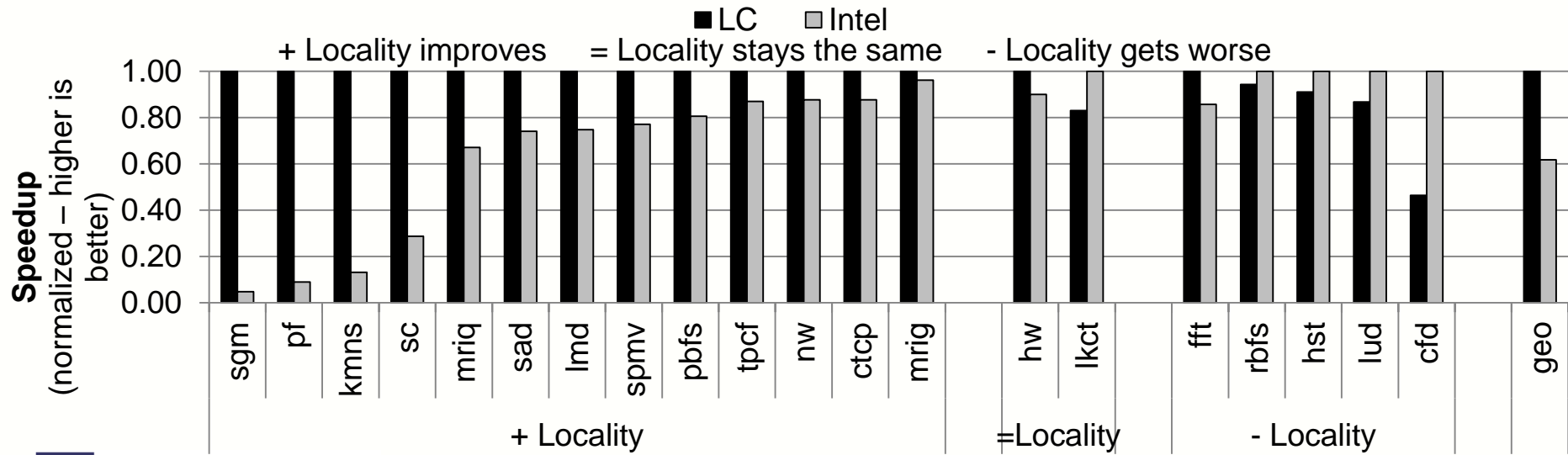
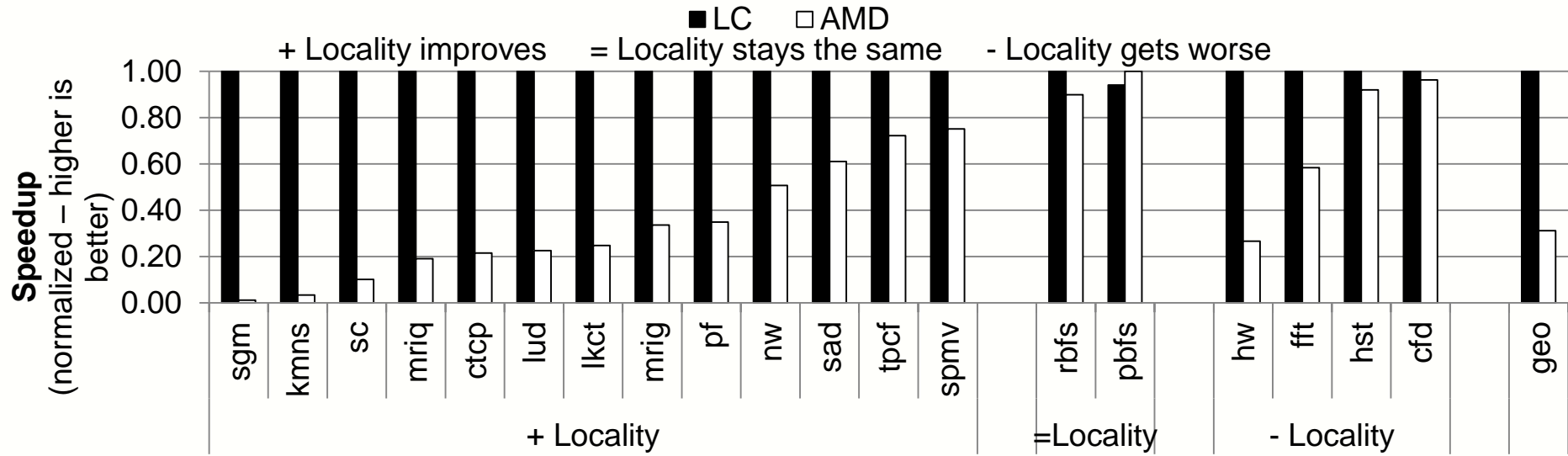


# Programming Heterogeneous Computing Systems

- Performance portability and locality with imperative parallel programming
  - OpenCL with MxPA (UIUC Research, now MulticoreWare Product)
- Algorithm selection, scalability, and efficiency with intentional parallel programming
  - Triolet [Rodrigues, et al, PPoPP 2014]



# MxPA Results: Comparison to Intel/AMD



# Programming in Triolet Python

Nonuniform FT (real part)

$$y_i = \sum_{j=0}^{n-1} x_j \cos(r_i k_j) \quad \text{for all } 0 \leq i < m$$



# Programming in Triolet Python

Nonuniform FT (real part)

Inner loop

$$y_i = \sum_{j=0}^{n-1} x_j \cos(r_i k_j) \quad \text{for all } 0 \leq i < m$$

Inner loop

```
sum(x * cos(r*k) for (x, k) in zip(xs, ks))
```



# Programming in Triolet Python

Nonuniform FT (real part)

Inner loop

Outer loop

$$y_i = \sum_{j=0}^{n-1} x_j \cos(r_i k_j) \quad \text{for all } 0 \leq i < m$$

```
ys = [ sum(x * cos(r*k) for (x, k) in zip(xs, ks)) for r in par(rs)]
```

# Programming in Triolet Python

Nonuniform FT (real part)

Inner loop

Outer loop

$$y_i = \sum_{j=0}^{n-1} x_j \cos(r_i k_j) \quad \text{for all } 0 \leq i < m$$

```
ys = [ sum(x * cos(r*k) for (x, k) in zip(xs, ks)) for r in par(rs)]
```

- “map and reduce” style programming—no new paradigm to learn
- Parallel details are implicit—easy to use
- Automated data partitioning, MPI rank generation, MPI messaging, OpenMP, etc.



# Productivity

## Triolet

```
ys = [sum(x * cos(r*k) for (x, k) in zip(xs, ks))
      for r in par(rs)]
```

- Library functions factor out data decomposition, parallelism, and communication

## C with MPI+OpenMP

```

1"#$1X8"#$8()*+,1X8"'"-.
1/01'2)X8'+345/01'26//7689;1;X8"#$8()*+<.
1/01'2)X8"=(#5/01'26//7689;1;X8"'"<-.
1*)#5!1"#$1(=#>?@1X8"'"<-!@!?.
1/01'A'+55;+34'B,1C,1/01'1DE,1?,1/01'26//7689;<.
1/01'A'+55;+34'>,1C,1/01'1DE,1?,1/01'26//7689;<.
1*)#5!1"#$1*FGH'+34'B!@!*"H-'IS+34'B,1X8"#$8()*+<.
1*)#5!1"#$18=-4-4"34'B!@!*FGH'+34'B!1X8"#$8()*+<.
1KH)=S1>+,1ZB+.
1"KIS(=#><1L
1!>+!@"#8G$">+.
1!B+!@"#8G$"B+.
1M
14H+1L
1!>+!@!X=HH)"5+34'>131+34)KSKH)=<<<.
1!B+!@!X=HH)"5+34'>131+34)KSKH)=<<<.
1M
1KH)=S1C(+*"FGH"@!X=HH)"5*FGH'+34'B!131+34)KSKH)=<<<.
1KH)=S13N+*"FGH"@!X=HH)"5*FGH'+34'B!131+34)KSKH)=<<<.
1"KIS(=#><1L
1!"#$180)>4(+!@!X8"'"#$8()*+PC.
1!1/01'84Q64+S13(4Q+@!X=HH)"5HO)>4(+131R131+34)K5/01'84Q64+<<<.
1!"#$10.
1!1K(150!@!?.101S180)>4(+.10TT<1L
1!1!"#$10)>4(+!@!10TC.
1!111/01'1+4B-5+,1+34',1/01'U96VE,10)>4(+!"->.
1!1111111111111?,1/01'26//7689;1;(4Q+MXX<.
1!111/01'1+4B-5B+,1+34',1/01'U96VE,10)>4(+!"->.
1!1111111111111?,1/01'26//7689;1;(4Q+M0)>4(+!TOX<.
1!111/01'1+4B-5(+!110)>4(+!"->3*FGH'+34'B,1*FGH'+34'B,1/01'U96VE,10)>4(+!"->.
1!1111111111111?,1/01'26//7689;1;(4Q+M4Y3HO)>4(+!TOX<.
1!1M
1!184Z*8N5(+*"FGH',1(+,1*FGH'+34'B!131+34)KSKH)=<<<.
1!1/01'7="S=HMSHO)>4(+3R,1(4Q+,1/01'ZEVE[Z'1]D68\<.
1!1K(44S(4Q+<.
1M
14H+1L
1!1/01'84*15+*,1+34',1/01'U96VE,1?,
1!11111111111?,1/01'26//7689;1/01'ZEVE[Z'1]D68\<.
1!1/01'84*15B+,1+34',1/01'U96VE,1?,
1!11111111111?,1/01'26//7689;1/01'ZEVE[Z'1]D68\<.
1!1/01'84*15(+*"FGH',1*FGH'+34'B,1/01'U96VE,1?,
1!11111111111?,1/01'26//7689;1/01'ZEVE[Z'1]D68\<.
1M
1!1
1!1!"#$1".
1!(=1)X818=(#HH4HK)(1+*F4-GH45+S5**<
1!1K(15!@!?.1"1S1*FGH'+34'B.1"TT<1L
1!1111KH)=S1+!@!?.
1!11!"#$1a.
1!111K(15!@!?.1a!S1+34'>.1aTT<
1!11111+!T@!B+M4X!11)*+K5(+*"FGH'+34'B,1/01'U96VE,1?,
1!111N+*"FGH'+34'B!1+<.
1!1M
1!1M
1/01']=SF4(SN+*FGH',1*FGH'+34'B,1/01'U96VE,1N+,1*FGH'+34'B,1/01'U96VE,
1!11111111111?,1/01'26//7689;<.
1K(44S(+*"FGH'+<.
1K(44SN+*FGH'+<.
1"KIS(=#><1L
1!1K(44S+<.
1M
14H+1L
1!1K(44SB+<.
1!1K(44S+<.
1M

```

Setup

```

1!1
1!1!"#$1".
1!(=1)X818=(#HH4HK)(1+*F4-GH45+S5**<
1!1K(15!@!?.1"1S1*FGH'+34'B.1"TT<1L
1!111KH)=S1+!@!?.
1!11!"#$1a.
1!111K(15!@!?.1a!S1+34'>.1aTT<
1!11111+!T@!B+M4X!11)*+K5(+*"FGH'+34'B,1/01'U96VE,1?,
1!111N+*"FGH'+34'B!1+<.
1!1M
1!1M
1/01']=SF4(SN+*FGH',1*FGH'+34'B,1/01'U96VE,1N+,1*FGH'+34'B,1/01'U96VE,
1!11111111111?,1/01'26//7689;<.
1K(44S(+*"FGH'+<.
1K(44SN+*FGH'+<.
1"KIS(=#><1L
1!1K(44S+<.
1M
14H+1L
1!1K(44SB+<.
1!1K(44S+<.
1M

```

Cleanup

# Productivity and Performance

## Triolet

$ys = [\text{sum}(x * \cos(r*k) \text{ for } (x, k) \text{ in } \text{zip}(xs, ks)) \text{ for } r \text{ in } \text{par}(rs)]$

- Library functions factor out data decomposition, parallelism, and communication

128-way Speedup (16 cores × 8 nodes)	Triolet	C with MPI+OpenMP
	99	115

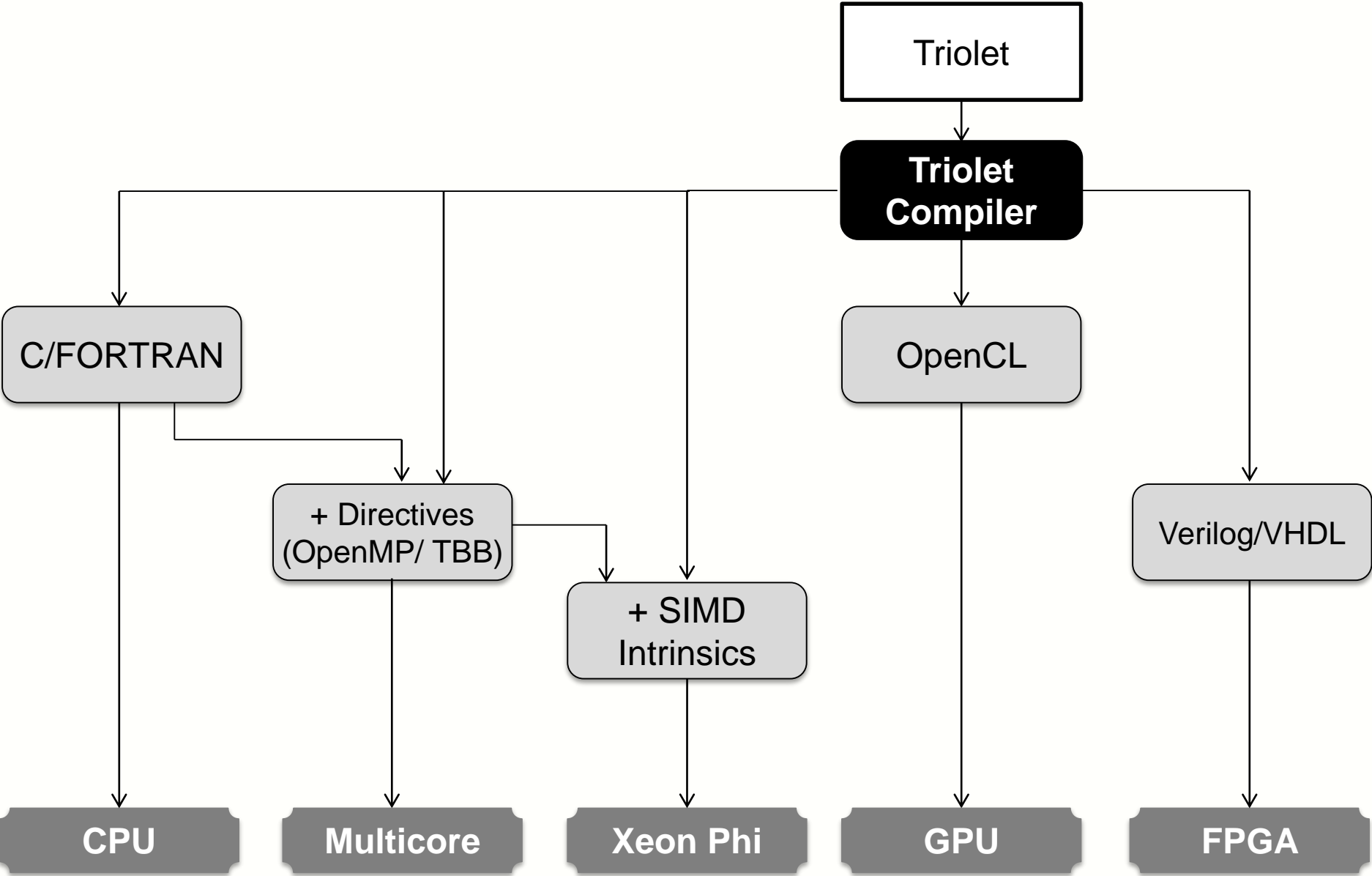
## C with MPI+OpenMP

```
int main() {
    int n = 100000000;
    double *A = malloc(n * sizeof(double));
    double *B = malloc(n * sizeof(double));
    double *C = malloc(n * sizeof(double));
    double *x = malloc(n * sizeof(double));
    double *k = malloc(n * sizeof(double));
    double *r = malloc(n * sizeof(double));
    double *xs = malloc(n * sizeof(double));
    double *ks = malloc(n * sizeof(double));
    double *rs = malloc(n * sizeof(double));
    double *y = malloc(n * sizeof(double));
    for (int i = 0; i < n; i++) {
        x[i] = 1.0 / (i + 1);
        k[i] = 1.0 / (i + 1);
        r[i] = 1.0 / (i + 1);
        xs[i] = 1.0 / (i + 1);
        ks[i] = 1.0 / (i + 1);
        rs[i] = 1.0 / (i + 1);
    }
    for (int i = 0; i < n; i++) {
        y[i] = 0.0;
    }
    for (int i = 0; i < n; i++) {
        for (int j = 0; j < n; j++) {
            y[i] += x[j] * cos(r[i] * k[j]);
        }
    }
    return 0;
}
```

Setup

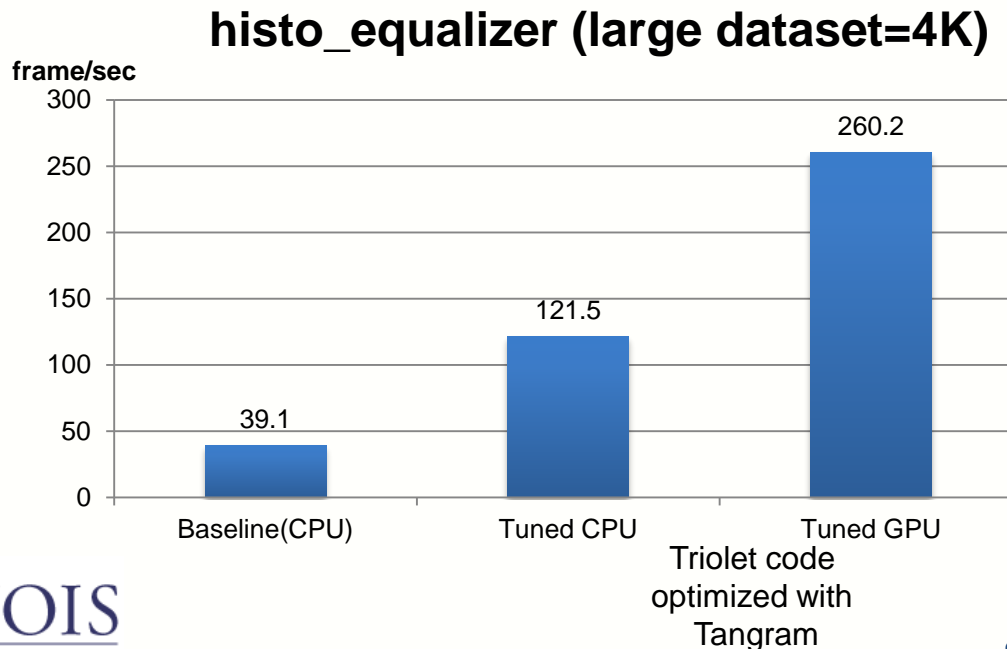
```
int main() {
    int n = 100000000;
    double *A = malloc(n * sizeof(double));
    double *B = malloc(n * sizeof(double));
    double *C = malloc(n * sizeof(double));
    double *x = malloc(n * sizeof(double));
    double *k = malloc(n * sizeof(double));
    double *r = malloc(n * sizeof(double));
    double *xs = malloc(n * sizeof(double));
    double *ks = malloc(n * sizeof(double));
    double *rs = malloc(n * sizeof(double));
    double *y = malloc(n * sizeof(double));
    for (int i = 0; i < n; i++) {
        x[i] = 1.0 / (i + 1);
        k[i] = 1.0 / (i + 1);
        r[i] = 1.0 / (i + 1);
        xs[i] = 1.0 / (i + 1);
        ks[i] = 1.0 / (i + 1);
        rs[i] = 1.0 / (i + 1);
    }
    for (int i = 0; i < n; i++) {
        for (int j = 0; j < n; j++) {
            y[i] += x[j] * cos(r[i] * k[j]);
        }
    }
    return 0;
}
```

Cleanup



# Triolet equalize\_frames Example

- Baseline is parallel code taken from DARPA PERFECT benchmark suite
- Tuned code outperforms baseline on both architectures

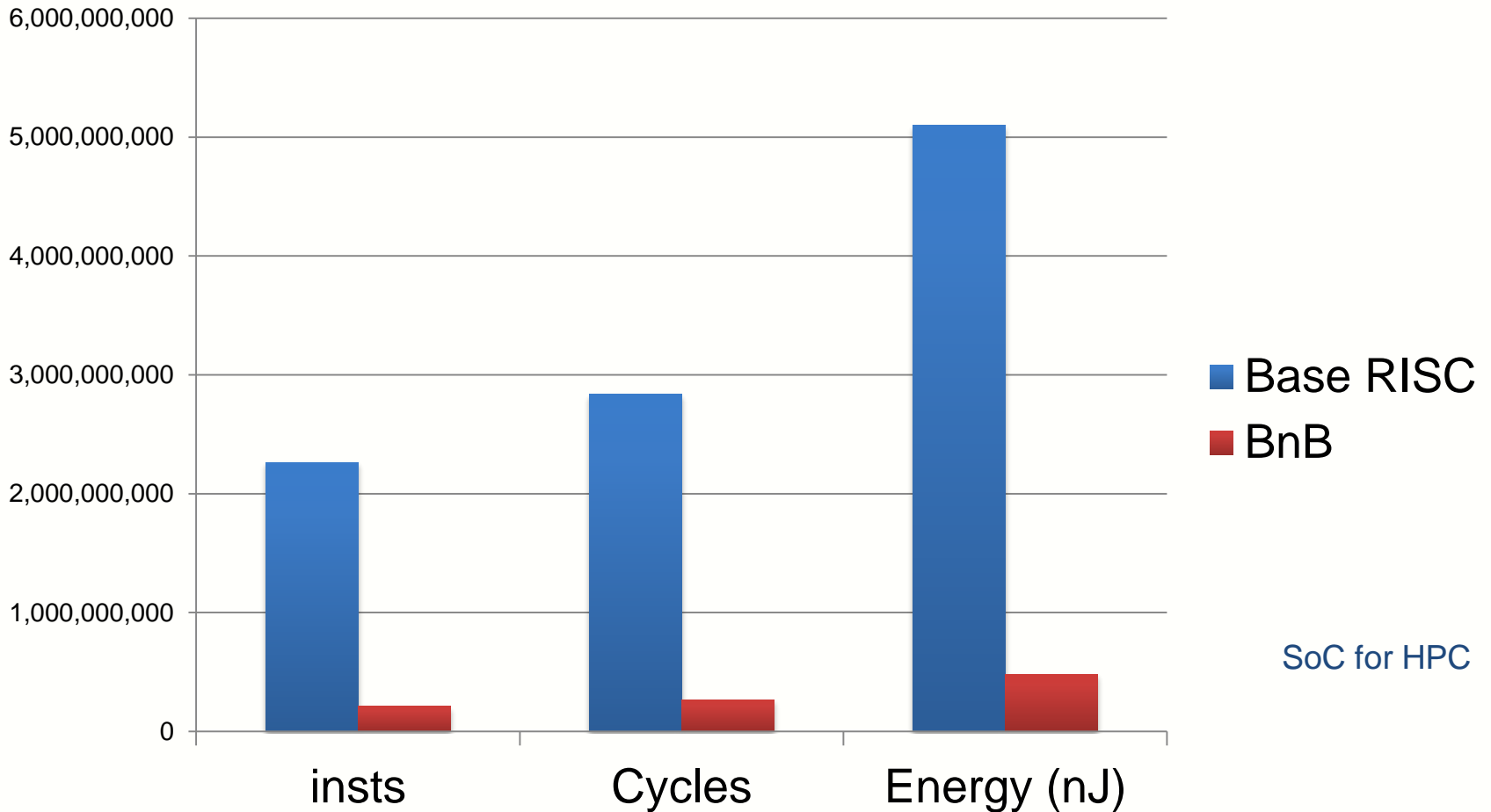


CPU time measured on  
Intel Xeon E5520  
(4-core, hyperthreaded)

GPU time measured on  
Tesla C2050



# Triolet Compiled Conv2D on BnB



~10.5x improvement across the board  
1024x1024 image, HMC memory system





# Conclusion

- Packing and architecture innovations enable continued improvement in performance and power
  - Hardware APIs for streamlined data movement
  - Hardware APIs for I/O-side and memory-side compute
- Intentional programming for system-level software mapping and transformation
  - Triolet Python brings intentional programming into mobile devices, HPC clusters, and distributed computing

# Essential work in writing efficient parallel code.

## Planning how to execute an algorithm

- Distribute computation across
  - cores,
  - hardware threads, and
  - vector processing elements
- Distribute data across
  - discrete GPUs or
  - clusters
- Orchestrate communication for
  - reductions,
  - variable-size list creation,
  - stencils, etc.

## Implementing the plan

- Rearrange data for locality
- Fuse or split loops
- Map loop iterations onto hardware
- Migrate code across SW components
  
- Allocate memory
- Partition data
- Insert data movement code
  
- Reduction trees
- Array packing
- Boundary cell communication



# Current State of Programming Heterogeneous Systems

CPU

Multicore

Xeon Phi

GPU

FPGA



# Current State of Programming Heterogeneous Systems

C/FORTRAN



CPU

Multicore

Xeon Phi

GPU

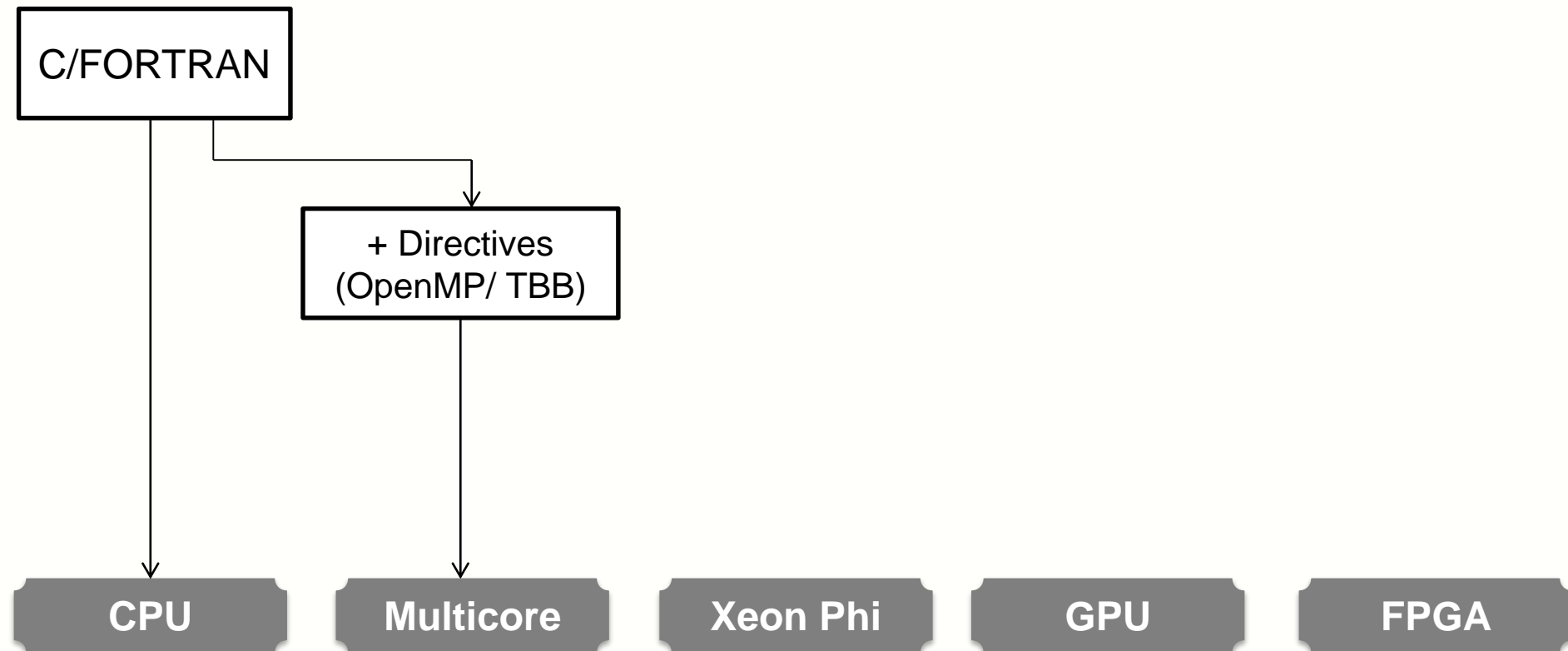
FPGA



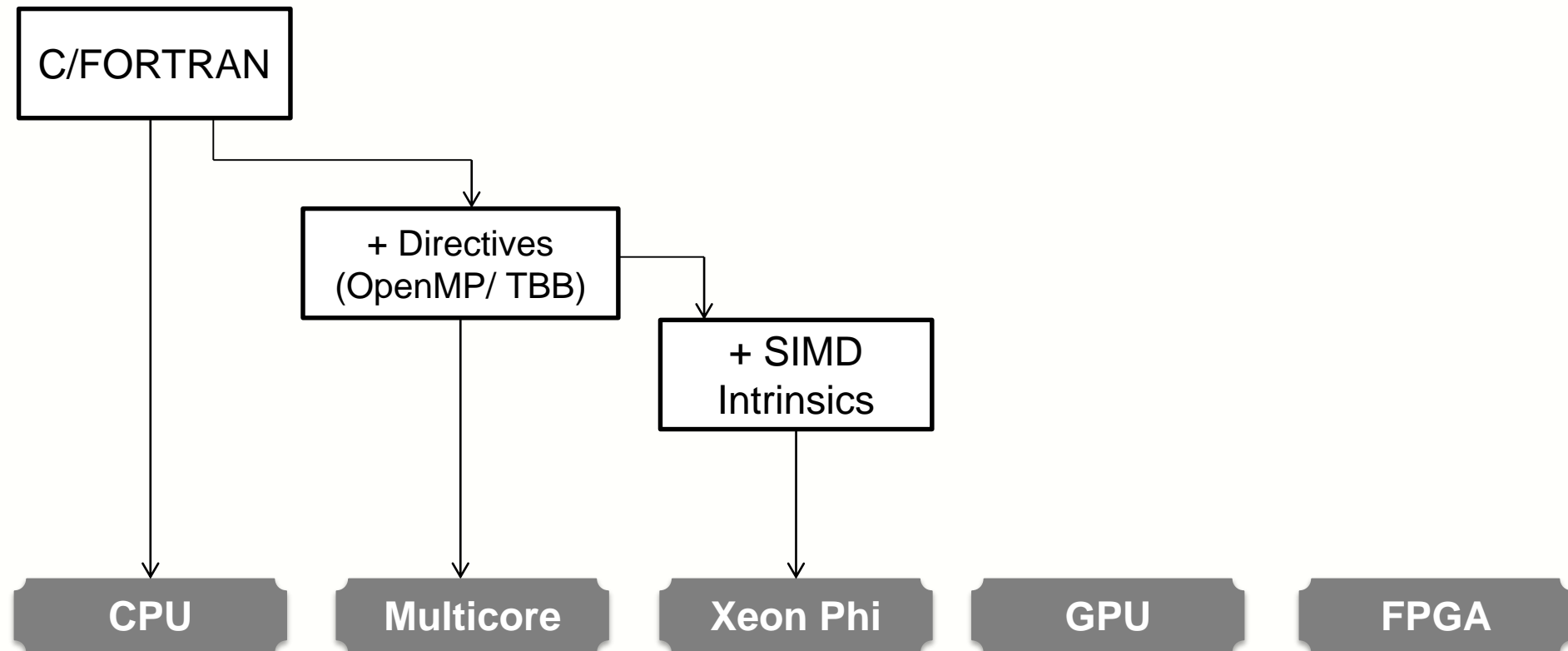
ILLINOIS  
UNIVERSITY OF ILLINOIS AT URBANA-CHAMPAIGN

SoC for HPC

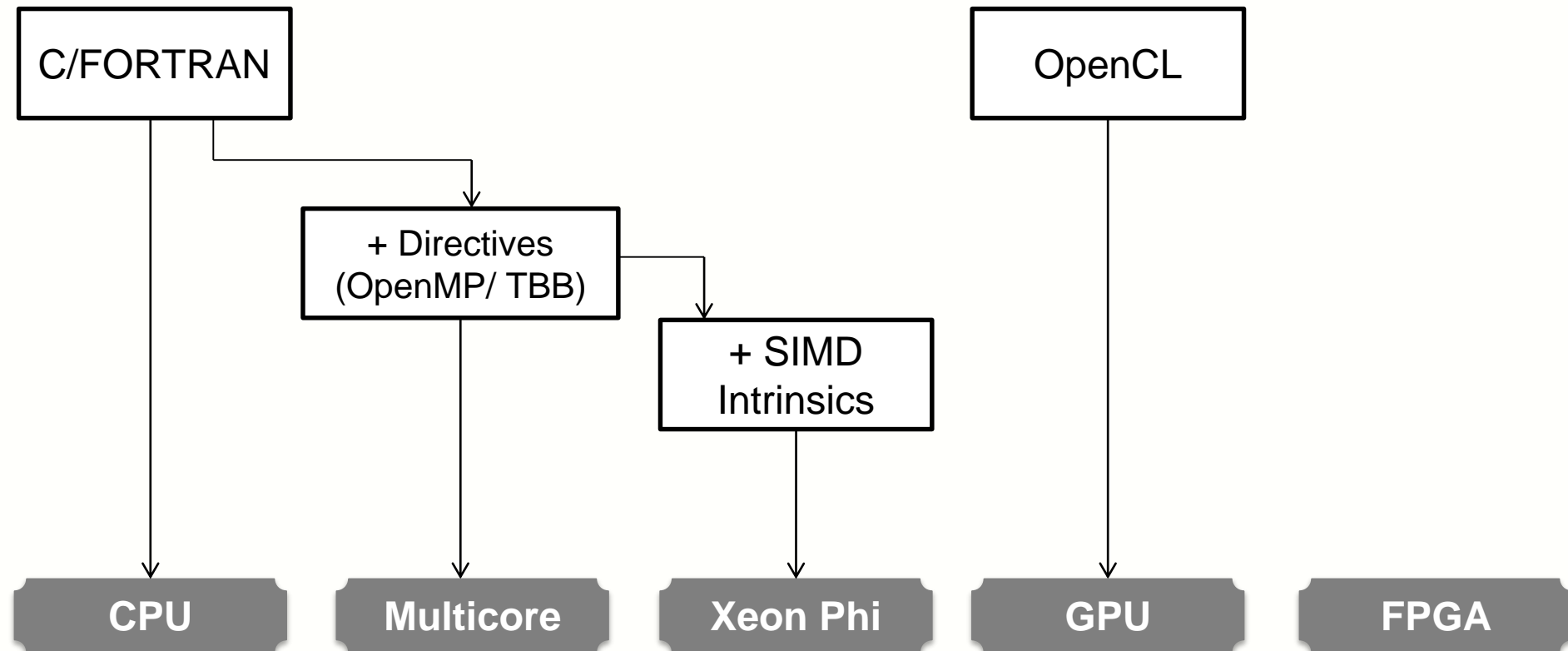
# Current State of Programming Heterogeneous Systems



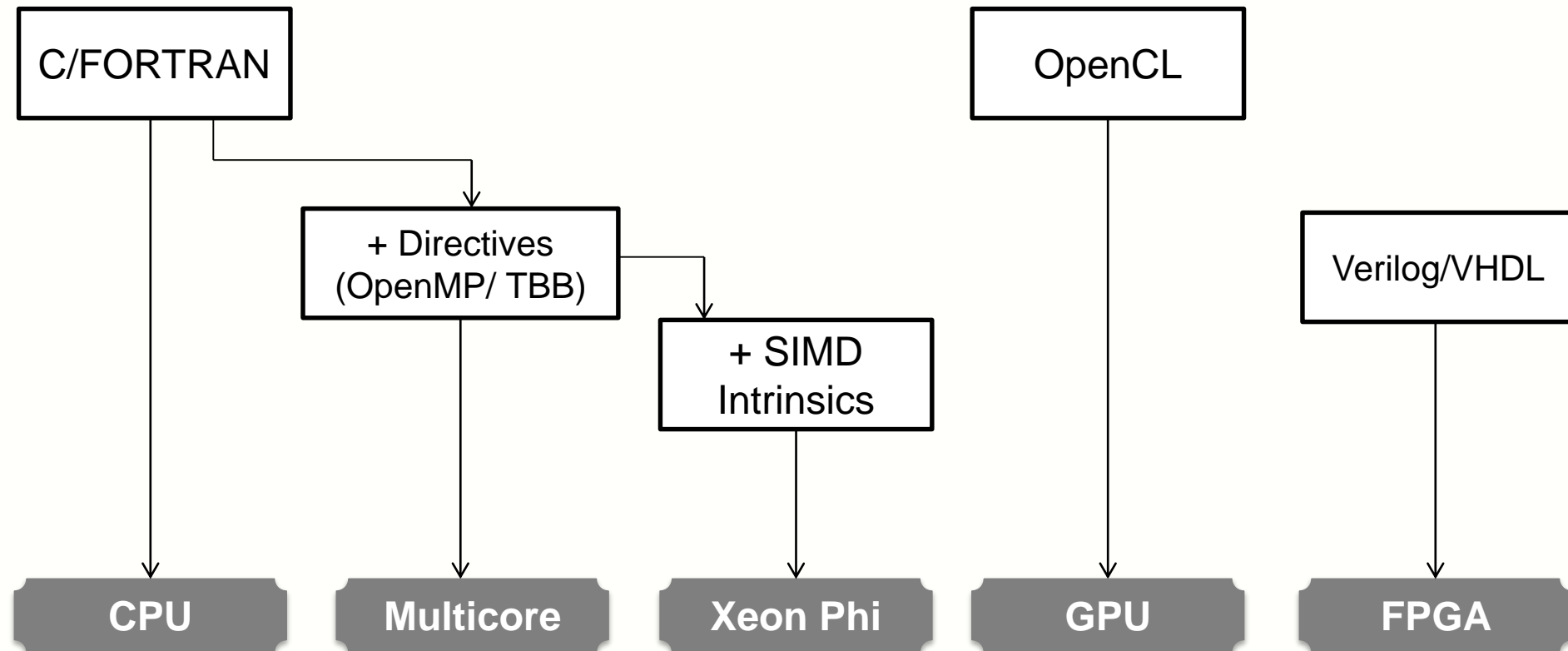
# Current State of Programming Heterogeneous Systems



# Current State of Programming Heterogeneous Systems



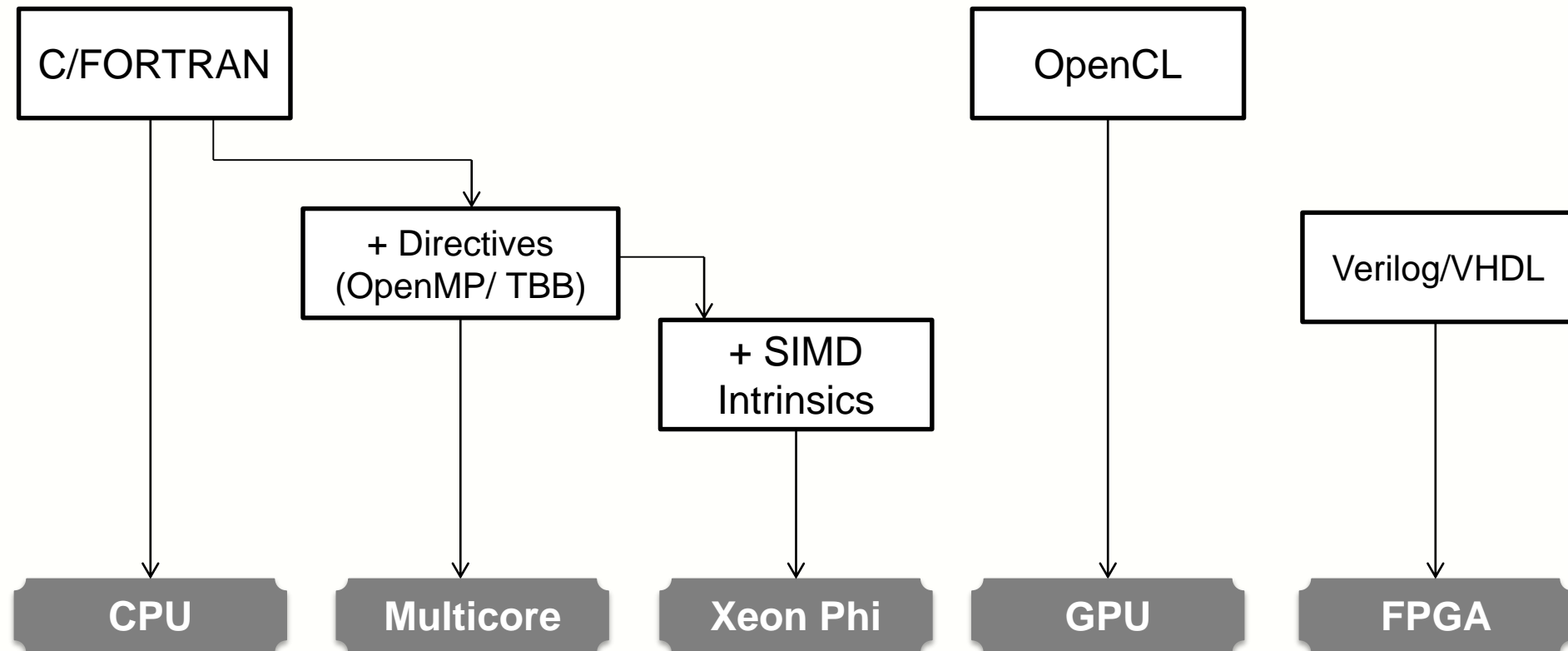
# Current State of Programming Heterogeneous Systems





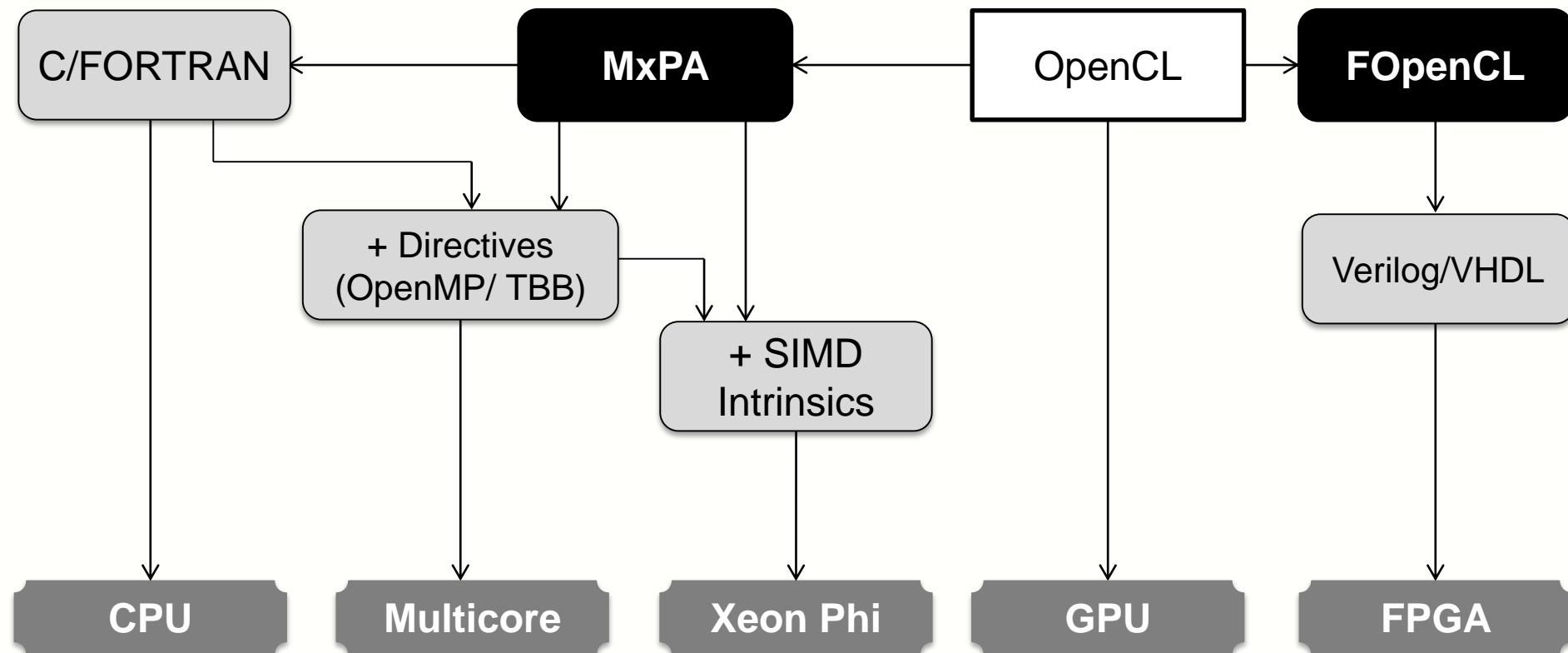
# Current State of Programming Heterogeneous Systems

Programming heterogeneous systems requires too many versions of code!



# Productivity in Programming Heterogeneous Systems

**Step #1: Keep only one of the versions and use portability tools to generate the others.**

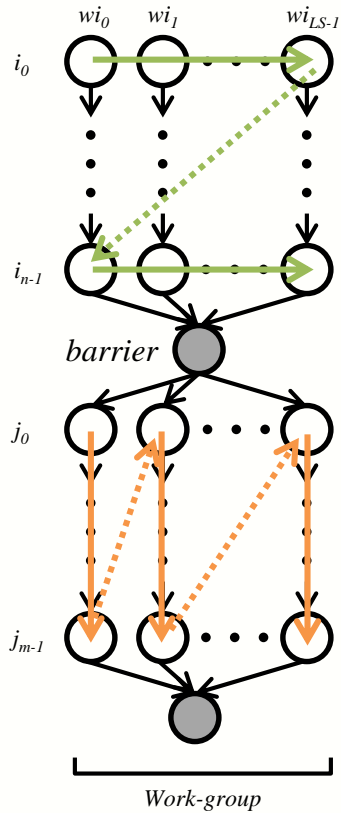


# MxPA: Overview

- Optimizes scheduling of work-item execution for **locality** and **vectorization** for the hardware
- Locality-centric compile-time scheduling selects between:
  - *BFO (Breadth-First Order)*: issue each memory access for all work-items first
  - *DFO (Depth-First Order)*: issue all memory accesses for a single work-item first
- Kernel fusion run-time scheduling reduces memory traffic for common data flow between kernels
- Dynamic vectorization maintains vector execution in spite of apparent control flow divergence



# MxPA: Locality-centric Scheduling



Dependency in executing OpenCL kernels

```
CLKernel() {
  // e.g. SOA
  for (i = 0..N) {
    .. = A[c0*i + wid];
  }
  barrier();
  // e.g. AOS
  for (j = 0..M) {
    .. = B[c1*wid + j];
  }
}
```

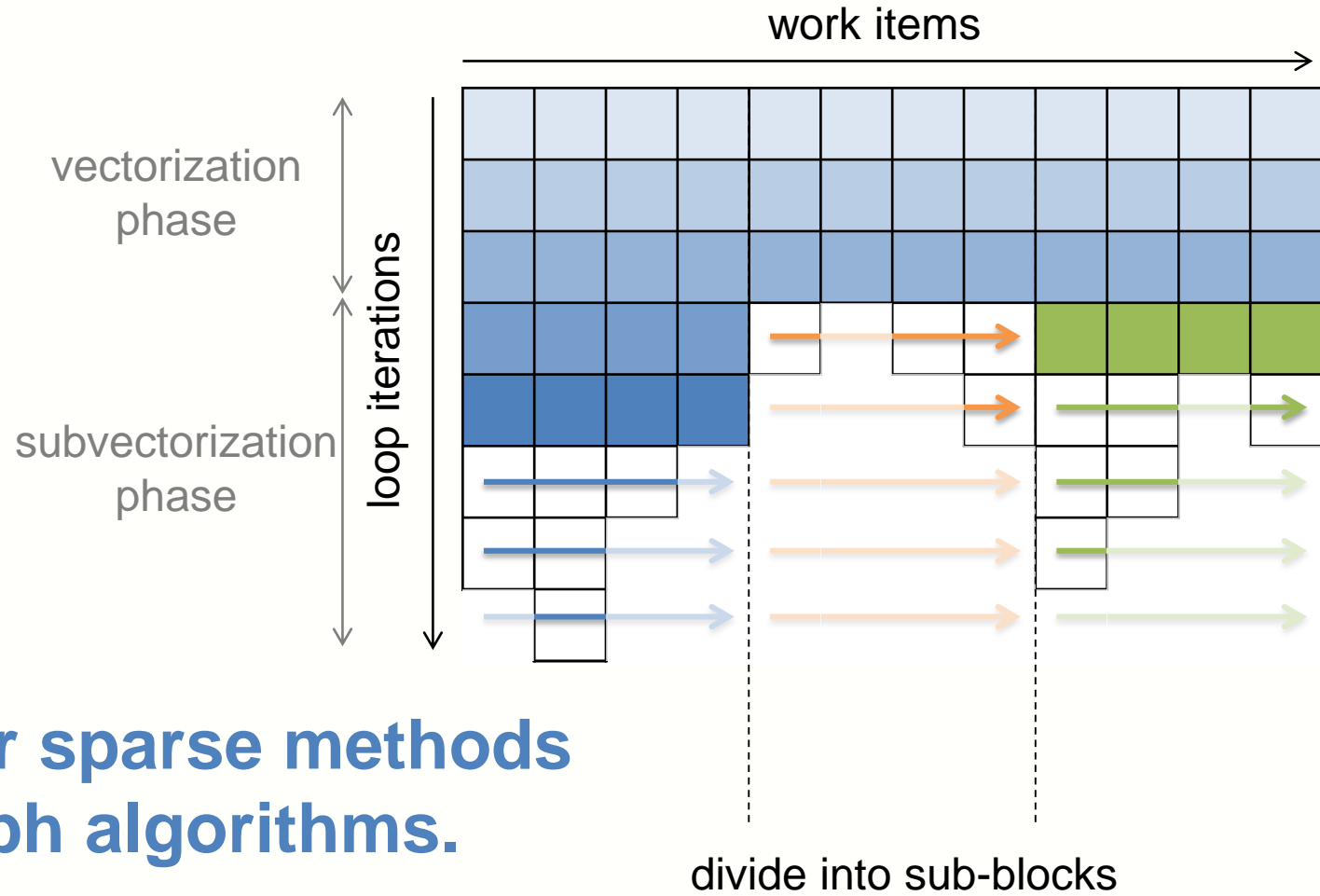


```
CLKernel_MXPA() {
  // BFO
  for (i = 0..N) {
    for (wid = 0..LS) {
      .. = A[c0*i + wid];
    }
  }
  // DFO
  for (wid = 0..LS) {
    for (j = 0..M) {
      .. = B[c1*wid + j];
    }
  }
}
```

An example OpenCL code

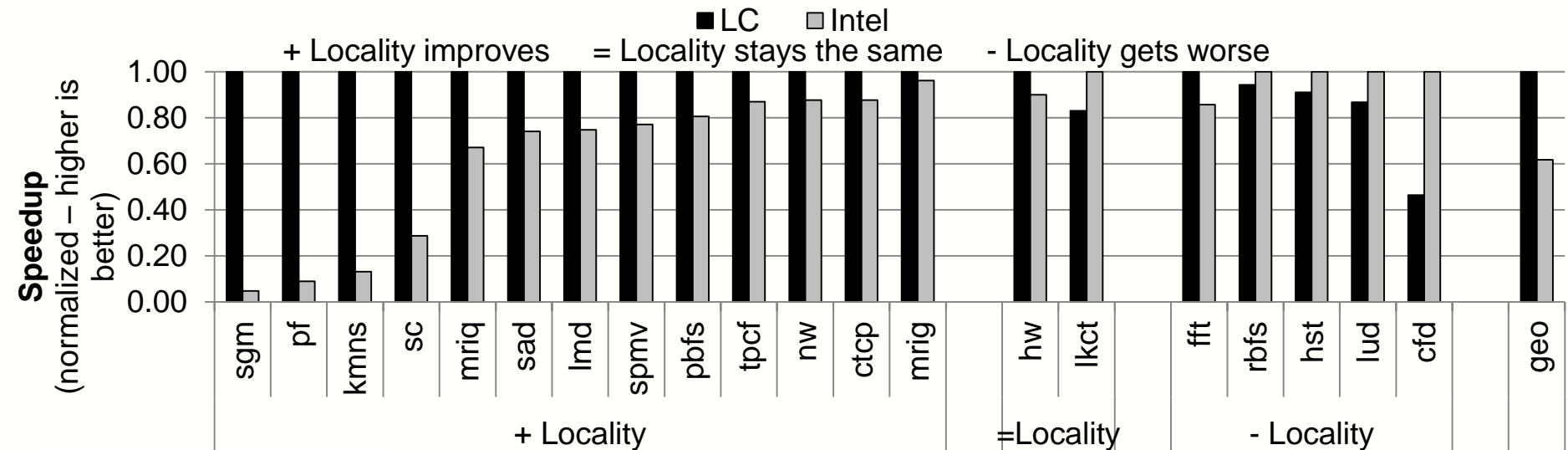
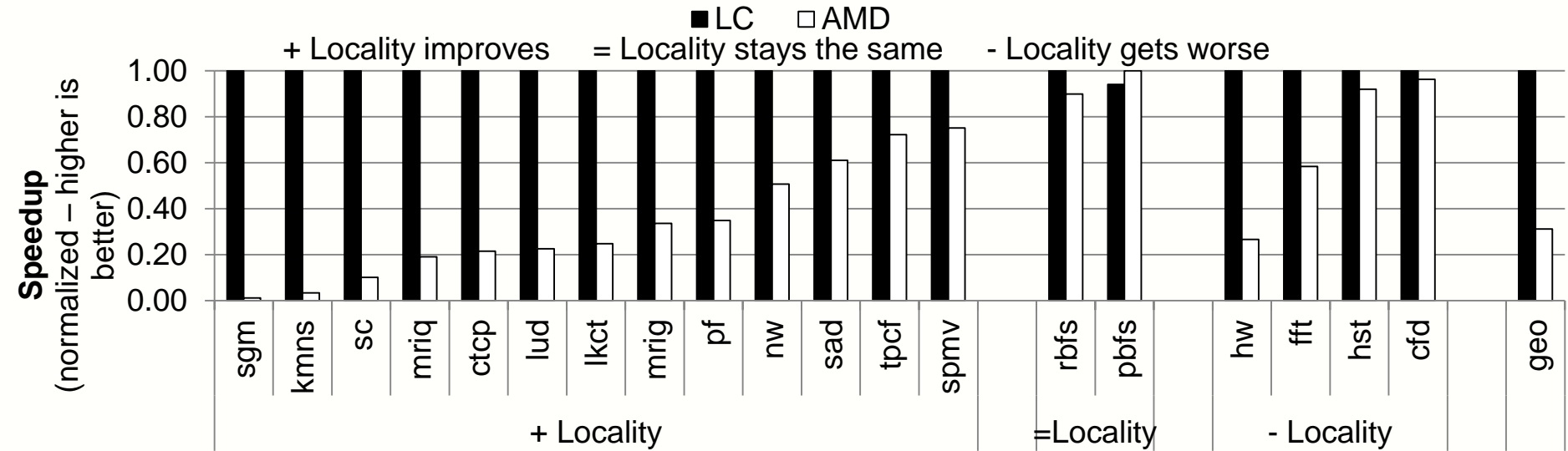
MXPA-translated code

# MxPA: Dynamic Vectorization of Control Divergent Loops



Effective for sparse methods  
and graph algorithms.

# MxPA Results: Comparison to Industry



# MxPA Results: Comparison to Industry

Metric	LC/AMD	LC/Intel
Speedup	3.20x	1.62x
L1 Data Cache Misses	0.12x	0.36x
Data TLB Misses	0.23x	0.33x
LLC Misses	0.91x	0.97x

# HIGH-LEVEL INTERFACE



# Who does the hard work in parallelization?

- General-purpose language + parallelizing compiler
  - Requires a very intelligent compiler
  - Limited success outside of regular array algorithms
- Delite - Domain-specific language + domain-specific compiler
  - Simplify compiler's job with language restrictions and extensions
  - Requires customizing a compiler for each domain
- Triolet - Parallel library + optimizing compiler
  - Library makes parallelization decisions
  - Uses a rich transformation, library aware compiler
  - Extensible—just add library functions

