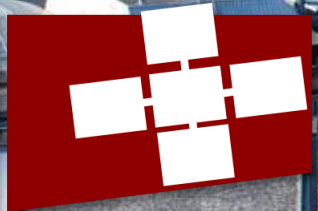


Tal Ben-Nun, Johannes de Fine Licht, Alexandros-Nikolaos Ziogas, Timo Schneider, Torsten Hoefler

# Stateful Dataflow Multigraphs: A Data-Centric Model for Performance Portability on Heterogeneous Architectures

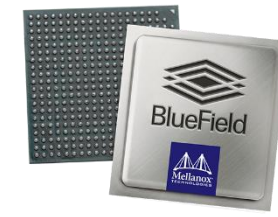


This project has received funding from the European Research Council (ERC) under grant agreement "DAPP (PI: T. Hoefler)".





# Motivation



# Computational Scientist

**Domain Scientist**

**Performance Engineer**

# Optimization Techniques

## ■ Multi-core CPU

- Tiling for complex cache hierarchies
- Register optimizations
- Vectorization

## ■ Many-core GPU

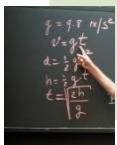
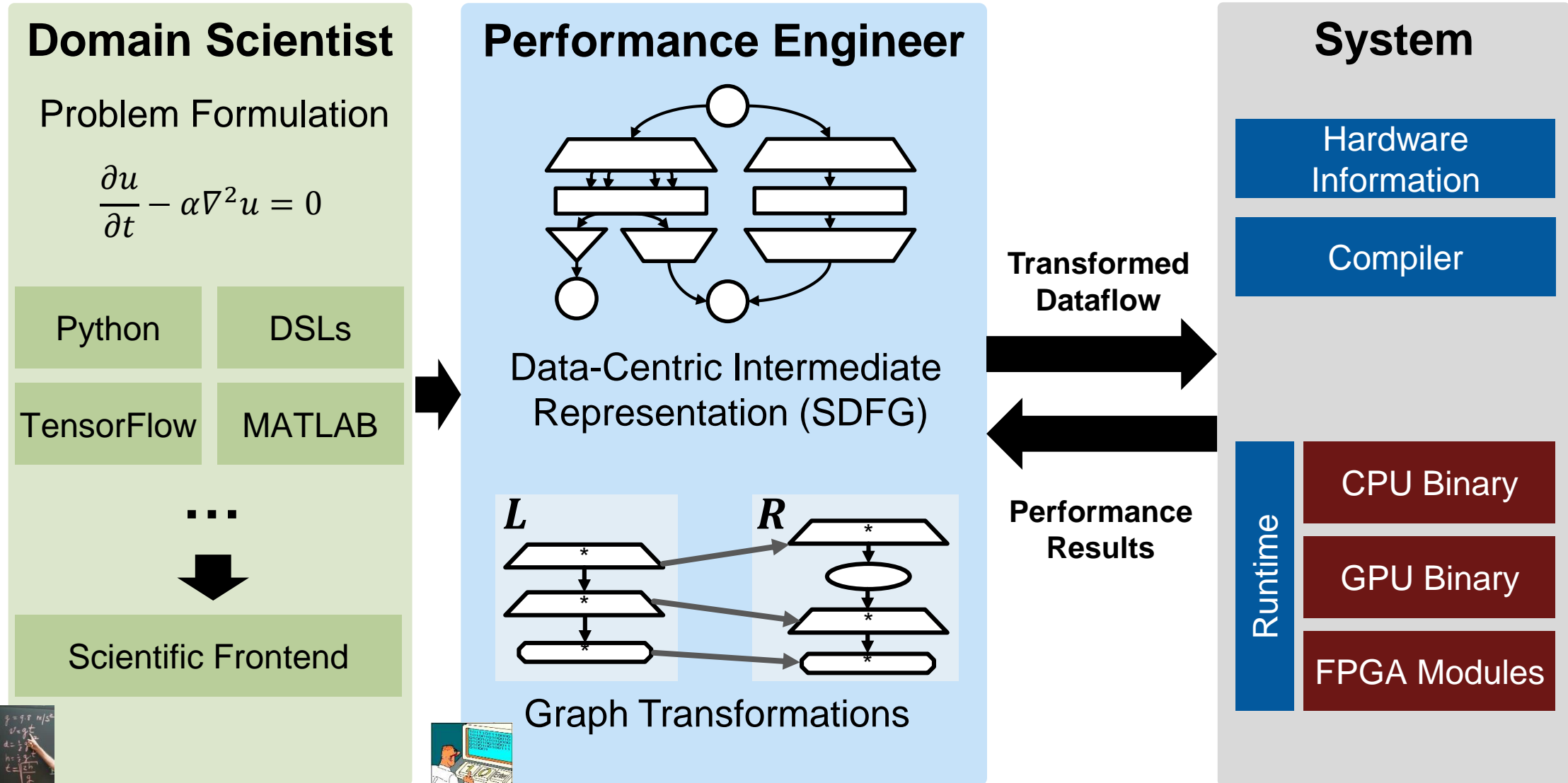
- Coalesced memory access
- Warp divergence minimization, register tiling
- Task fusion

## ■ FPGA

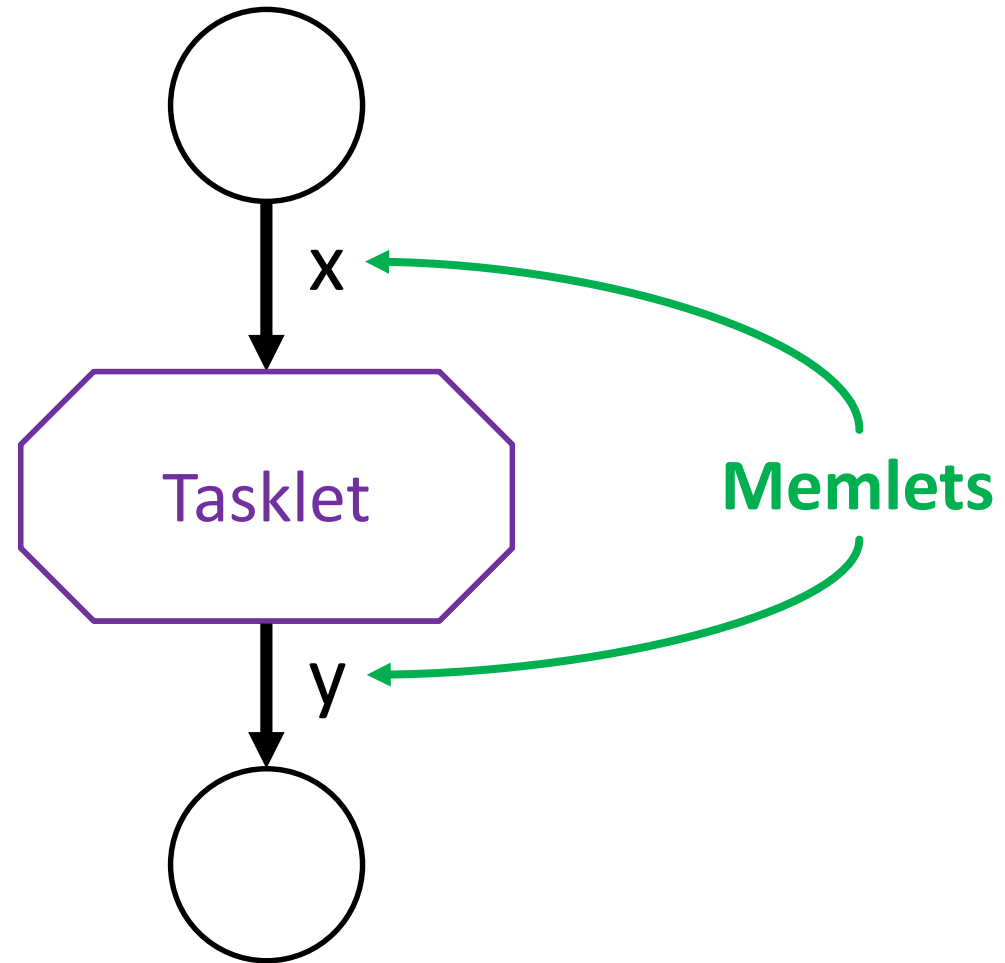
- Maximize resource utilization (logic units, DSPs)
- Streaming optimizations, pipelining
- Explicit buffering (FIFO) and wiring



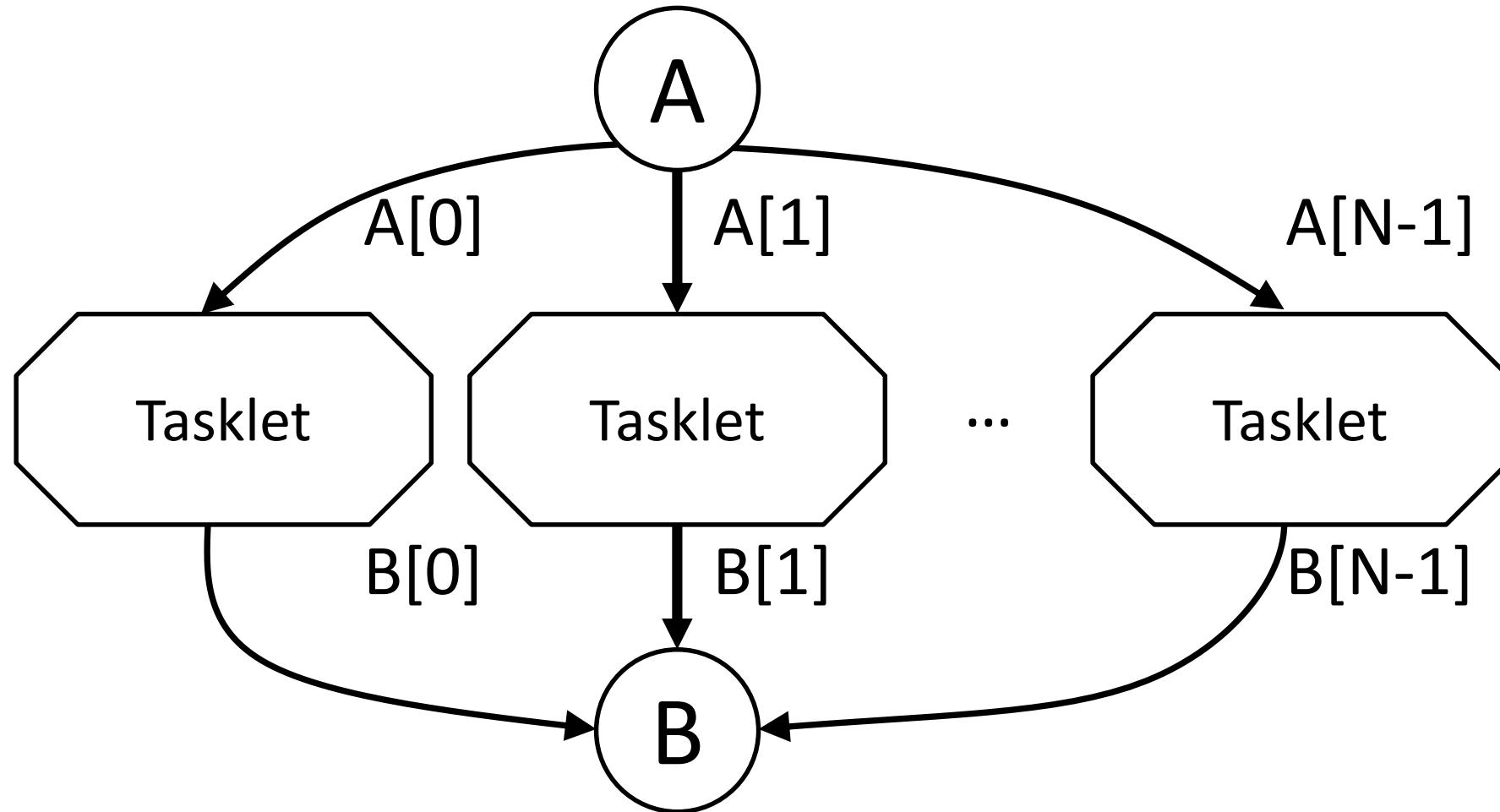
# DaCe Overview



# Dataflow Programming in DaCe

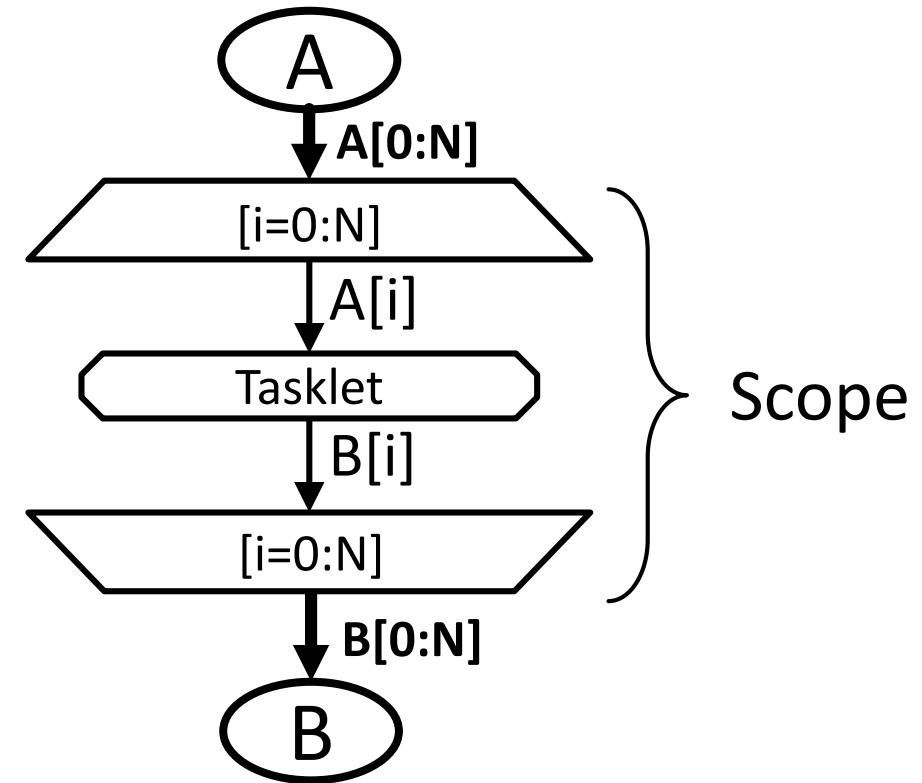
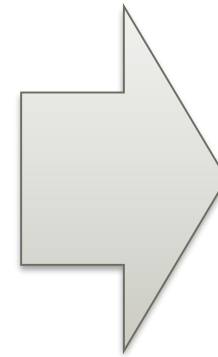
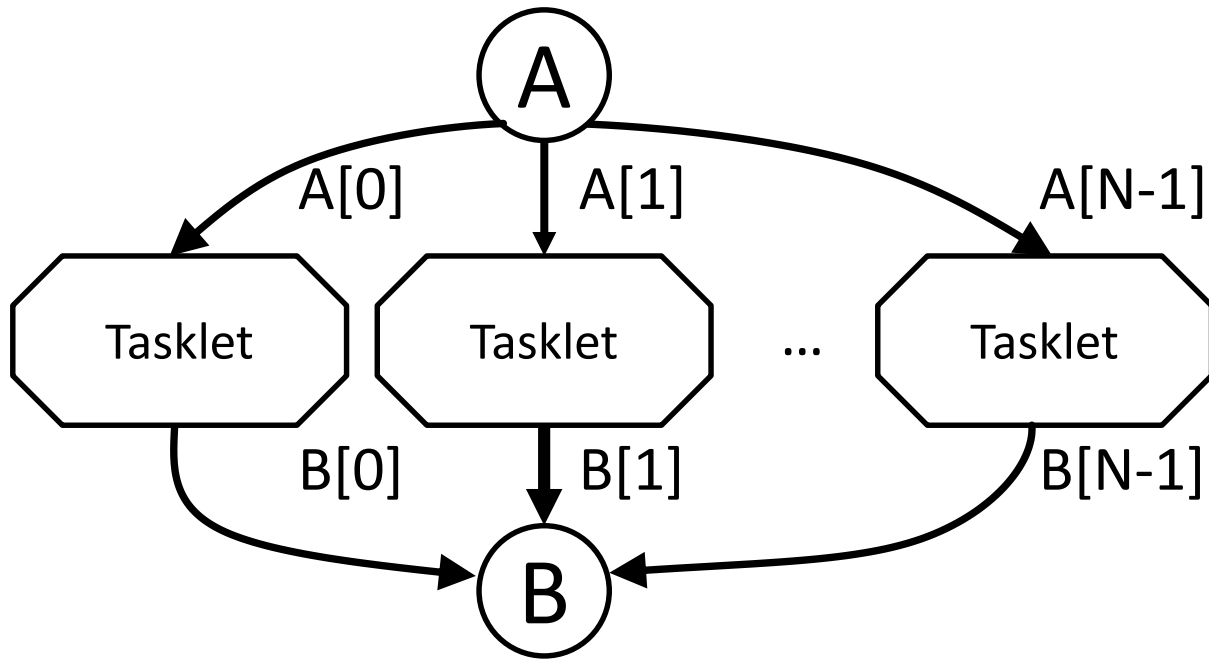


# Parallel Dataflow Programming

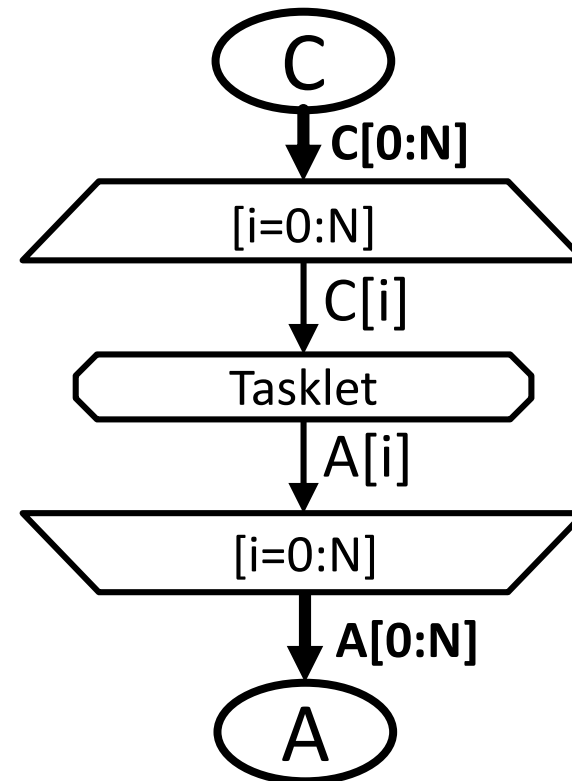
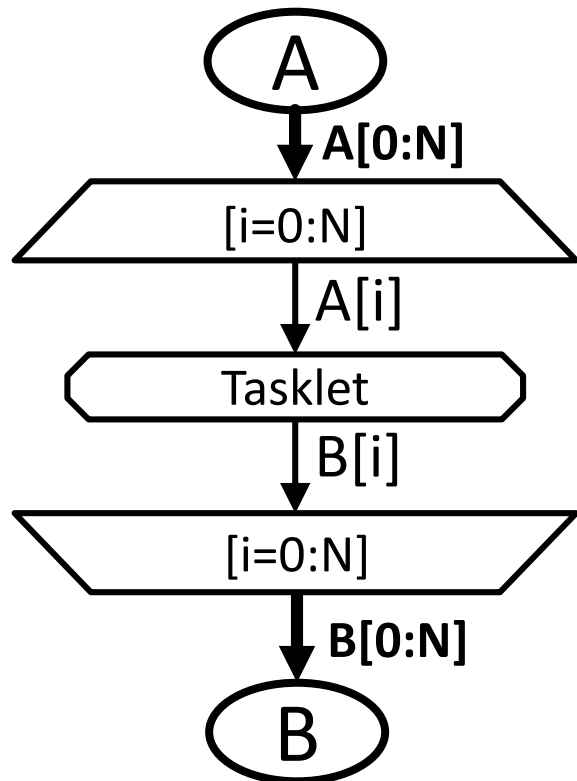




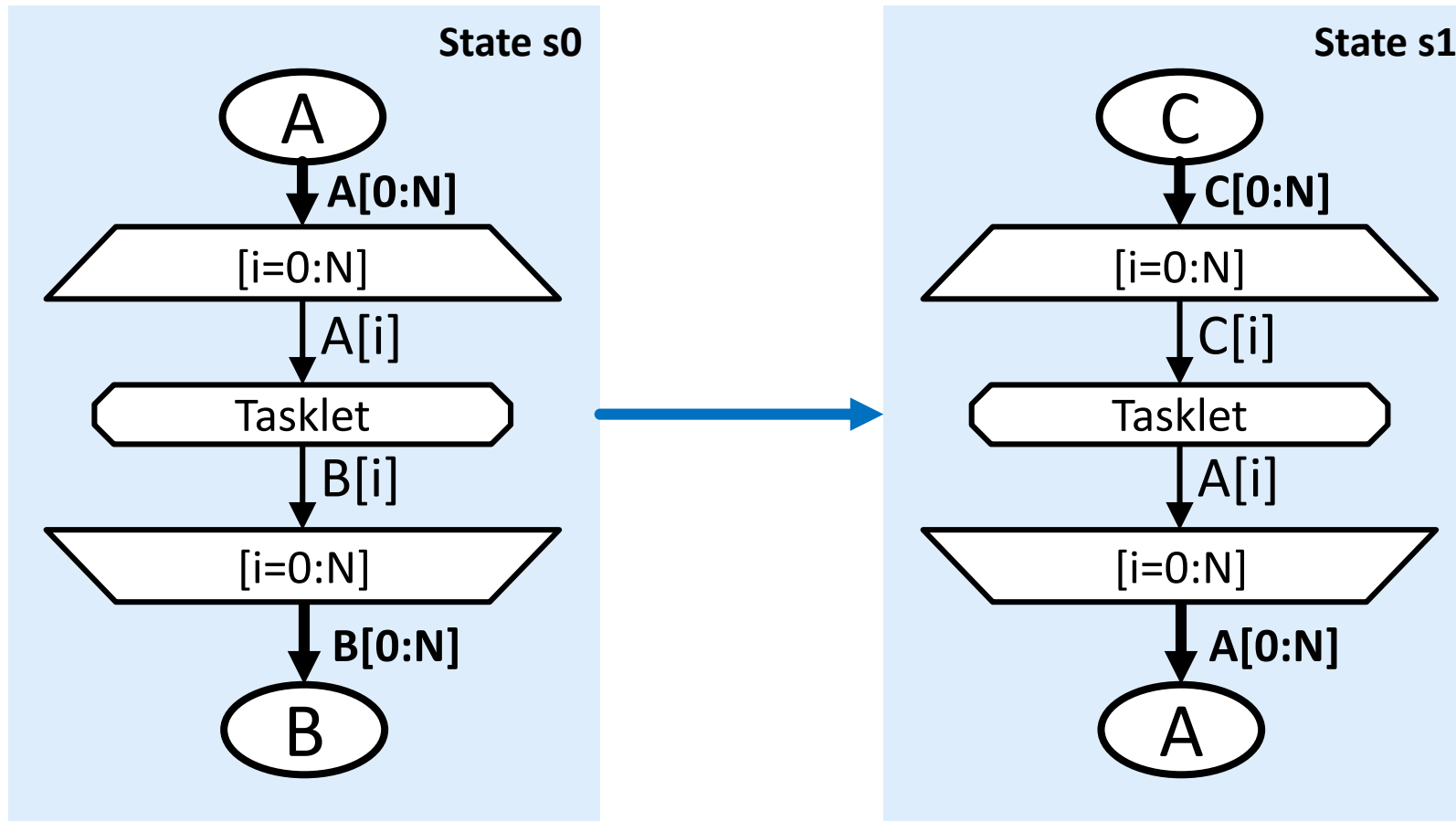
# Parallel Dataflow Programming



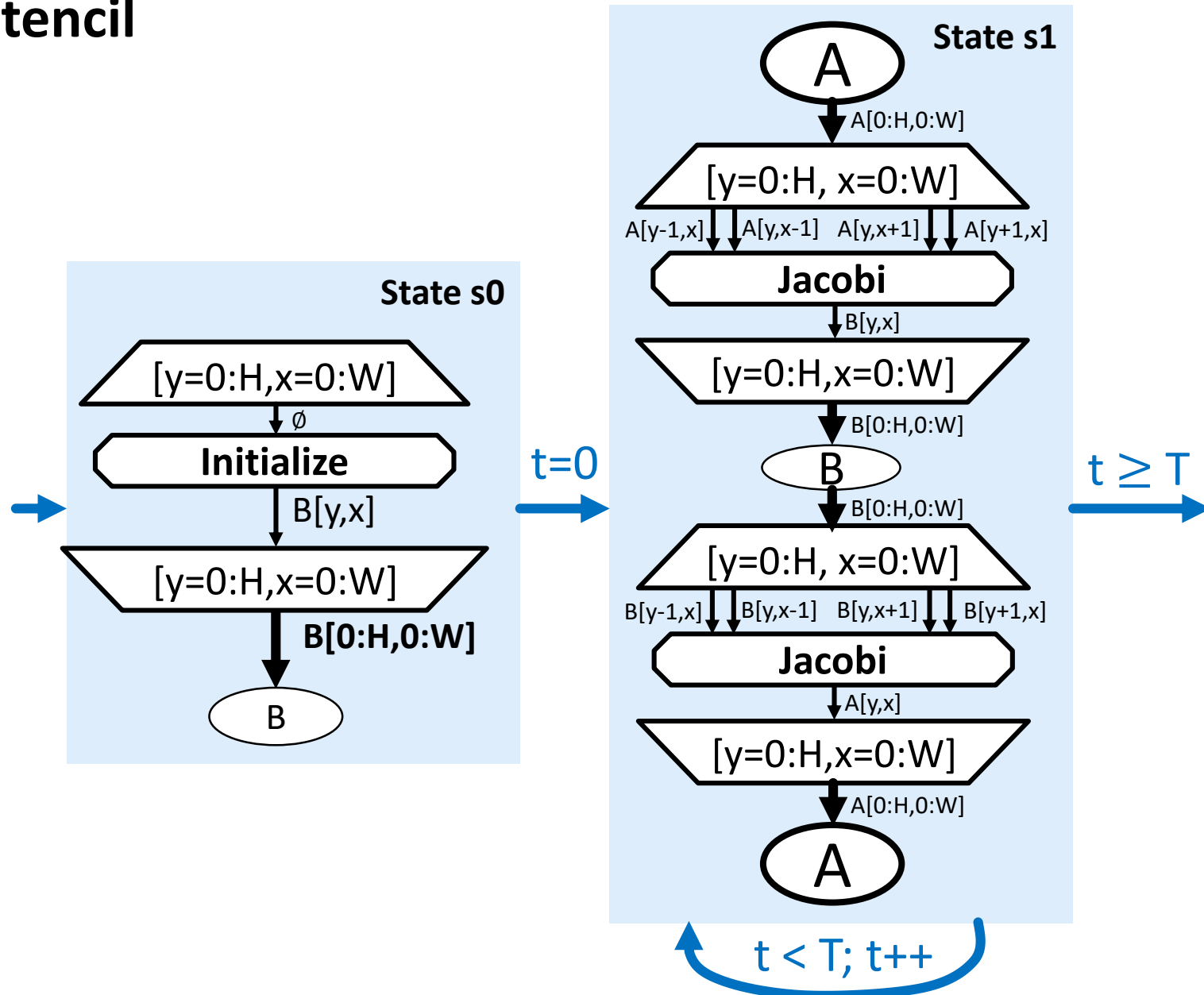
# Stateful Parallel Dataflow Programming



# Stateful Parallel Dataflow Programming



# Example: 2D Stencil



# Meet the Nodes

**State**

State machine element

**Tasklet**

Fine-grained computational block

**Array**

N-dimensional data container

**Map**

**Exit**

Parametric graph abstraction for parallelism

**Stream**

Streaming data container

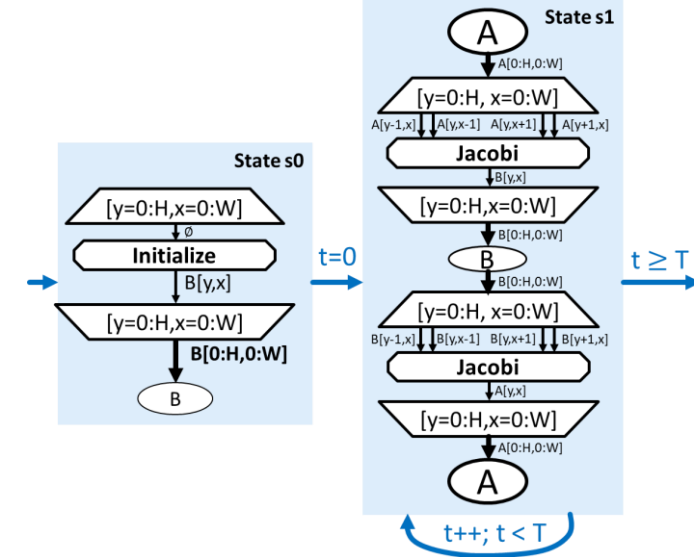
**Consume**

**Exit**

Dynamic mapping of computations on streams

Conflict Resolution

Defines behavior during conflicting writes





# Meet the Nodes

## Stateful Dataflow Multigraphs

form *symbol = expression*. Once a state finishes execution, all outgoing state transitions of that state are evaluated in an arbitrary order, and the destination of the first transition whose condition is true is the next state which will be executed. If no transition evaluates to true, the program terminates. Before starting the execution of the next state, all assignments are performed, and the left-hand side of assignments become symbols.

### A.2 Operational Semantics

**A.2.1 Initialization.** *Notation* We denote collections (sets/lists) as capital letters and their members with the corresponding lowercase letter and a subscript, i.e., in an SDFG  $G = (S, T, s_0)$  the set of states  $S$  as  $s_i$ , with  $0 \leq i < |S|$ . Without loss of generality we assume  $s_0$  to be the start state. We denote the value stored at memory location  $a$  as  $M[a]$ , and assume all basic types are size-one elements to simplify address calculations.

The state of execution is denoted by  $\rho$ . Within the state we carry several sets: *loc*, which maps names of data nodes and transients to memory addresses; *sym*, which maps symbol names (identifiers) to their current value; and *vis*, which maps connectors to the data visible at that connector in the current state of execution.

We define a helper function *size()*, which returns the product of all dimensions of the data node or element given as argument (using  $\rho$  to resolve symbolic values). Furthermore, *id()* returns the *name* property of a data or transient node, and *offs()* the offset of a data element relative to the start of the memory region it is stored in. The function *copy()* creates a copy of the object given as argument, i.e., when we modify the copy, the original object remains the same.

**Invocation** When an SDFG  $G$  is called with the data arguments  $A \equiv [a_i = p_i]$  ( $a_i$  is an identifier,  $p_i$  is an address/pointer) and symbol arguments  $Z \equiv [z_i = v_i]$  ( $z_i$  is an identifier,  $v_i$  is an integer) we initialize the configuration  $\rho$ :

- (1) For all symbols  $z_i$  in  $Z$ :  $sym[z_i] \leftarrow v_i$ .
- (2) For all data and stream nodes  $d_i \in G$  without incoming edges s.t.  $id(d_i) = a_i$ :  $loc(d_i) \leftarrow p_i$ ,  $vis(d_i.data) \leftarrow M[p_i, \dots, p_i + size(d_i)]$ .
- (3) Set *current* to a copy of the start state of  $G$ ,  $s_0$ .
- (4) Set *state* to  $id(s_0)$ .
- (5) Set  $qsize[f_i]$  to zero for all stream nodes  $f_i \in G$ .

This can be expressed as the following rule:

```

G = (S, T),
start_state(s_0)
D: data nodes in G,
F: stream nodes in G
-----
(InitG, A, Z), \rho \leftarrow \rho
state \leftarrow id(s_0)
current \leftarrow copy(s_0)
\forall a_i \in A: loc[a_i] \leftarrow p_i
\forall z_i \in Z: sym[z_i] \leftarrow v_i
\forall d_i \in D: vis[d_i.data] \leftarrow M[p_i, \dots, p_i + size(d_i)]
\forall f_i \in G.S: qsize[f_i] \leftarrow 0
    
```

**A.2.2 Propagating Data in a State.** Execution of a state entails propagating data along edges, governed by the rules defined below.

SC '19, November 17–22, 2019, Denver, CO, USA

**Element Processing** In each step, we take one element  $q$  (either a memlet or a node) of *current*, for which all input connectors have visible data, then:

If  $q$  is a **memlet** (*src*, *dst*, *subset*, *reindex*, *wcr*), update  $vis[dst]$  to  $wcr(reindex(subset(vis[src])))$ :

```

q = memlet(src, dst, subset, reindex, wcr),
(vis[src], \rho) \neq \emptyset,
(wcr(reindex(subset(vis[src])), \rho)) \rightarrow [d_0, \dots, d_n]
(q, \rho) \rightarrow \rho[vis[dst] \leftarrow [d_0, \dots, d_n]]
    
```

If  $q$  is a **data node**, update its referenced memory for an input connector  $c_i$ :

$M[loc(id(q)) \dots loc(id(q)) + size(vis[q.data])] = vis[q.data]$ :

```

q = data[id, dims, bt, transient],
(\forall c_i \in C.in: vis[c_i], \rho) \neq \emptyset,
(\forall c_i \in C.in: vis[c_i], \rho) \rightarrow [d_0^i, \dots, d_{k_i}^i],
\forall d_j^i \in loc(id(q)) + offs[d_i] = d_j^i
(q, \rho) \rightarrow \rho[vis[c_i]: M[d_j^i] \dots d_j^i + size(d_j^i)] = d_j^i
    
```

If  $q$  is a **tasklet**, generate a prologue that allocates local variables for all input connectors  $c_i$  of  $q$ , initialized to  $vis[c_i]$  ( $P_1$ ), as well as output connectors ( $P_2$ ). Generate an epilogue  $Ep$  which updates  $\rho[vis[c_i] \rightarrow v_i]$  for each output connector  $c_i$  of  $q$  with the contents of the appropriate variable (declared in  $P_2$ ). Execute the concatenation of  $(P_1; P_2; code; Ep)$ :

```

q = tasklet(Cin, Cout, code),
(\forall c_i \in Cin: vis[c_i], \rho) \neq \emptyset,
(P_1 = [\forall c_i \in Cin: \text{update}(id(c_i) = vis[c_i]), \rho],
P_2 = [\forall c_i \in Cout: \text{update}(id(c_i) = v_i)],
(Ep = [\forall c_i \in Cout: \text{update}(id(c_i) = v_i)], \rho),
(q, \rho) \rightarrow \rho[exec(P_1; P_2; code; Ep)]
    
```

If  $q$  is a **mapentry** node with range  $y = R$  ( $y$  is the identifier) and scope  $o \in V$ : Remove  $o$  from *current*. Remove  $q$  and the corresponding map exit node from  $o$ . For each element in  $R$ , replicate  $o$ , resolve any occurrence of  $y$  to  $r_i$ , connect incoming connectors of  $q$  and  $p$  in *state*.

```

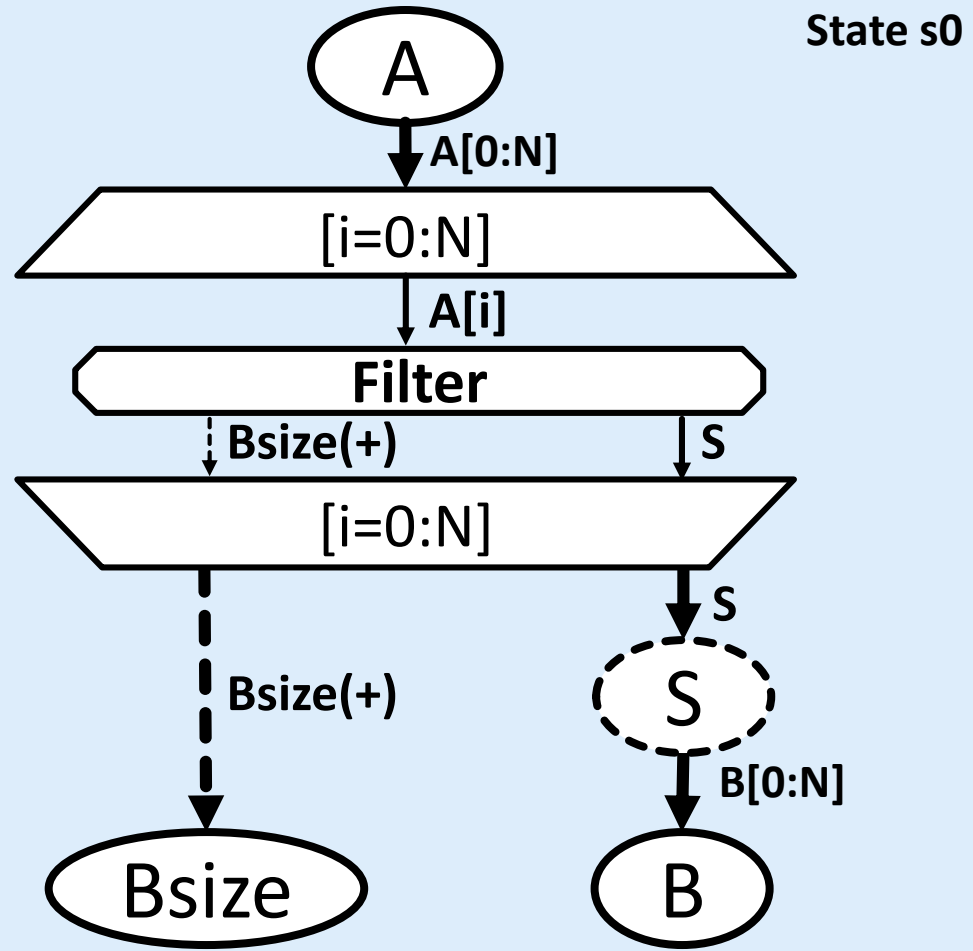
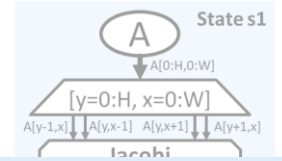
q = mapentry(Cin, Cout, R),
\forall r_i \in R: vis[r_i] \neq \emptyset,
o = scope(q), o' = scope(q) \setminus \{q, memlet(q)\},
NewSym = [r_i \in R: r_i],
(q, \rho) \rightarrow \rho[
current \leftarrow o' \setminus \{r_i\},
\forall r_i \in NewSym: \rho[r_i] \leftarrow vis[r_i]
]
    
```

If  $q$  is a **consume-entry** node, defined by (*range*, *cond*, *cin*, *cout*), replace  $q$  with a mapentry and do the same for the corresponding consume exit node. Then we create a new SDFG *new*, which contains the contents of the original consume scope  $scope(q)$ . *new* consists of one state  $s_0$ , and a single state transition to the same state with the condition *cond*, defined by  $(s_0, s_0, cond, [])$ . Finally, we replace  $scope(q)$  in *current* with an invoke node for *new* and reconnect the appropriate edges between the entry and exit nodes.

```

q = consume(range, cond, cin, cout),
newSDFG = SDFG(scope(q) \setminus \{q, exit(q)\}, (s_0, s_0, cond, []))
\forall v \in invokable(newSDFG)
mem = mapentry(range, cin, cout)
mem = mapexit(cond, cin, cout, mem)
(q, \rho) \rightarrow \rho[
current \leftarrow current \setminus \{q, exit(q)\} \cup \{mem, mem\}
]
    
```

If  $q$  is a **reduce** node defined by the tuple (*cin*, *cout*, *range*), we create a mapentry node *men* with the same range, a mapexit node *mex*, and a tasklet  $o = 1$ . We add these nodes to the node set of



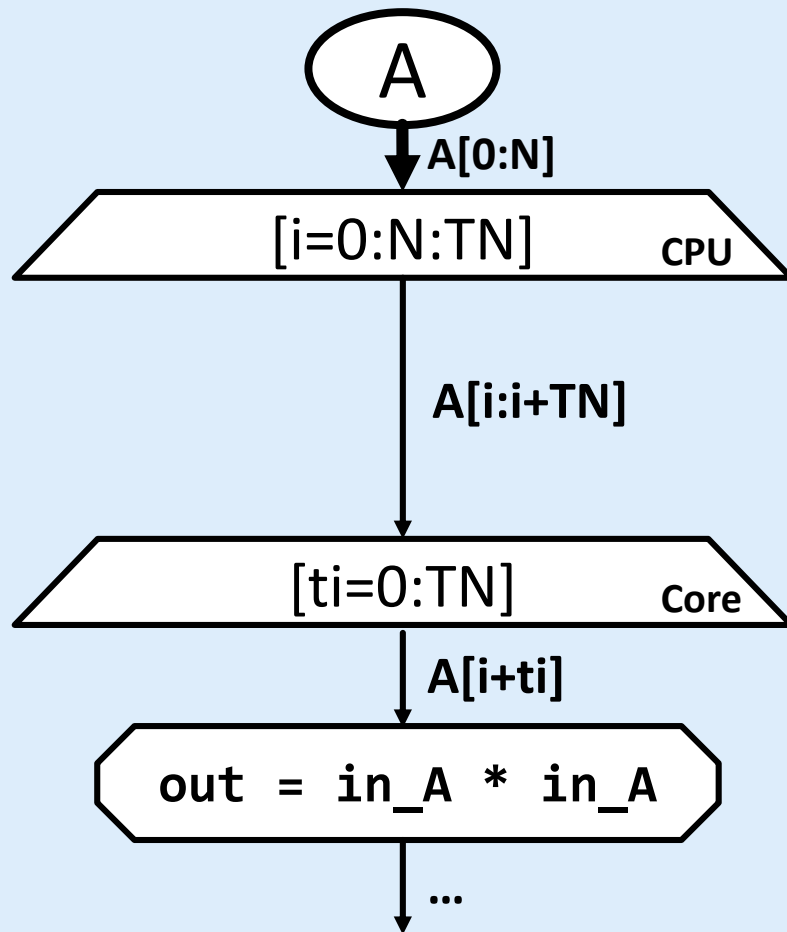
Machine element  
ed computa  
ional data co  
c graph abstr  
g data contain  
mapping of c

Conflict Resolution

Defines behavior during conflicting writes

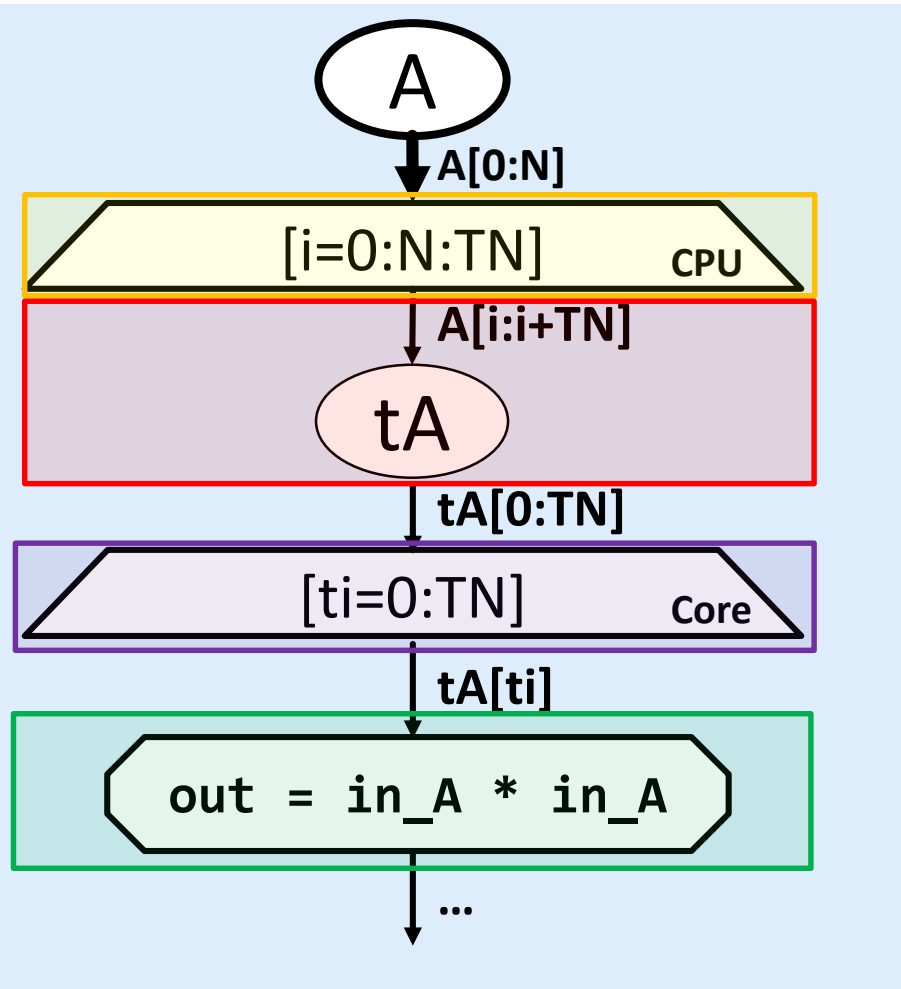
# Hierarchical Parallelism and Heterogeneity

- Maps have schedules, arrays have storage locations



# Hierarchical Parallelism and Heterogeneity

- Maps have schedules, arrays have storage locations

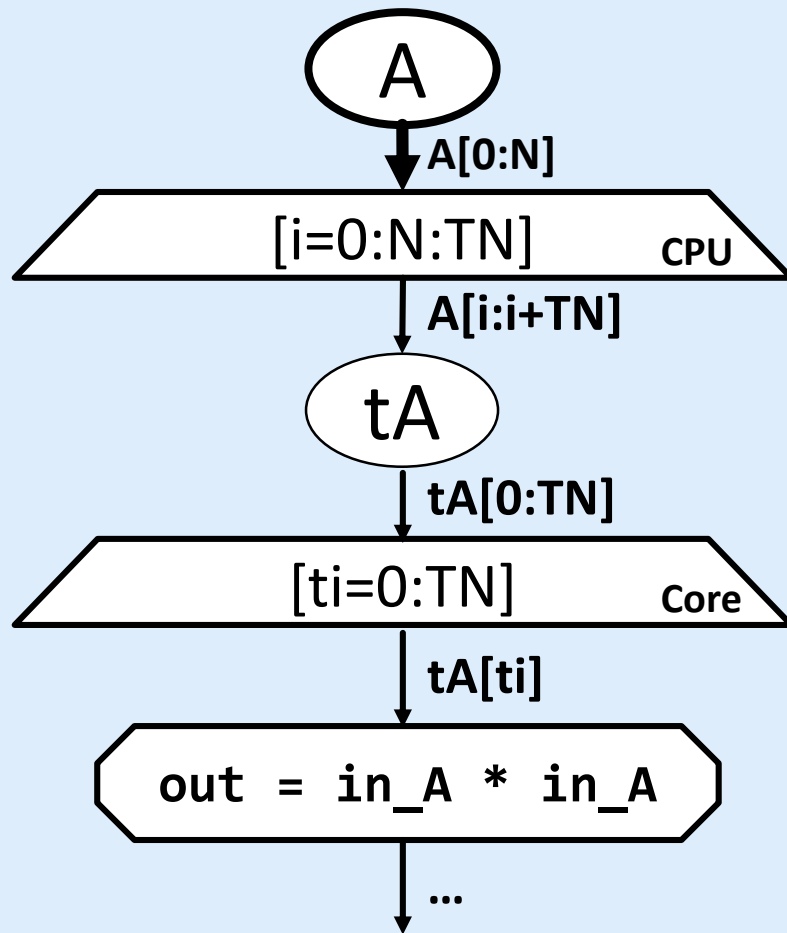


```
// ...
#pragma omp parallel for
for (int i = 0; i < N; i += TN) {
    vec<double, 4> tA[TN];
    Global2Stack_1D<double, 4, 1> (
        &A[i], min(N - i, TN), tA);

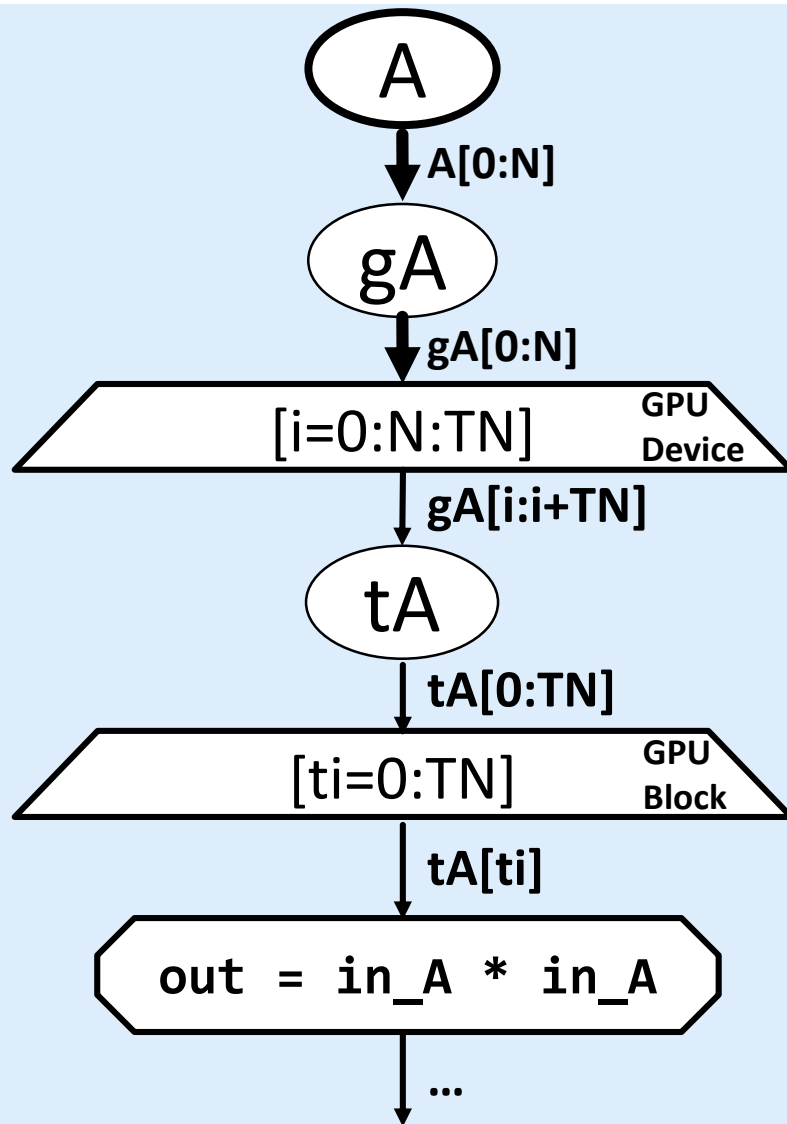
    for (int ti = 0; ti < TN; ti += 1) {

        vec<double, 4> in_A = tA[ti];
        auto out = (in_A * in_A);
        tC[ti] = out;
    }
}
```

# Hierarchical Parallelism and Heterogeneity



# Hierarchical Parallelism and Heterogeneity



```

__global__ void multiplication_1(...) {
    int i = blockIdx.x * TN;
    int ti = threadIdx.y + 0;
    if (i+ti >= N) return;

    __shared__ vec<double, 2> tA[TN];
    GlobalToShared1D<double, 2, TN, 1, 1, false>(gA, tA);

    vec<double, 2> in_A = tA[ti];
    auto out = (in_A * in_A);
    tC[ti] = out;
}
    
```



# Hardware Mapping: Load/Store Architectures

- **Recursive code generation (C++, CUDA)**

**Control flow:** Construct detection and gotos

- **Parallelism**

**Multi-core CPU:** OpenMP, atomics, and threads

**GPU:** CUDA kernels and streams

**Connected components** run concurrently

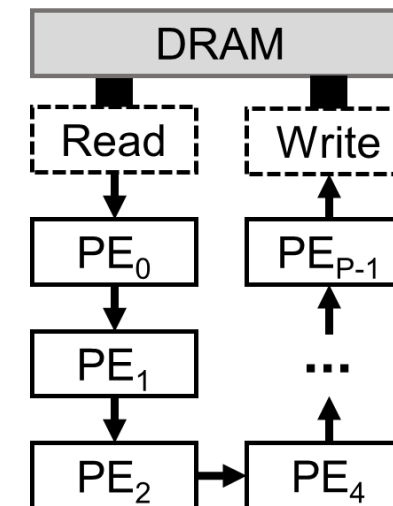
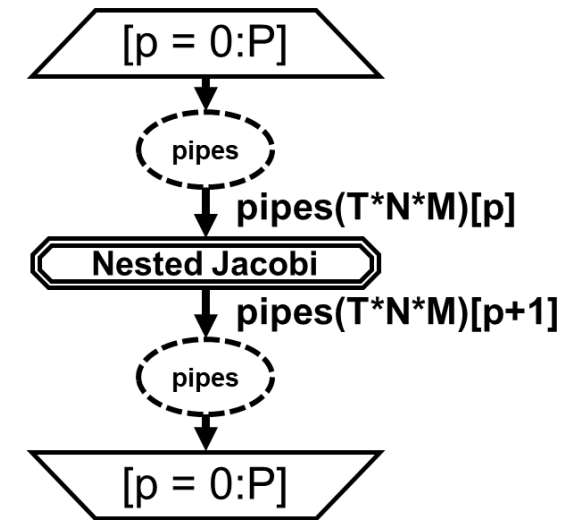
- **Memory and interaction with accelerators**

Array-array edges create intra-/inter-device copies

```
// ...  
#pragma omp parallel for  
for (int i = 0; i < N; i += TN) {  
    vec<double, 4> tA[TN];  
    Global2Stack_1D<double, 4, 1> (  
        &A[i], min(N - i, TN), tA);  
  
    for (int ti = 0; ti < TN; ti += 1) {  
  
        vec<double, 4> in_A = tA[ti];  
        auto out = (in_A * in_A);  
        tC[ti] = out;  
    }  
}
```

# Mapping to Reconfigurable Hardware

- Module generation with HDL and HLS**  
 Xilinx SDAccel  
 Intel FPGA (experimental)
- Parallelism**  
 Exploiting **temporal** locality: pipelines  
 Exploiting **spatial** locality: vectorization, replication
- Replication**  
 Enables parametric systolic array generation



# Data-centric Parallel Programming for Python

- **Programs** are integrated within existing codes

In Python, integrated functions in existing code

In MATLAB, separate .m files

In TensorFlow, takes existing graph

- **In Python: Implicit and Explicit Dataflow**

**Implicit:** numpy syntax

**Explicit:** Enforce **memory access** decoupling from **computation**

- **Output compatible with existing programs**

C-compatible SO/DLL file with autogenerated include file

```
@dace.program
def program_numpy(A, B):
    B[:] = np.transpose(A)
```

```
@dace.program
def program_explicit(A, B):

    @dace.map
    def transpose(i: _[0:N],
                 j: _[0:M]):
        a << A[i,j]
        b >> B[j,i]

    b = a
```

# Matrix Multiplication SDFG

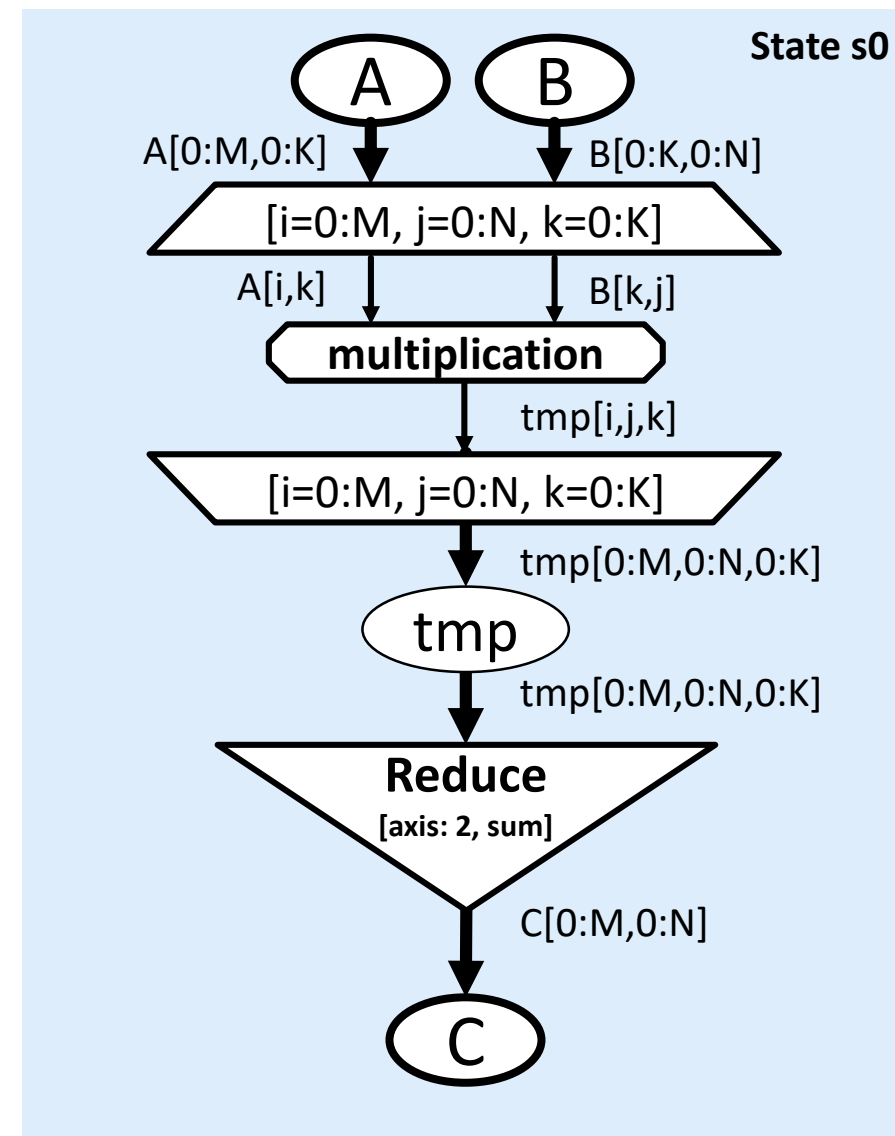
```

@dace.program
def gemm(A: dace.float64[M, K], B: dace.float64[K, N],
        C: dace.float64[M, N]):
    # Transient variable
    tmp = np.ndarray([M, N, K], dtype=A.dtype)

    @dace.map
    def multiplication(i: _[0:M], j: _[0:N], k: _[0:K]):
        in_A << A[i,k]
        in_B << B[k,j]
        out >> tmp[i,j,k]

        out = in_A * in_B

    dace.reduce(lambda a, b: a + b, tmp, C, axis=2)
    
```



# Matrix Multiplication SDFG

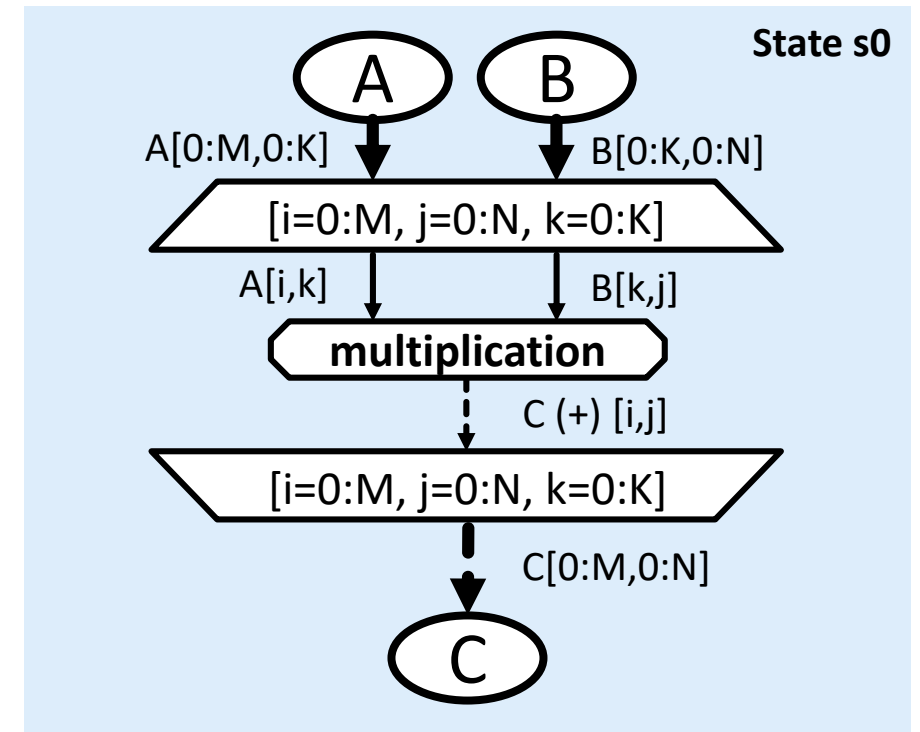
```

@dace.program
def gemm(A: dace.float64[M, K], B: dace.float64[K, N],
        C: dace.float64[M, N]):
    # Transient variable
    tmp = np.ndarray([M, N, K], dtype=A.dtype)

    @dace.map
    def multiplication(i: _[0:M], j: _[0:N], k: _[0:K]):
        in_A << A[i,k]
        in_B << B[k,j]
        out >> tmp[i,j,k]

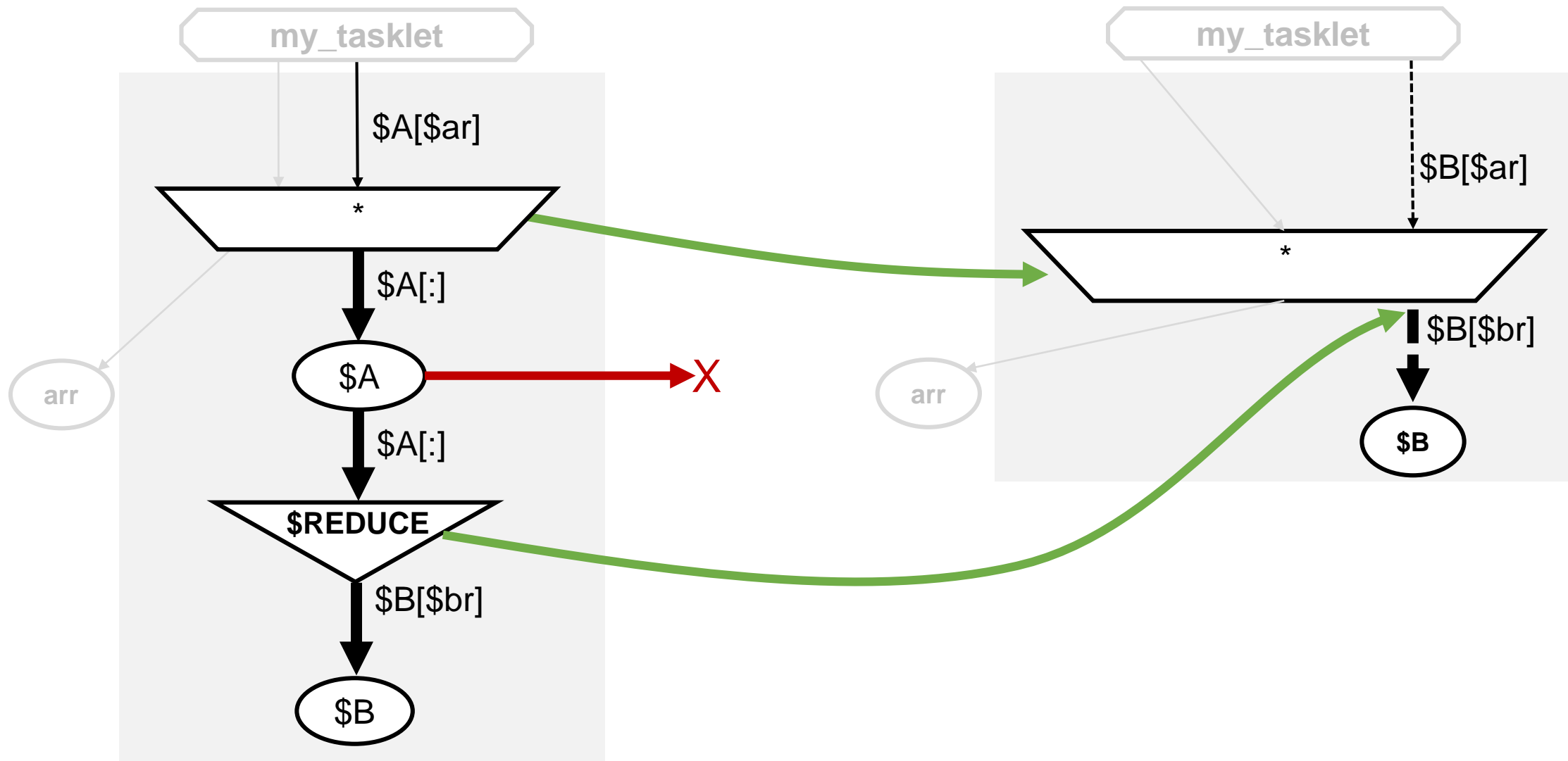
        out = in_A * in_B

    dace.reduce(lambda a, b: a + b, tmp, C, axis=2)
    
```

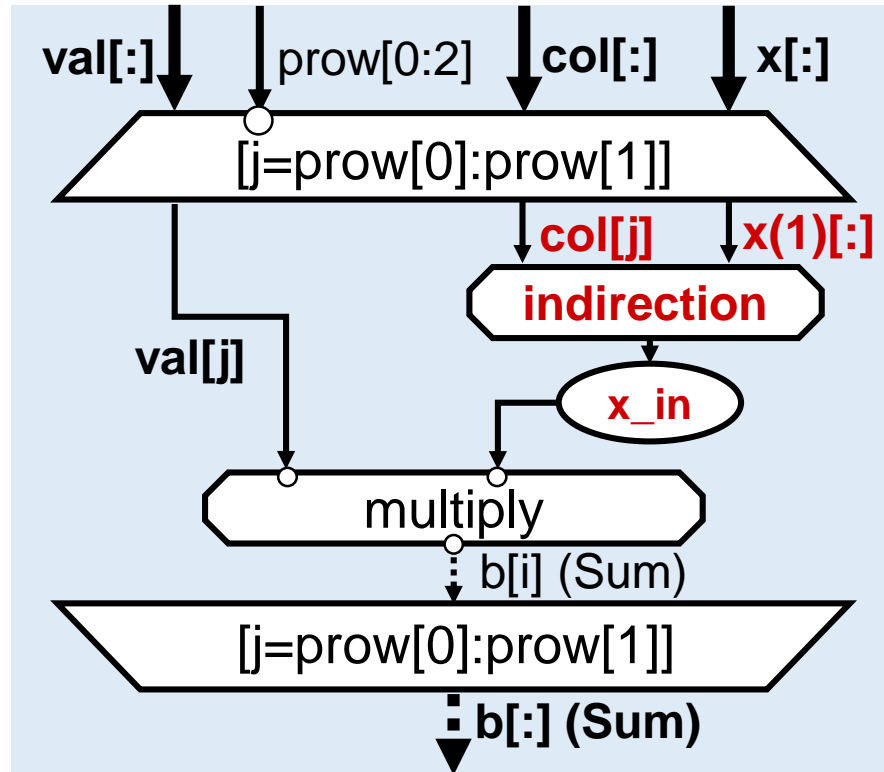




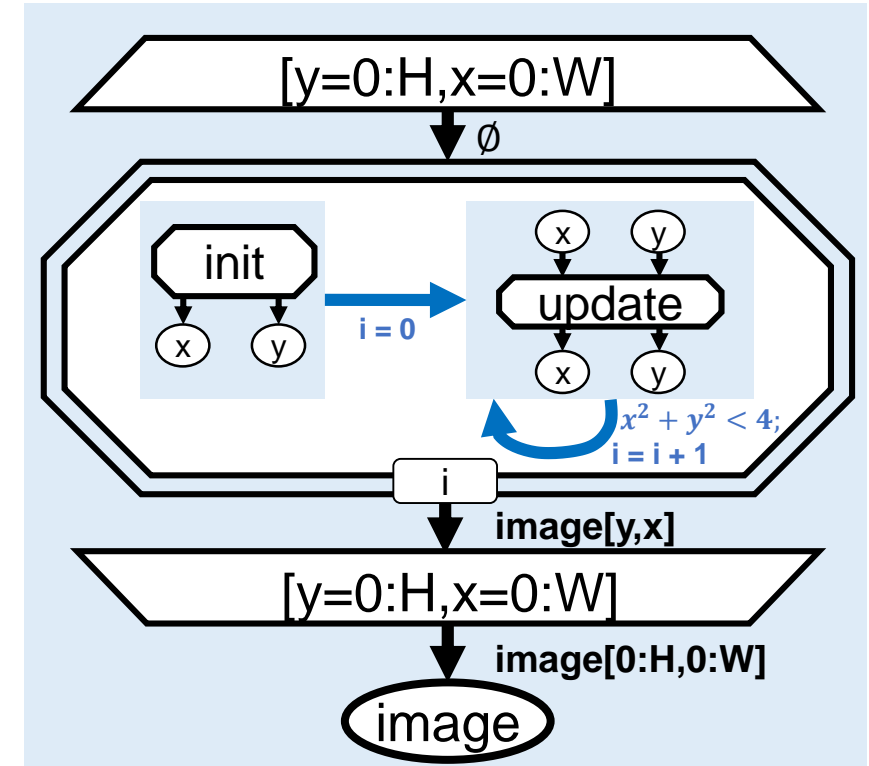
# MapReduceFusion Transformation



# Programming Model Challenges



Indirect memory access



Nested state machines

# DIODE (or: Data-centric Integrated Optimization Development Environment)

The screenshot displays the DIODE IDE interface with three main panels:

- CodeIn (Python):** Shows the high-level DACE program for a Jacobi iteration. It includes imports for `future`, `dace`, `numpy`, and `scipy`. The program defines a `@dace.program` block with a `def jacobi(A, iterations):` function. Inside, it uses `@dace.map` and `@dace.iterate` to perform the iterative computation on a grid of size `H` by `W`.
- CodeOutComponent (C++):** Shows the generated C++ code for the Jacobi iteration, including memory allocation and OpenMP parallelization for the inner loop.
- OptGraph for 'jacobi':** A dataflow graph showing the execution flow. It starts with an input `A`, followed by a `a2b` transformation (represented by a red trapezoid), then a `tmp` node, and finally a `b2a` transformation (represented by a white trapezoid). The graph is annotated with a blue box and arrows, and a message box states "No revertible transformation found."
- Transformation History for 'jacobi':** A list of applied transformations, including `FPGATransformMap`, `FPGATransformState`, `GPUTransformLocalStorage`, `GPUTransformMap`, and `MapExpansion`.
- Properties:** A configuration panel for the selected transformation, showing settings like `debuginfo`, `flatten`, `is_async`, `is_collapsed`, `label` (set to `a2b`), `params` (set to `["y", "x"]`), `range` (set to `Show`), `schedule` (set to `CPU_Multicore`), and `unroll`.

# DIODE (or: Data-centric Integrated Optimization Development Environment)

**Source Code**

```

1 # -*- coding: utf-8 -*-
2 from __future__ import print_function
3
4 import argparse
5 import dace
6 import numpy as np
7 from scipy import ndimage
8
9 W = dace.symbol('W')
10 H = dace.symbol('H')
11 MAXITER = dace.symbol('MAXITER')
12
13
14 @dace.program(dace.float32[H, W], dace.int32)
15 def jacobi(A, iterations):
16     # Transient variable
17     tmp = dace.define_local([H, W], dtype=A.dtype)
18
19     @dace.while_loop(0, iterations)
20     def reset():
21         out >> tmp[y, x]
22         out = 0.0
23
24     @dace.iterate(_[0:iterations])
25     def step(t):
26         @dace.map(_[1:H-1, 1:W-1])
27         def a2b(y, x):
28             in_N << A[y-1, x]
29             in_S << A[y+1, x]
30             in_W << A[y, x-1]
31             in_E << A[y, x+1]
32             in_C << A[y, x]
33             out >> tmp[y, x]
34
35             out = 0.2 * (in_C + in_N + in_S + in_W + in_E)
36
37     # Double buffering
38

```

**SDFG (malleable)**

**Transformations**

- ▲ FPGATransformMap
- FPGATransformMap
- FPGATransformMap\$1
- FPGATransformMap\$2
- FPGATransformSDFG
- ▲ FPGATransformState
- FPGATransformState
- FPGATransformState\$1
- ▲ GPUTransformLocalStorage
- GPUTransformLocalStorage
- GPUTransformLocalStorage\$2
- ▲ GPUTransformMap
- GPUTransformMap
- GPUTransformMap\$1
- GPUTransformMap\$2
- GPUTransformState
- ▲ MapExpansion
- MapExpansion

**Generated Code**

```

11 void __program_jacobi_internal(float * __restrict__ A, int iterations, int H, int
12 {
13     float *tmp = nullptr;
14     cudaMallocHost(&tmp, H * W * sizeof(float));
15     int t;
16     __state_jacobi_reset_tmp:
17     {
18         #pragma omp parallel for
19         for (int y = 1; y < H-1; y++)
20             for (int x = 1; x < W-1; x++)
21                 {
22                     auto __out = dace::ArrayViewOut<float, 0, 1, 1>(tmp + ((W * y
23                     dace::vec<float, 1> out;
24                     ///////////////
25                     // Top-level code (reset_tmp)
26                     out = 0.0;
27                     ///////////////
28                     __out.write(out);
29                 }
30

```

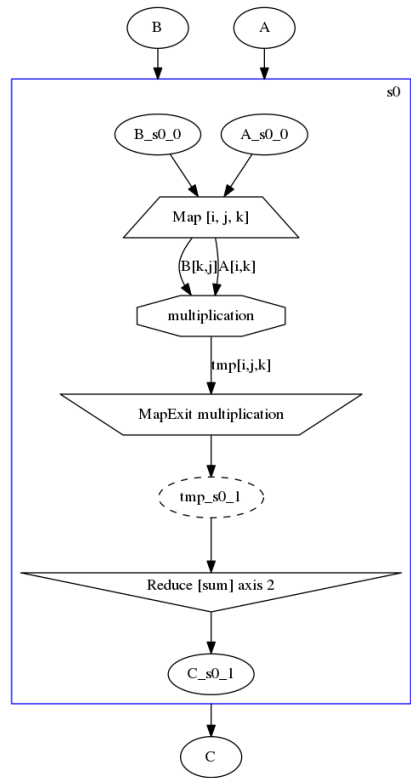
**Transformation History**

No revertible transformation found.

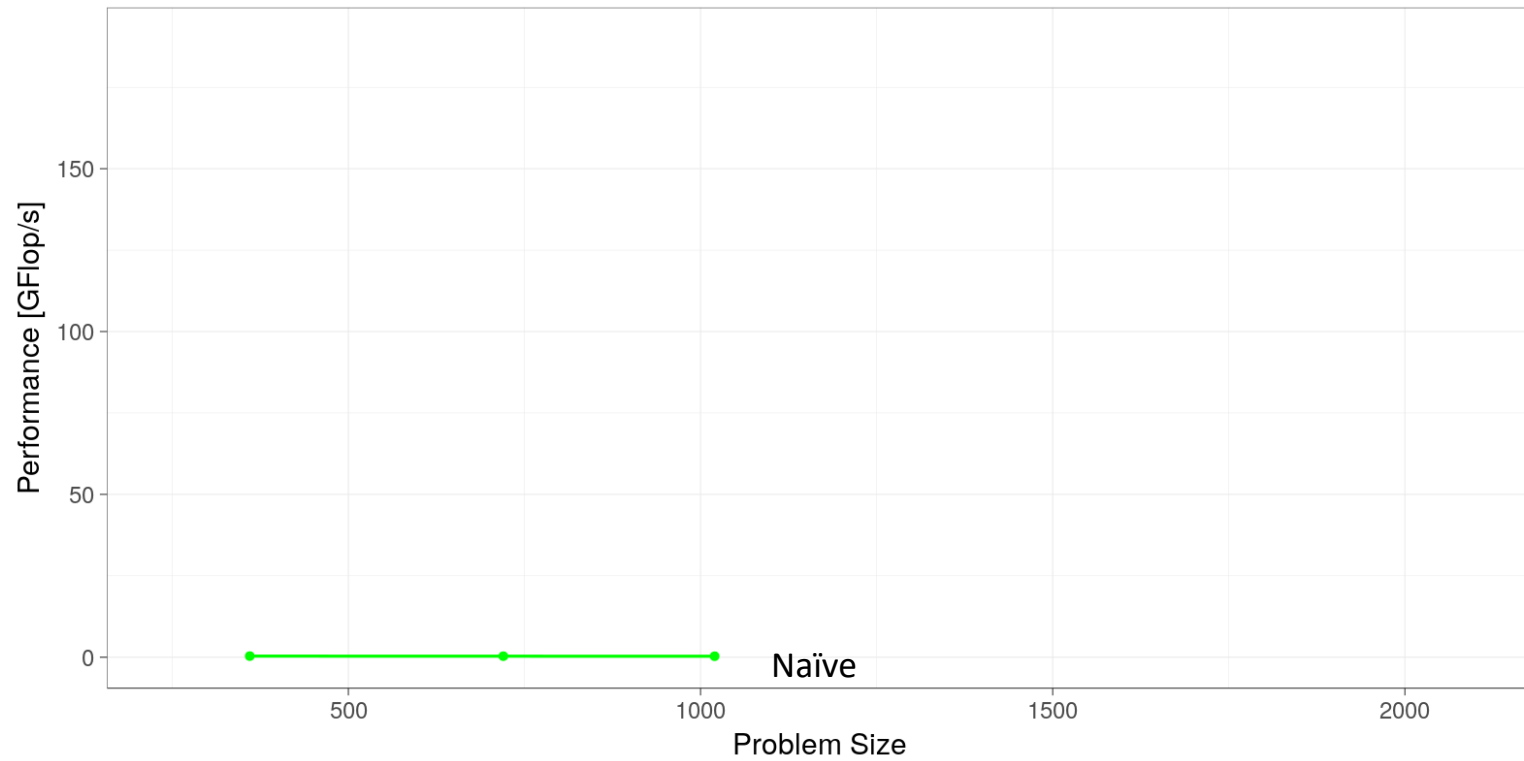
**SDFG Properties**

General	
debuginfo	N/A
fence_instrumentation	<input type="checkbox"/>
flatten	<input type="checkbox"/>
is_async	<input type="checkbox"/>
is_collapsed	<input type="checkbox"/>
label	a2
params	["Y", "x"]
range	Show
schedule	CPU_Multicore
unroll	<input type="checkbox"/>
MapEntry - General	
entry_in_connectors	["IN_1"]
entry_out_connectors	["OUT_1"]

# Performance

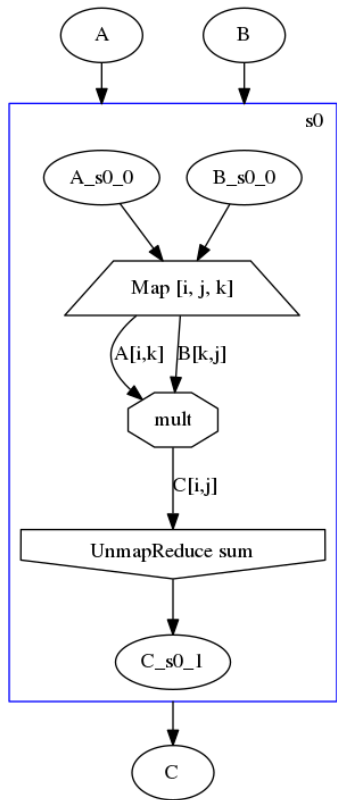


SDFG

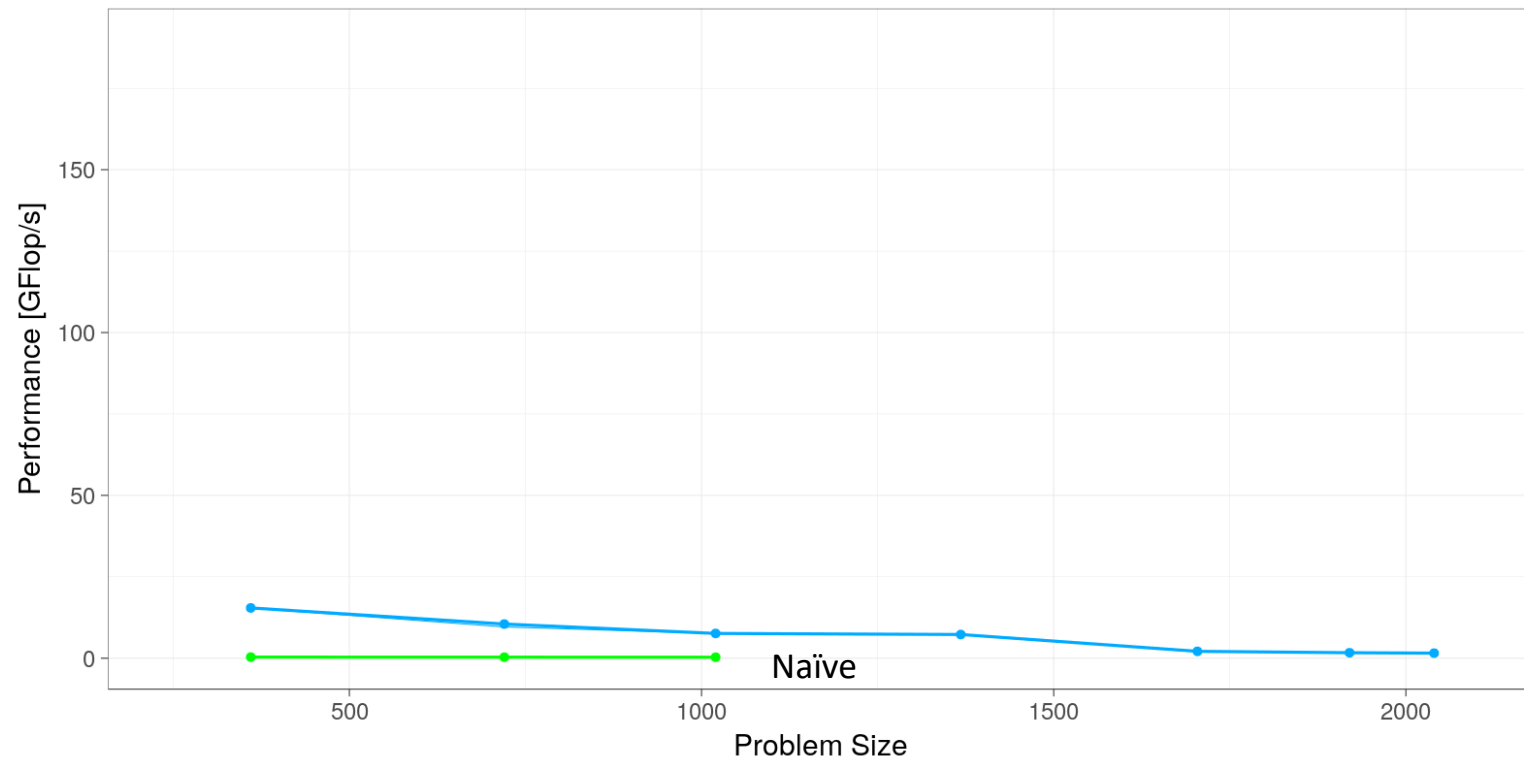




# Performance

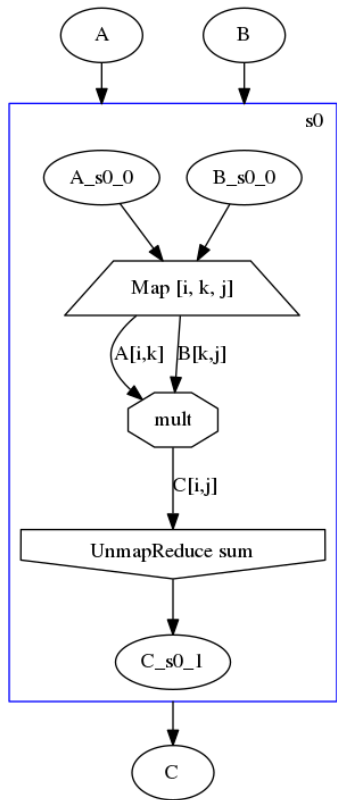


SDFG

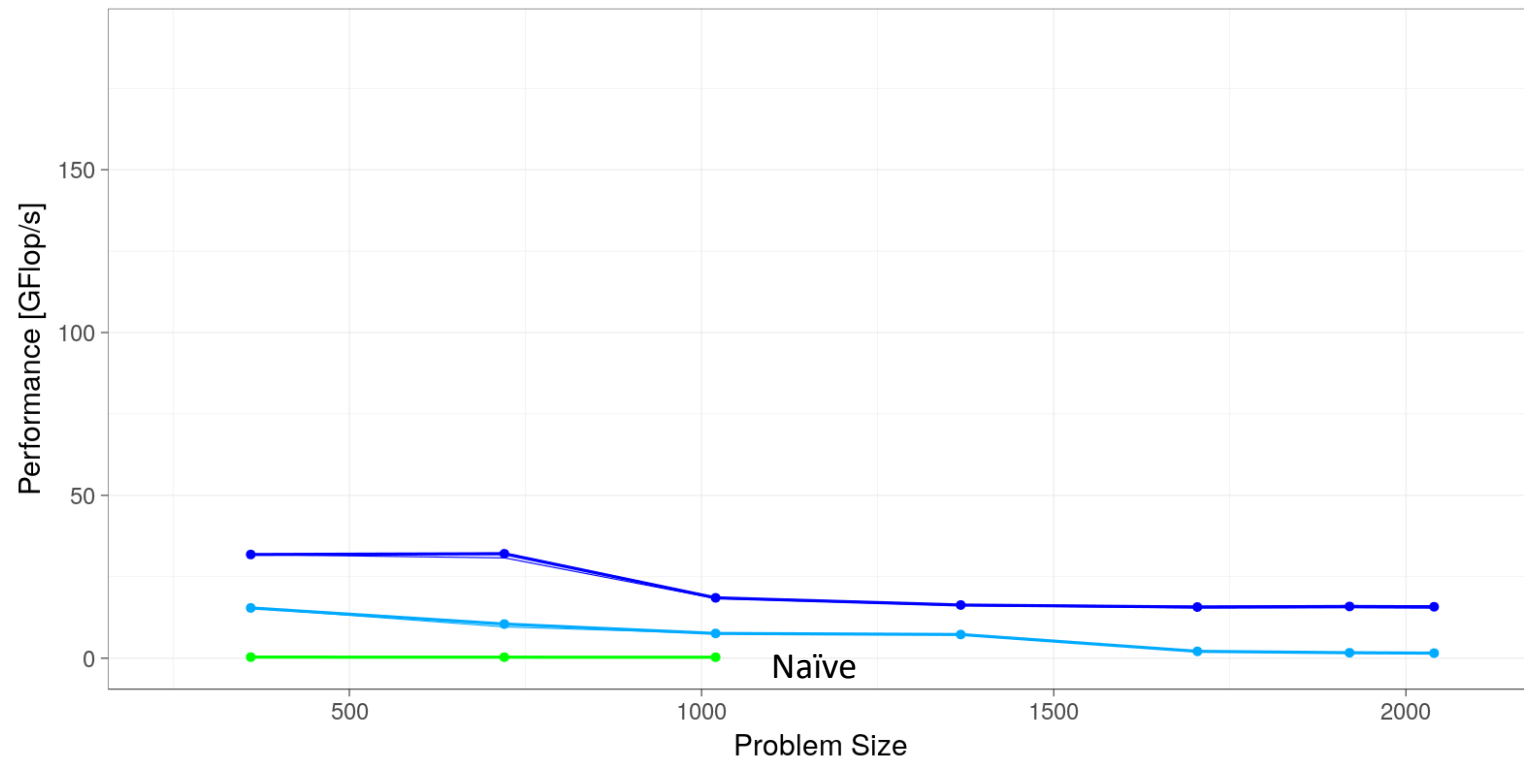


MapReduceFusion

# Performance

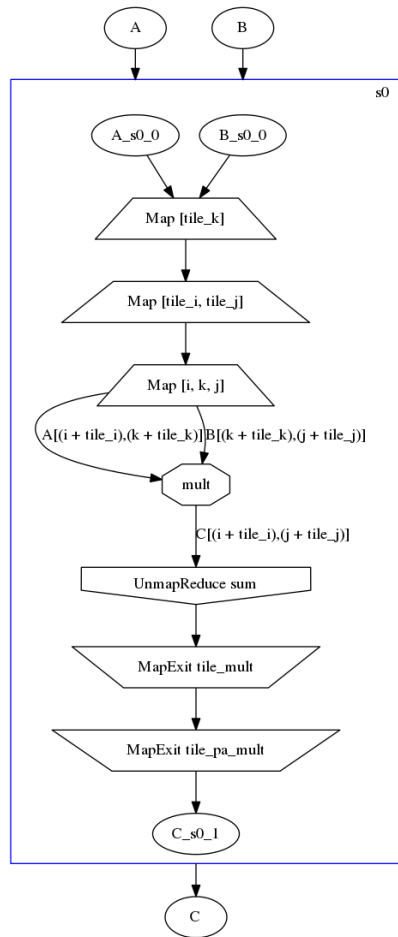


SDFG

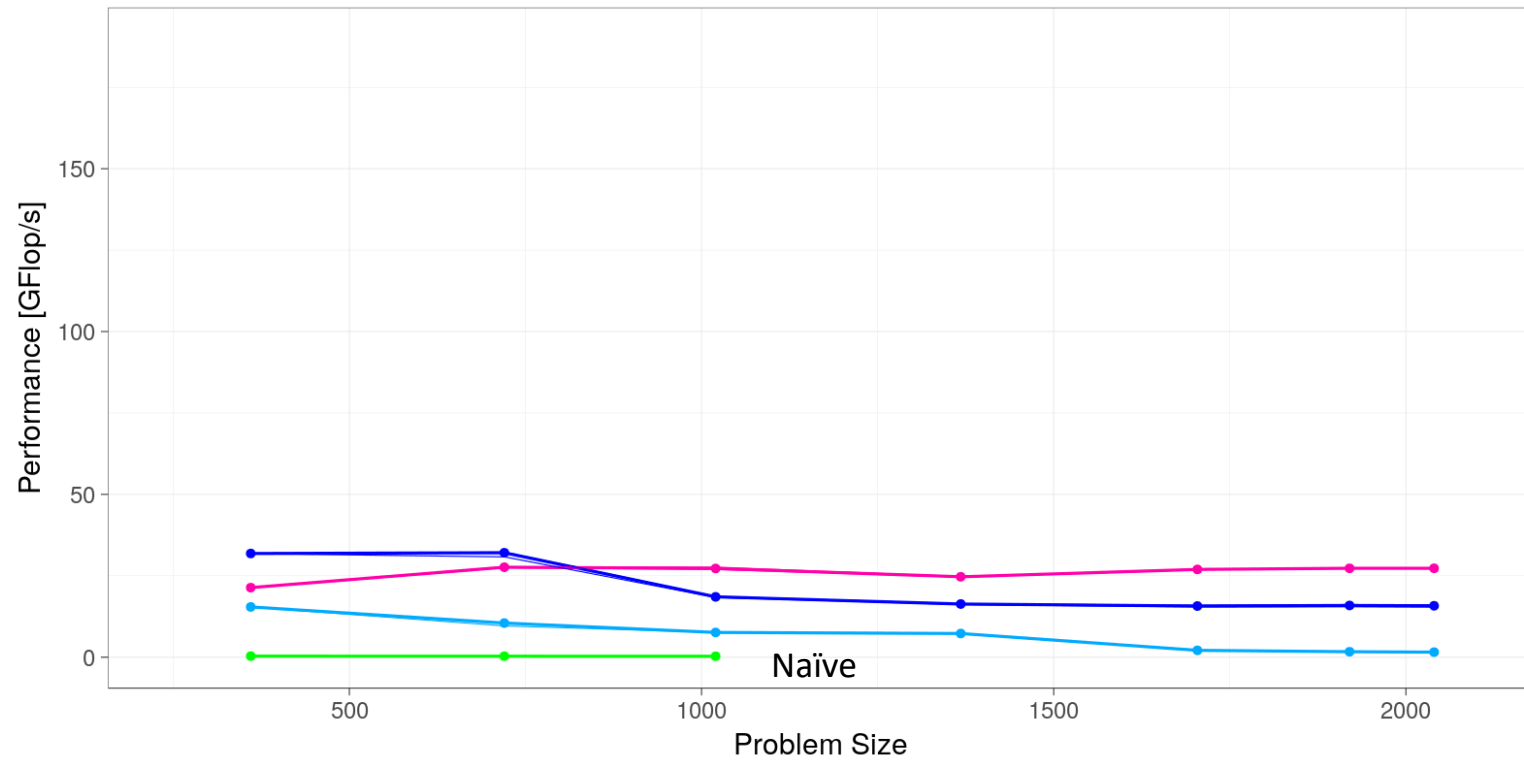


LoopReorder  
MapReduceFusion

# Performance

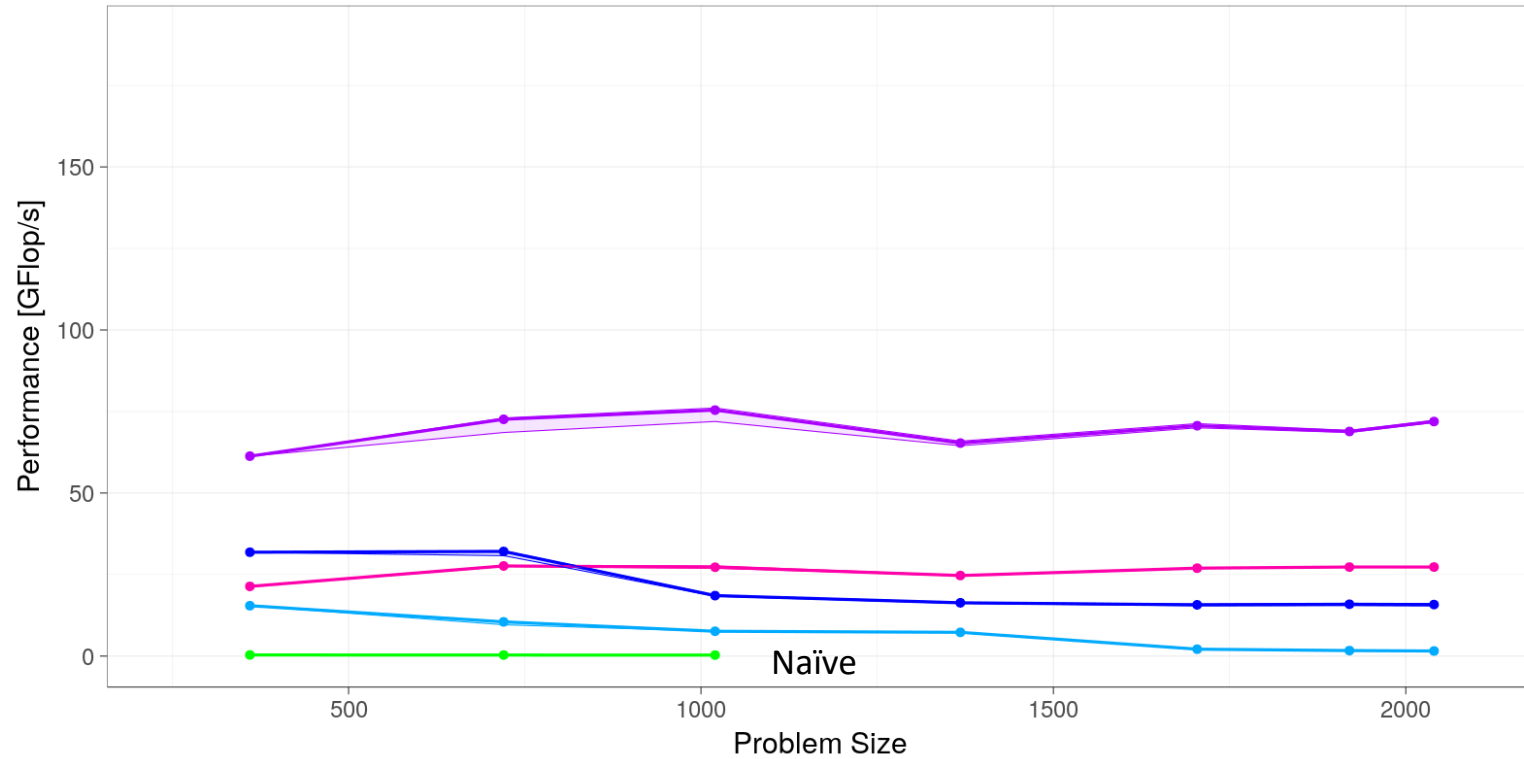
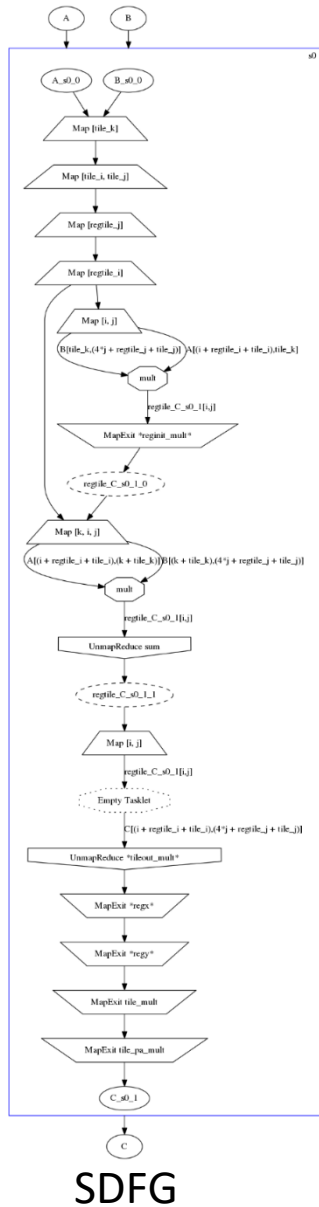


SDFG



BlockTiling  
LoopReorder  
MapReduceFusion

# Performance

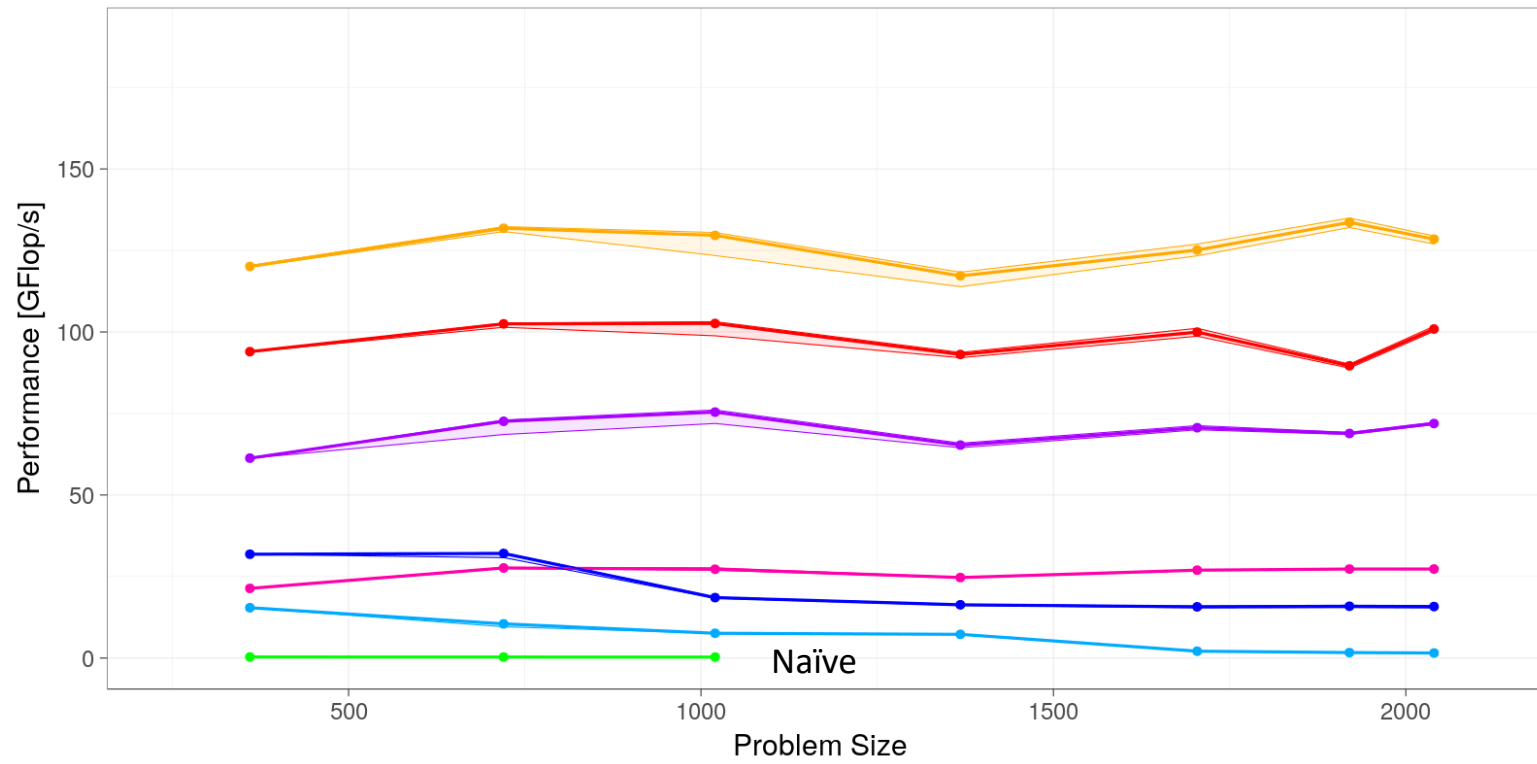
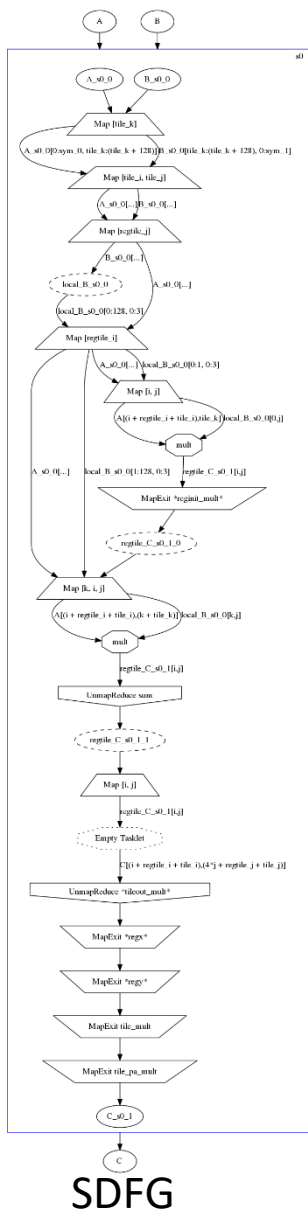


RegisterTiling

BlockTiling  
LoopReorder  
MapReduceFusion

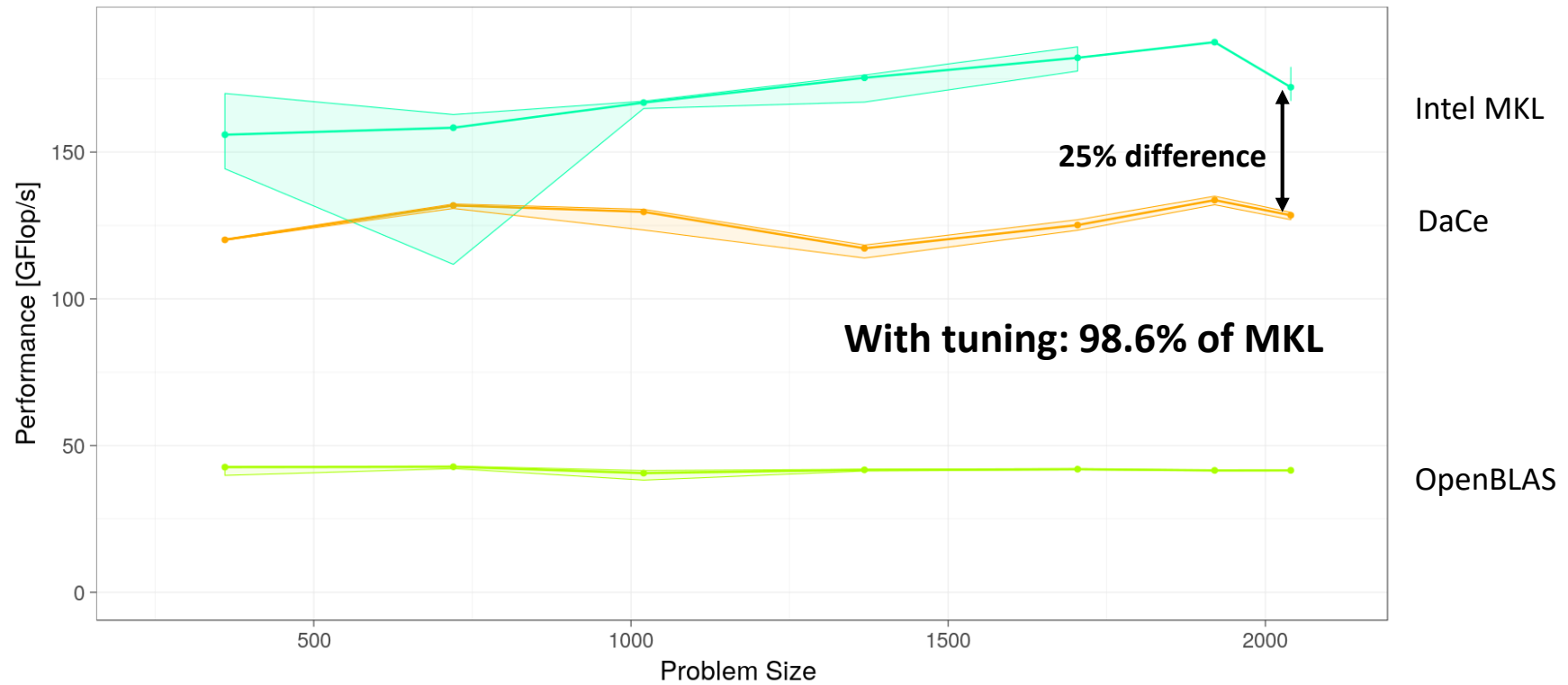


# Performance



- PromoteTransient
- LocalStorage
- RegisterTiling
- BlockTiling
- LoopReorder
- MapReduceFusion

# Performance







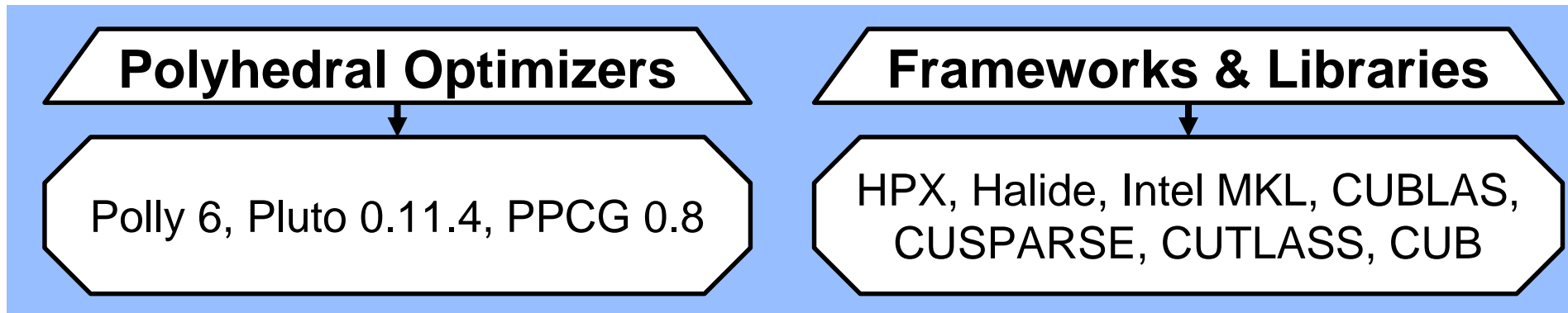
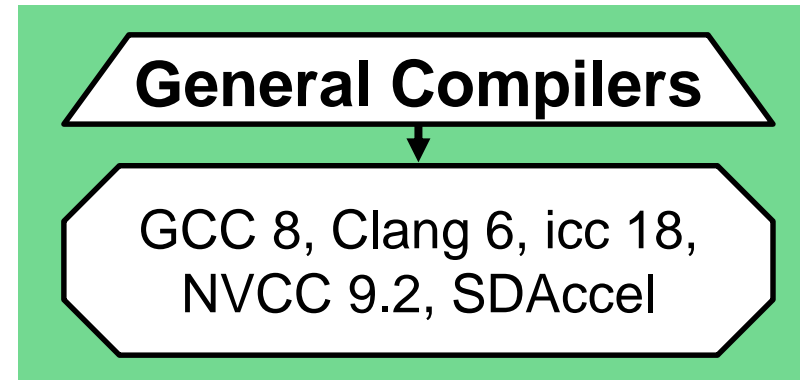
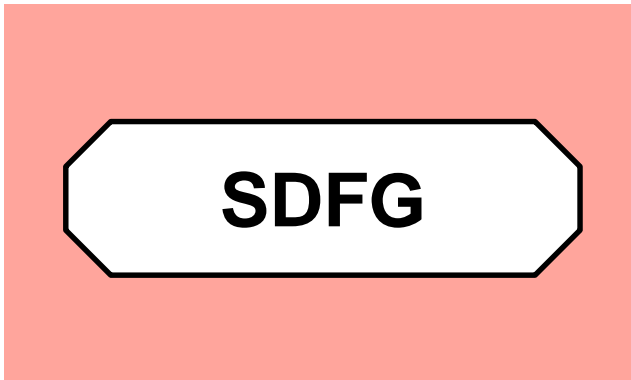
Intel Xeon E5-2650 v4



NVIDIA Tesla P100



Xilinx VU9P



# Performance Evaluation: Fundamental Kernels (CPU)

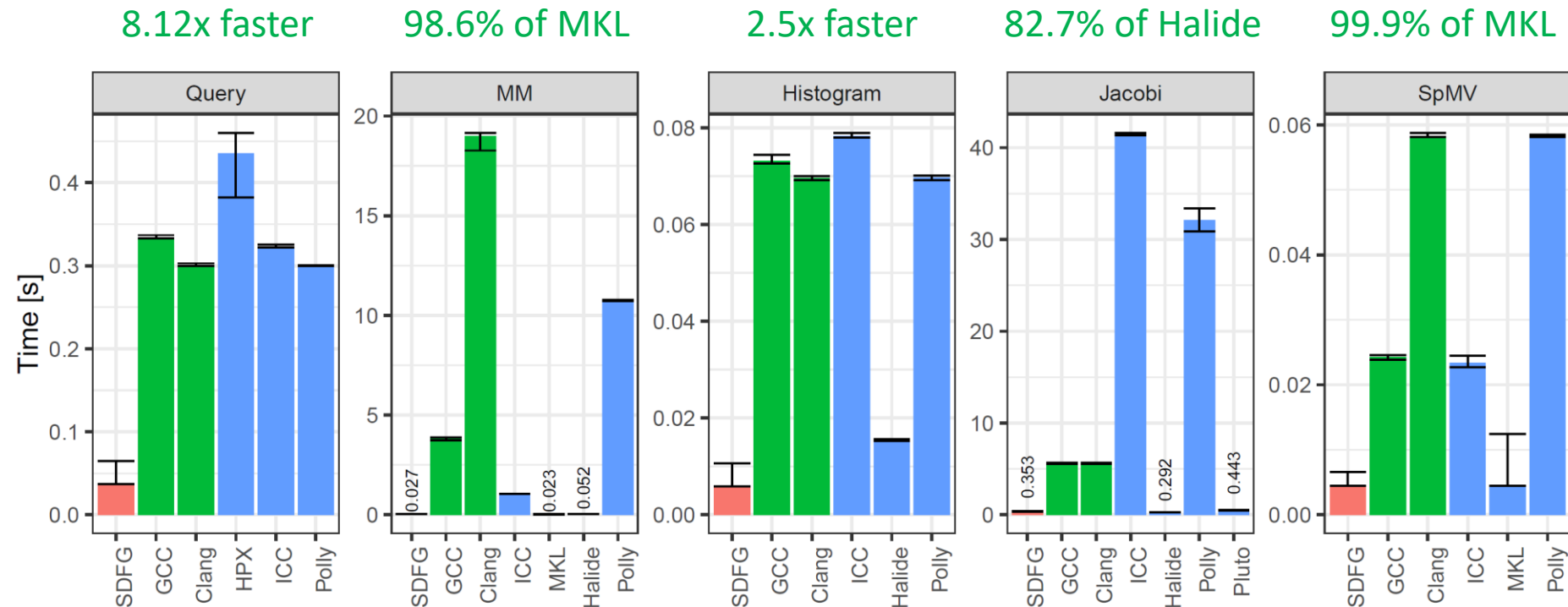
**Database Query:** roughly 50% of a 67,108,864 column

**Matrix Multiplication (MM):** 2048x2048x2048

**Histogram:** 8192x8192

**Jacobi stencil:** 2048x2048 for T=1024

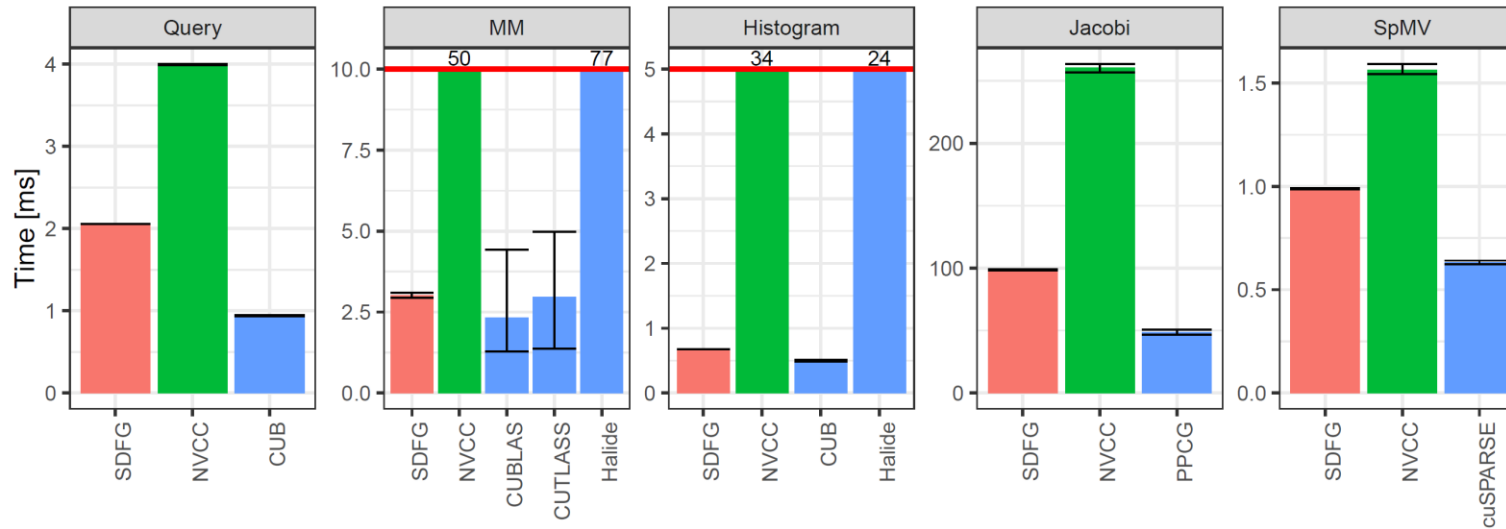
**Sparse Matrix-Vector Multiplication (SpMV):** 8192x8192 CSR matrix (nnz=33,554,432)



# Performance Evaluation: Fundamental Kernels (GPU, FPGA)

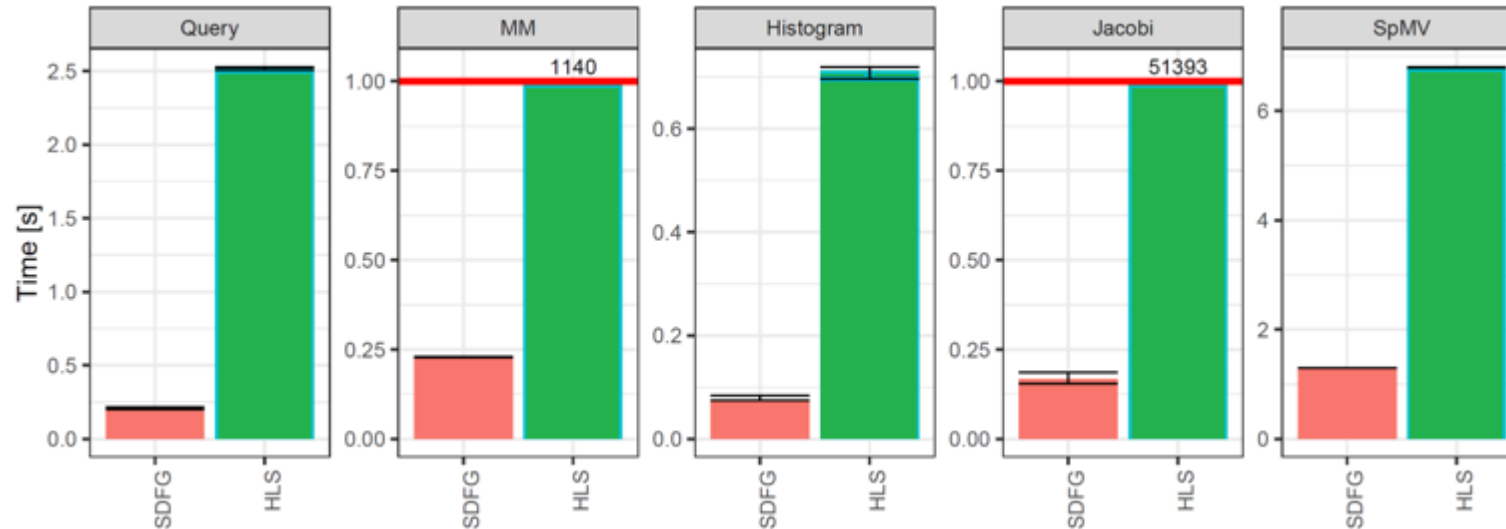
90% of CUTLASS

GPU



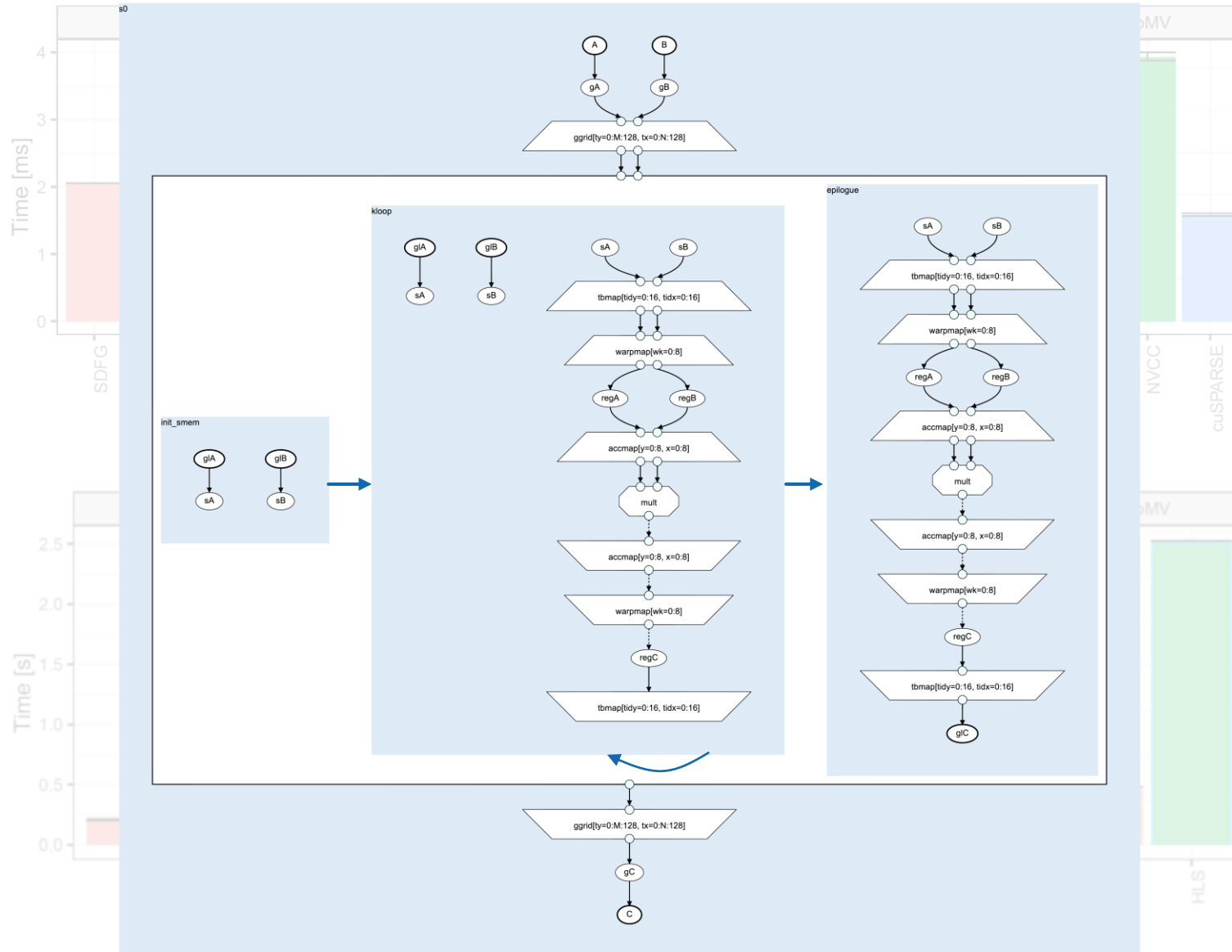
19.5x of Spatial

FPGA



# Performance Evaluation: Fundamental Kernels (GPU, FPGA)

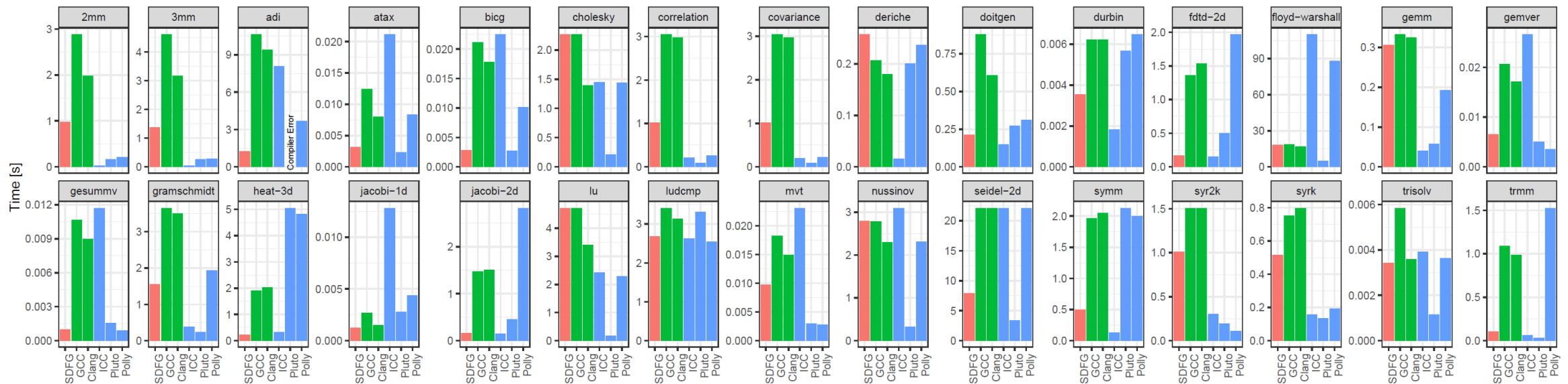
GPU



FPGA

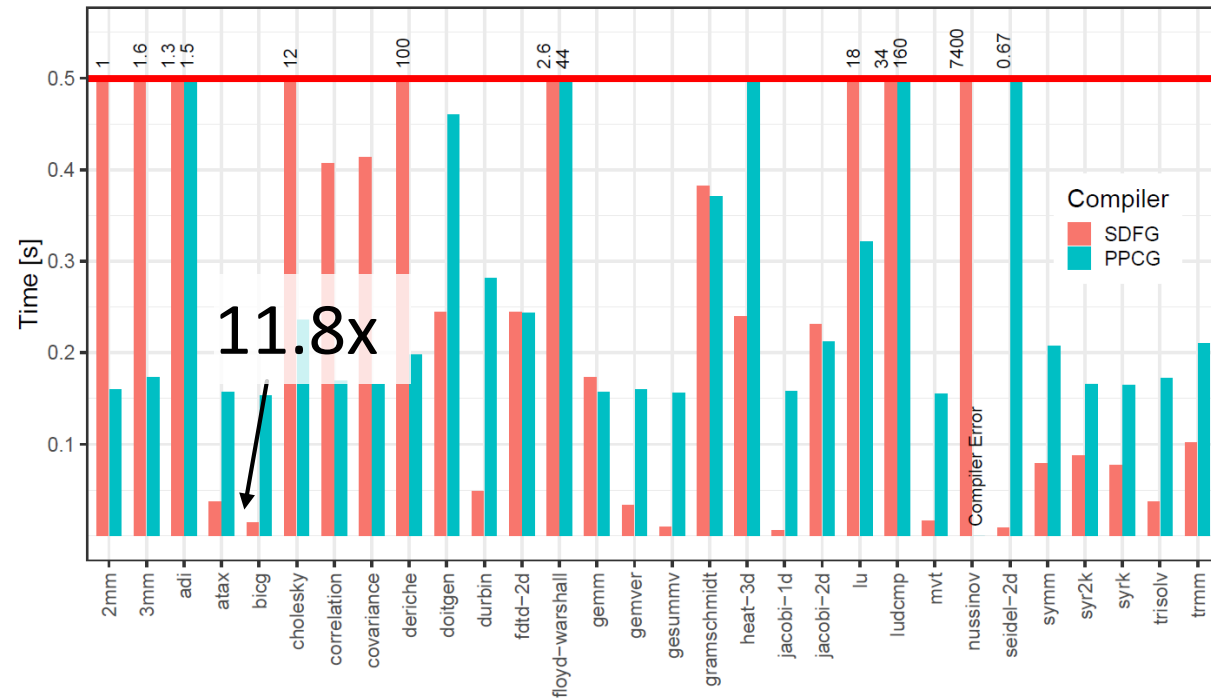
# Performance Evaluation: Polybench (CPU)

- Polyhedral benchmark with 30 applications
- Without any transformations, achieves 1.43x (geometric mean) over general-purpose compilers



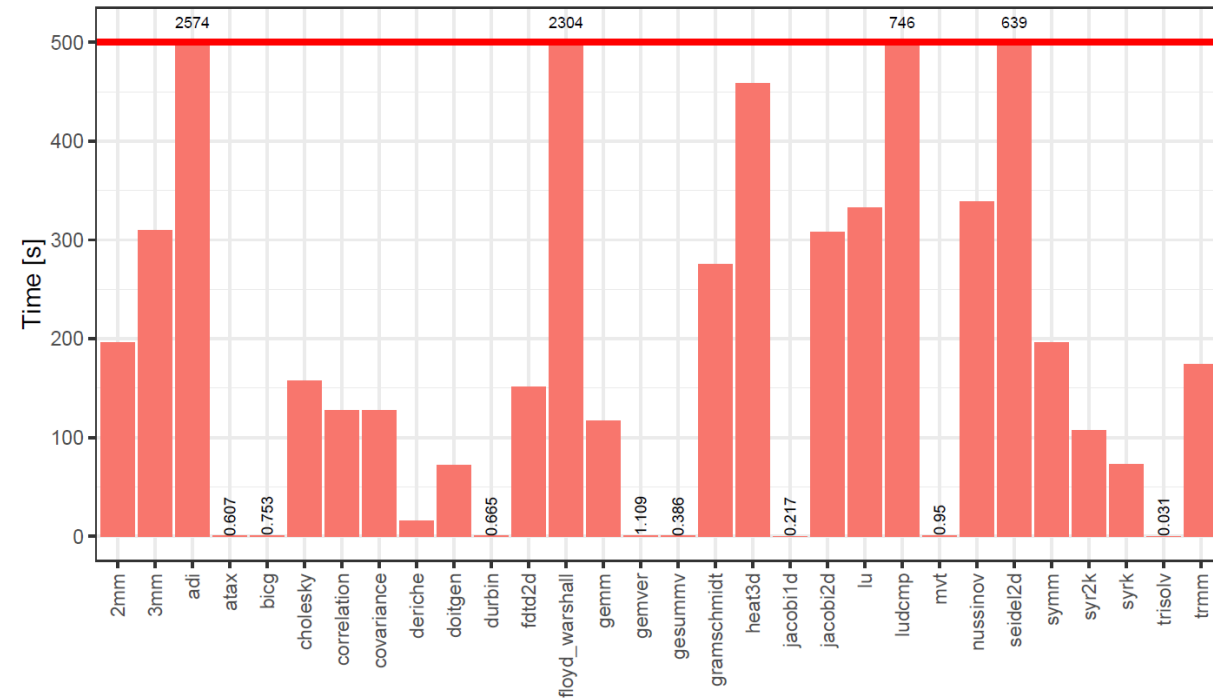
# Performance Evaluation: Polybench (GPU, FPGA)

- Automatically transformed from CPU code



**GPU**

(1.12x geomean speedup)



**FPGA**

The **first** full set of placed-and-routed Polybench

# Case Study: Parallel Breadth-First Search

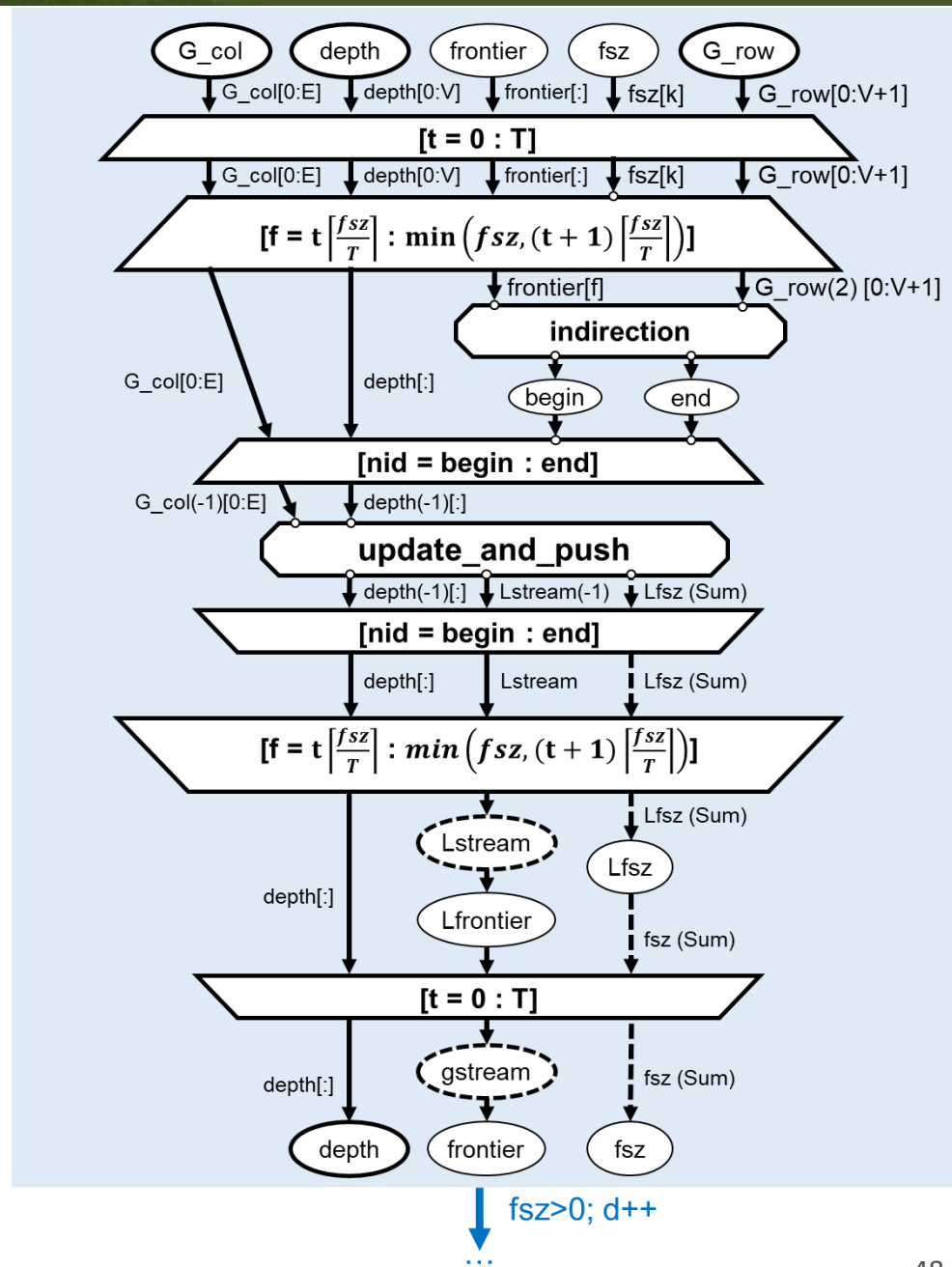
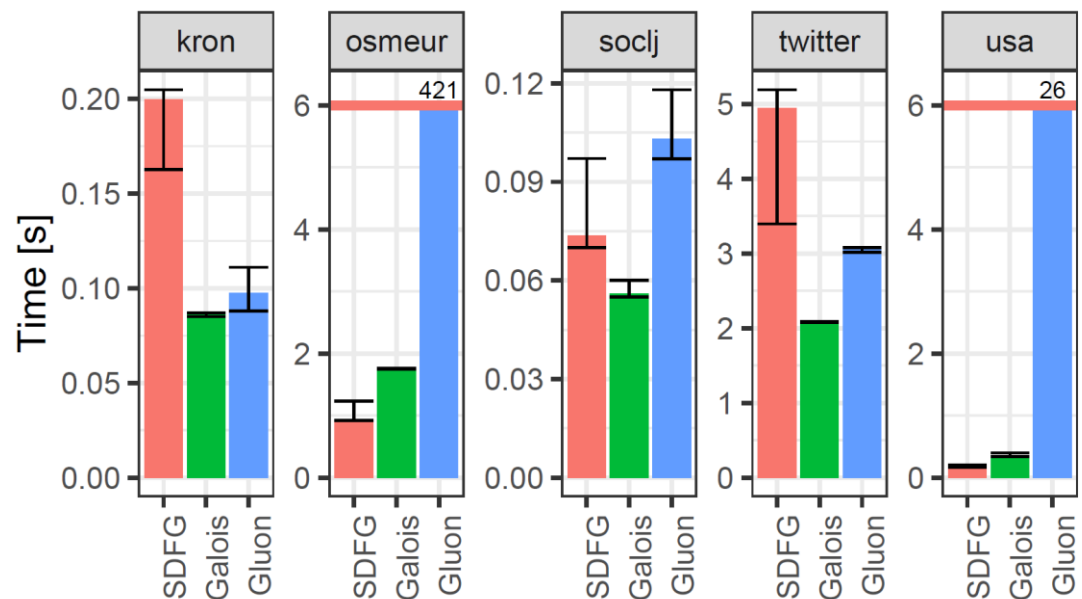
- Compared with Galois and Gluon

- Graphs:

Road maps: USA, OSM-Europe

Social networks: Twitter, LiveJournal

Synthetic: Kronecker Graphs

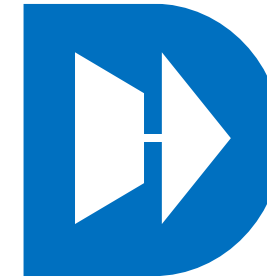




# Conclusions

<https://www.github.com/spcl/dace>

pip install dace



```
@dapp.program
def program(A, B):
    @dapp.map(_[0:N,0:M])
    def transpose(i, j):
        a << A[i,j]
        b >> B[j,i]
    ...
```

