**ETH**_zürich_

**D** INFK

**Alexandru Calotoiu**

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

Jan Kleine, Philipp Schaad and others at SPCL

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

# Motivation

Cray-1 Vector Processor

CPU

Multi-Core

GPU Computing Heterogeneous Systems

NVIDIA CUDA

OpenCL

OpenACC
More Science. Less Programming

AMD RADEON GRAPHICS

NVIDIA GEFORCE GTX

Specialization

QUALCOMM ZEROTH NPU

SAMSUNG Exynos 9

FPGAs and beyond

intel Stratix 10

KINTEX UltraSCALE

BlueField
Mellanox

intel inside

1976          1992      2002        2007           2012          2016    2017    Today
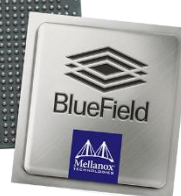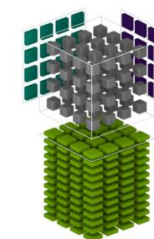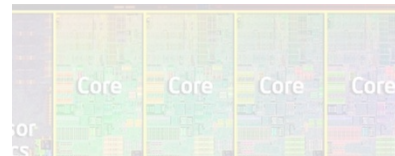
# Motivation

Cray-1 Vector
Processor

CPU   Multi-Core

GPU Computing
Heterogeneous Systems

Specialization

FPGAs and
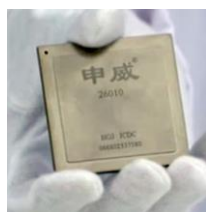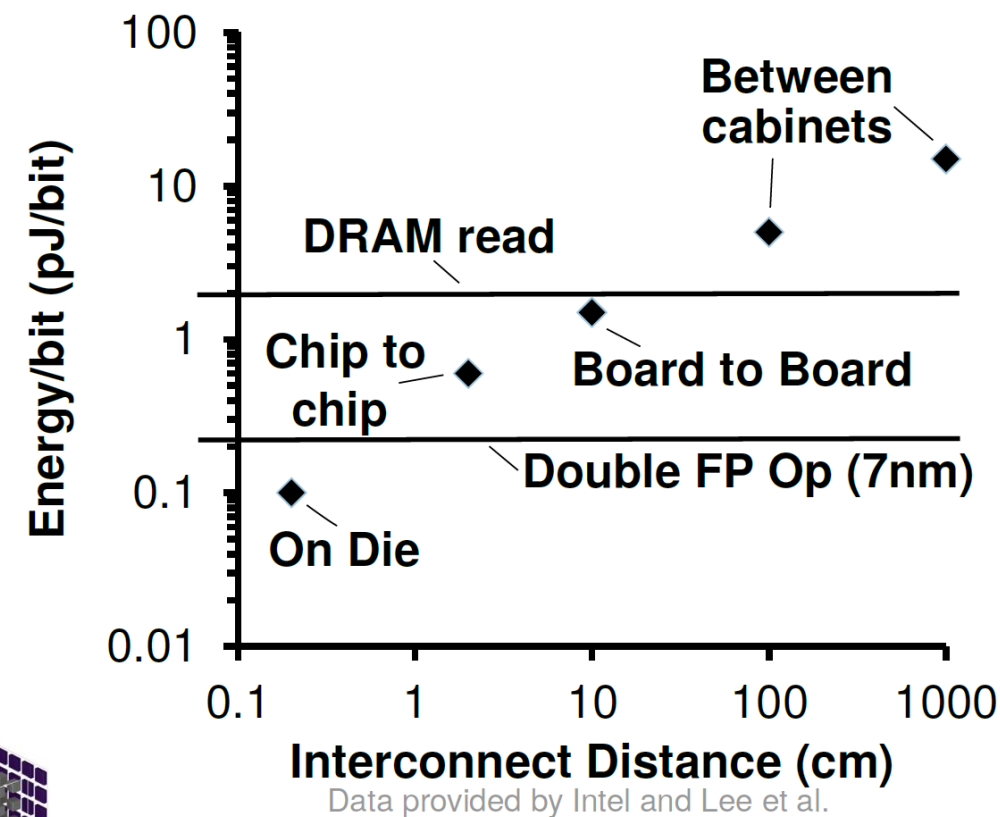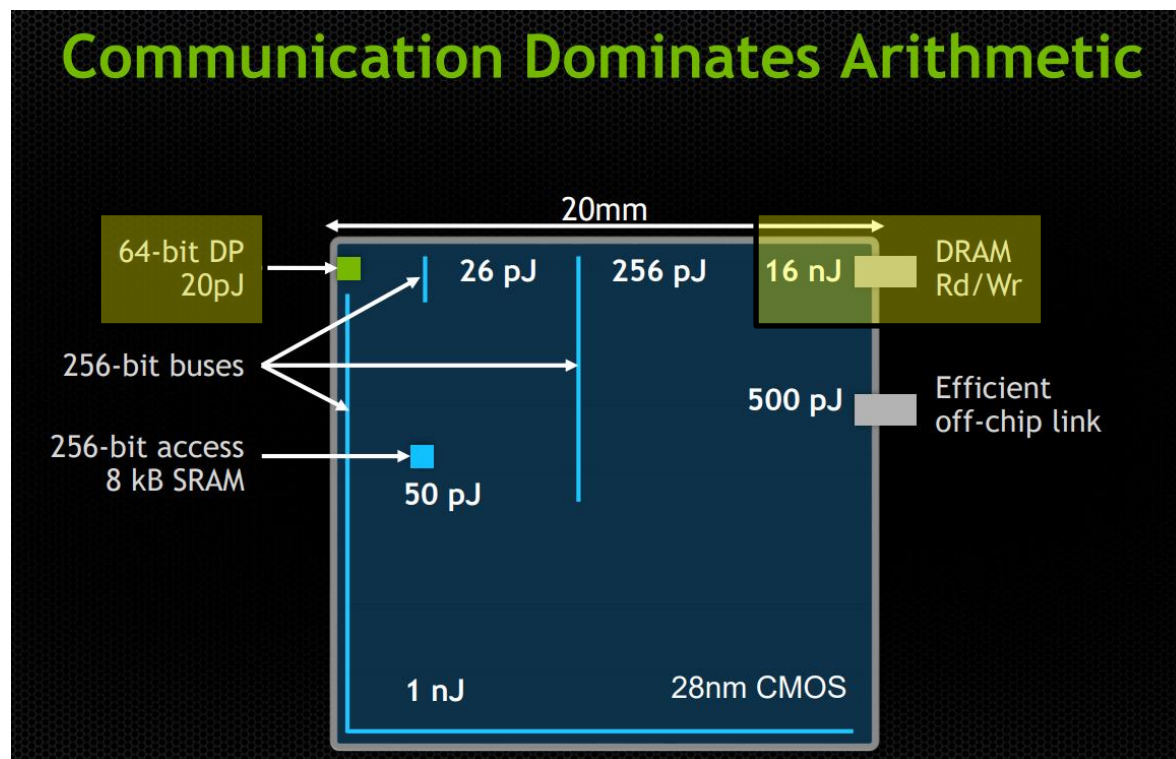beyond

# Newer computer ≠ faster application

| 1976 | 1992 | 2002 | 2007 | 2012 | 2016 | 2017 | Today |

# Motivation

## Communication Dominates Arithmetic

20mm

64-bit DP 20pJ

26 pJ     256 pJ     16 nJ     DRAM Rd/Wr

256-bit buses

256-bit access 8 kB SRAM

50 pJ

500 pJ     Efficient off-chip link

1 nJ     28nm CMOS

Slide courtesy of NVIDIA



**Energy/bit (pJ/bit)** vs **Interconnect Distance (cm)**

Between cabinets

DRAM read

Chip to chip     Board to Board

Double FP Op (7nm)

On Die

Data provided by Intel and Lee et al.

4

# Motivation



Communication Dominates Arithmetic

20mm

64-bit DP
20pJ          26 pJ    256 pJ    16 nJ     DRAM
                                           Rd/Wr

256-bit access
8 kB SRAM
                50 pJ              500 pJ    Efficient
                                            off-chip link

1 nJ                     28nm CMOS

Slide courtesy of NVIDIA

Computation costs decrease, data movement remains expensive!

Data provided by Intel and Lee et al.
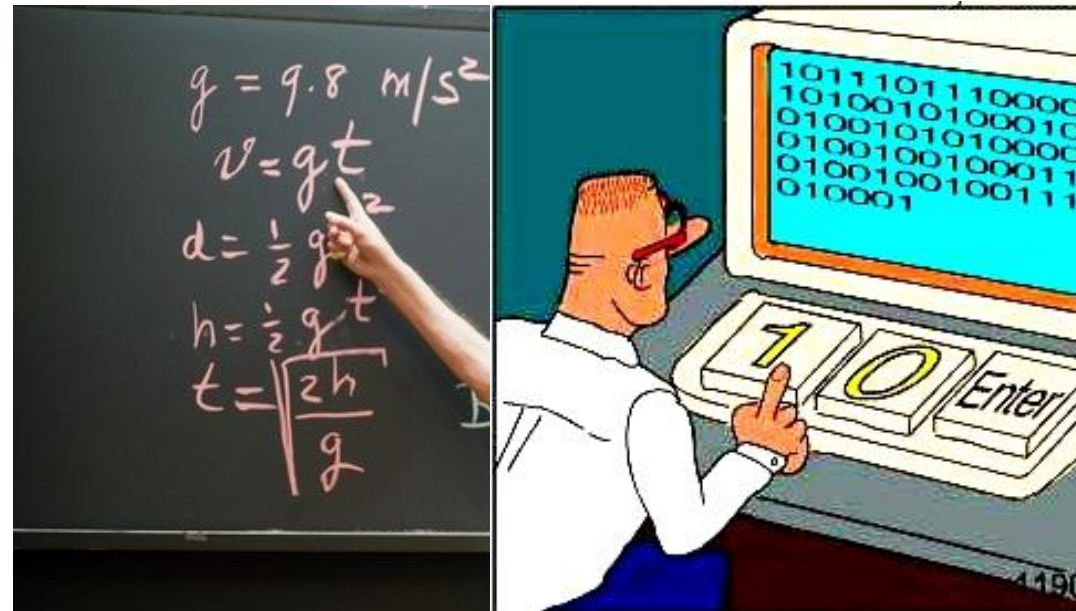
Source: US DoE

# Computational Scientist

**Domain Scientist** **Performance Engineer**

# Optimization Techniques

- **Multi-core CPU**
  - Tiling for complex cache hierarchies
  - Register optimizations
  - Vectorization

- Many-core GPU

## High-performance optimization = data movement reduction

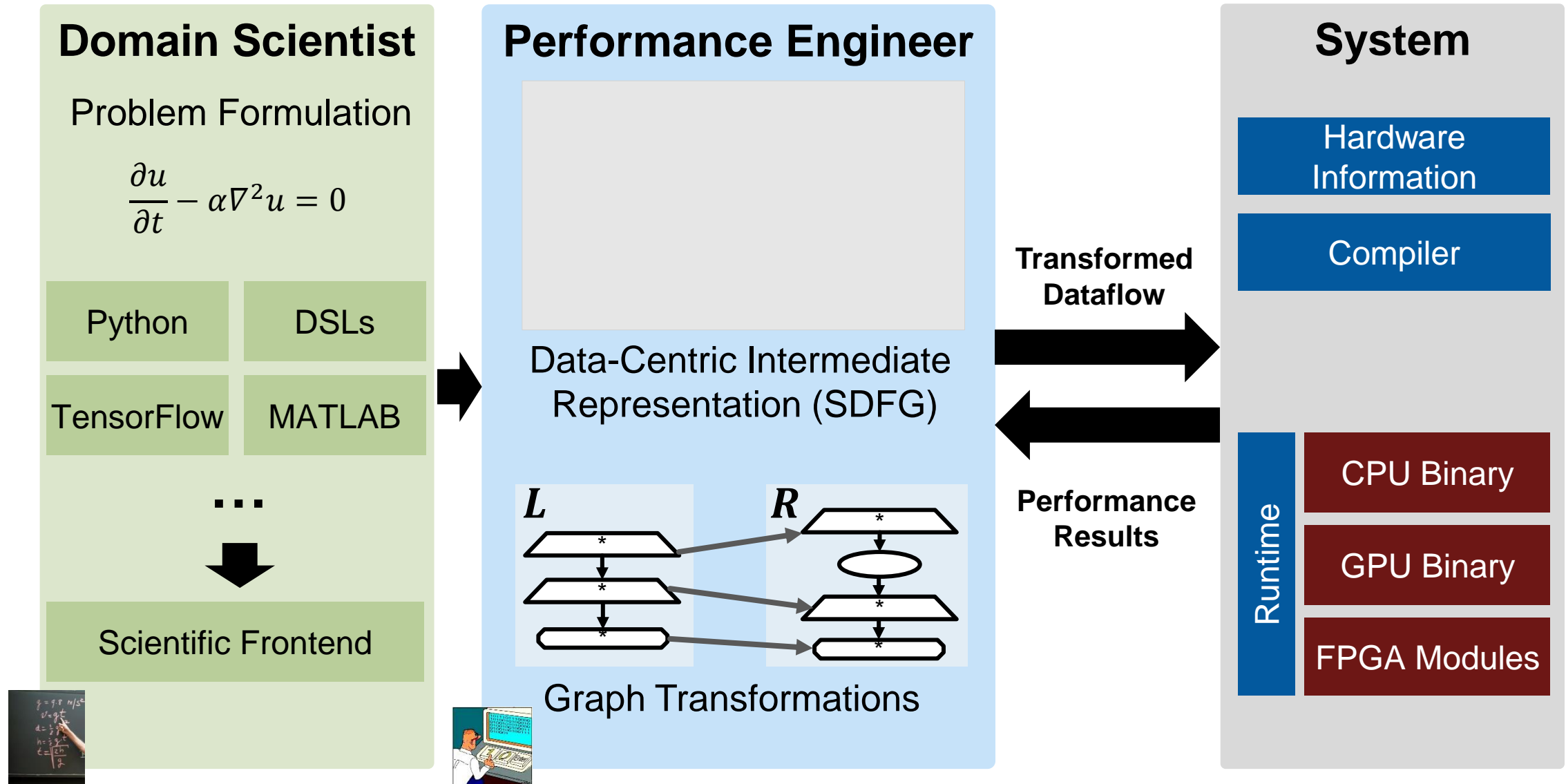  - Warp divergence minimization, register tiling
  - Task fusion

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

# ▶aCe Overview

## Domain Scientist

Problem Formulation

$$\frac{\partial u}{\partial t} - \alpha \nabla^2 u = 0$$

| Python | DSLs |
| --- | --- |
| TensorFlow | MATLAB |

· · ·

Scientific Frontend

## Performance Engineer

Data-Centric Intermediate Representation (SDFG)



*L*      *R*

Graph Transformations

**Transformed Dataflow**

**Performance Results**

## System

Hardware Information

Compiler

Runtime

CPU Binary

GPU Binary

FPGA Modules

# 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

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

- **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

```python
@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
```
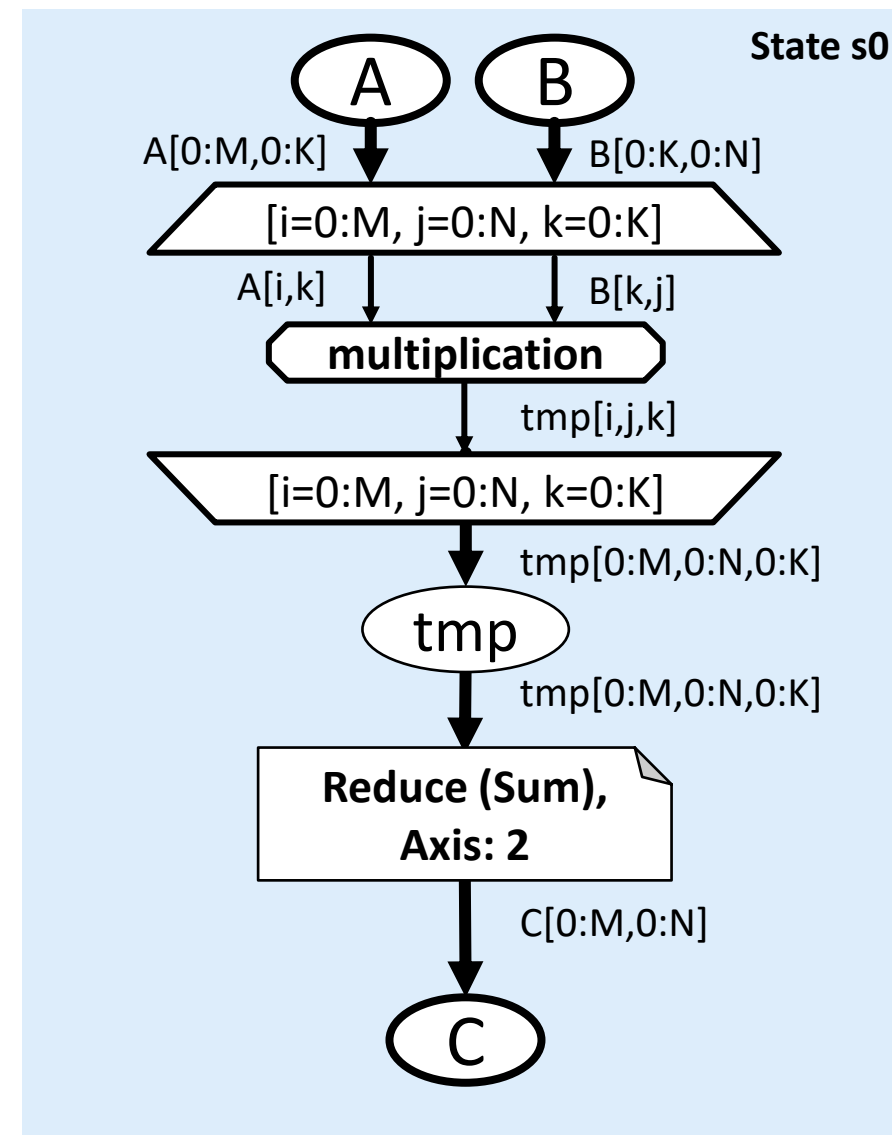
# Matrix Multiplication SDFG

```python
@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)
```



State s0

A  B

A[0:M,0:K]          B[0:K,0:N]

[i=0:M, j=0:N, k=0:K]

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

multiplication

tmp[i,j,k]

[i=0:M, j=0:N, k=0:K]

tmp[0:M,0:N,0:K]

tmp

tmp[0:M,0:N,0:K]

Reduce (Sum),
Axis: 2
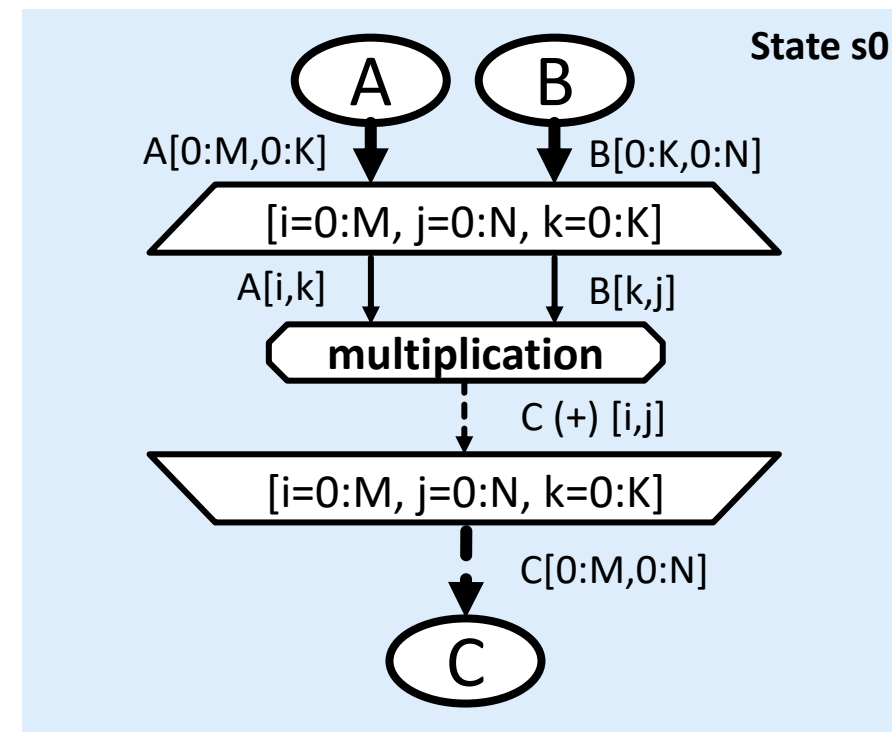
C[0:M,0:N]

C

# 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)
```
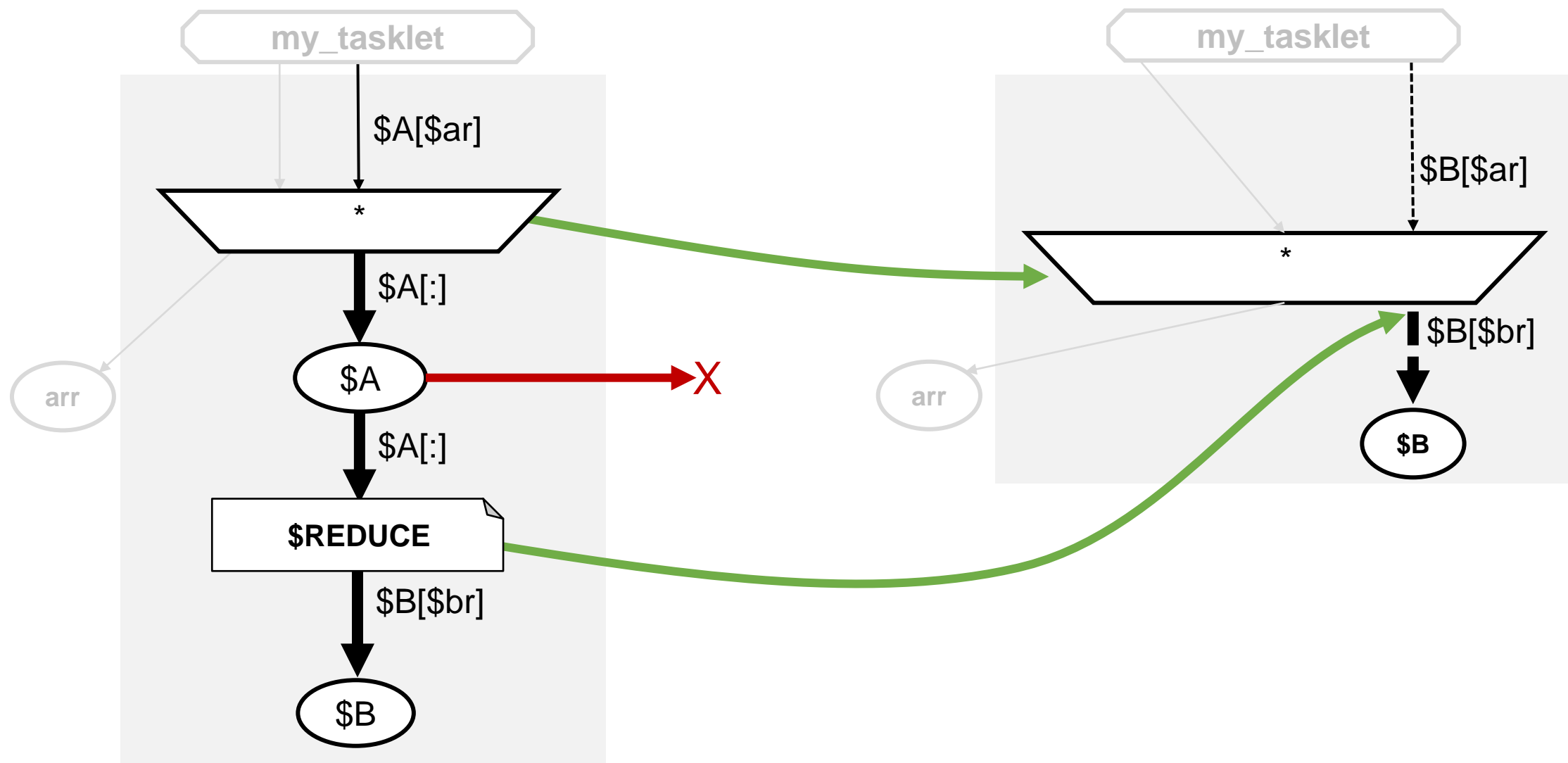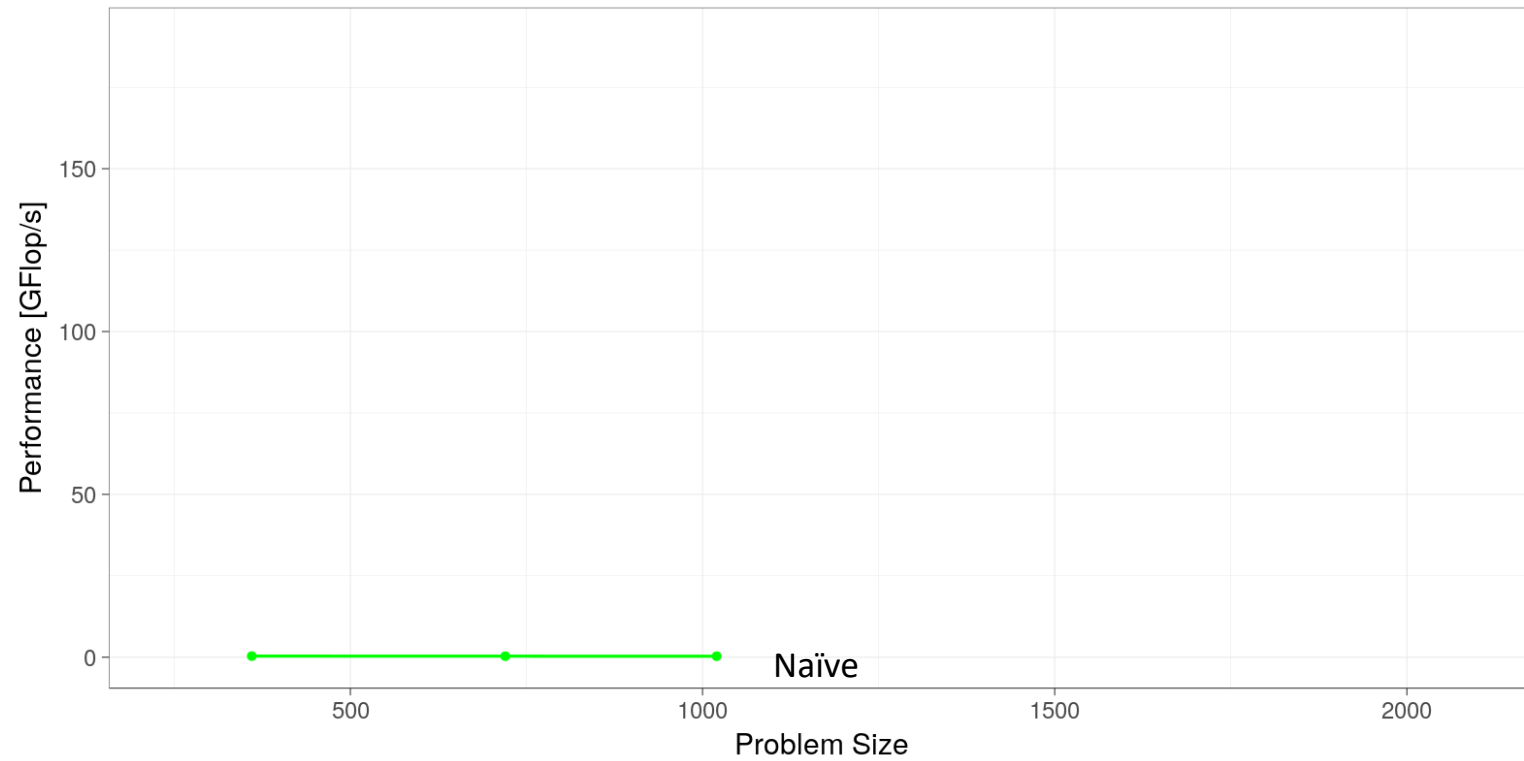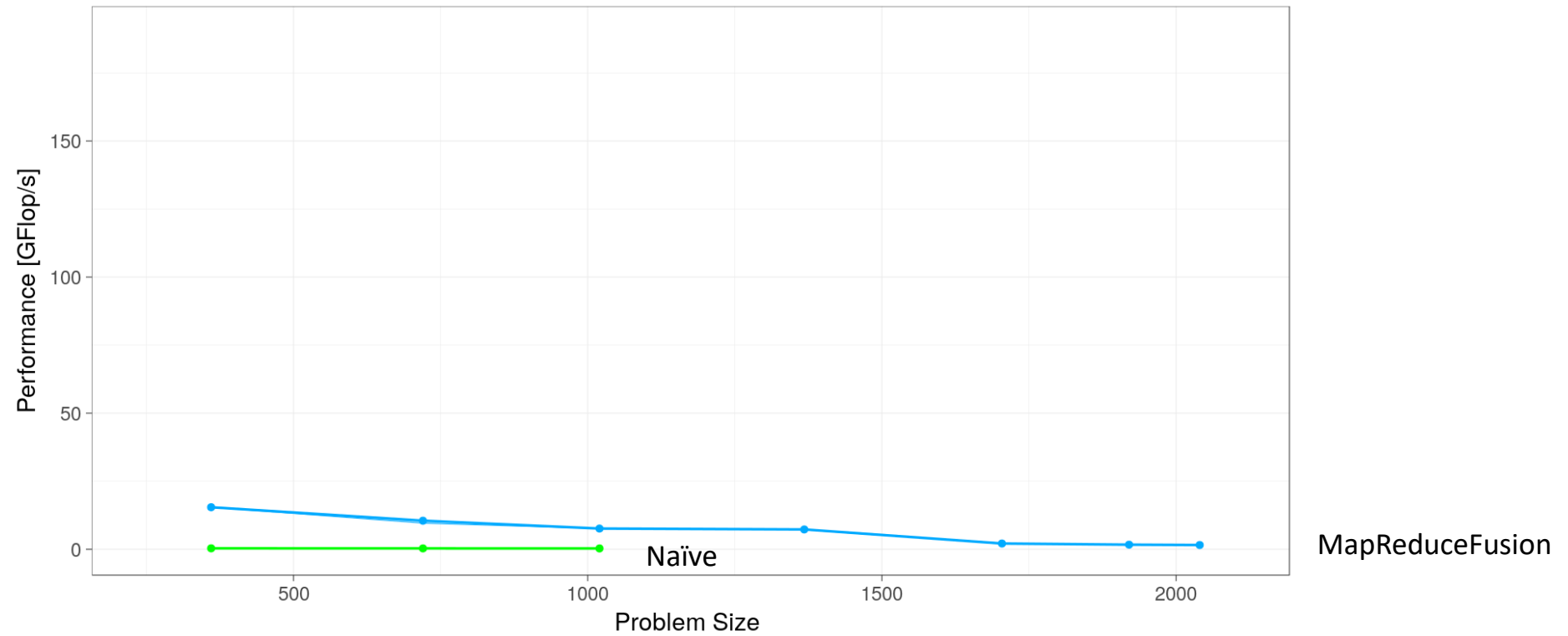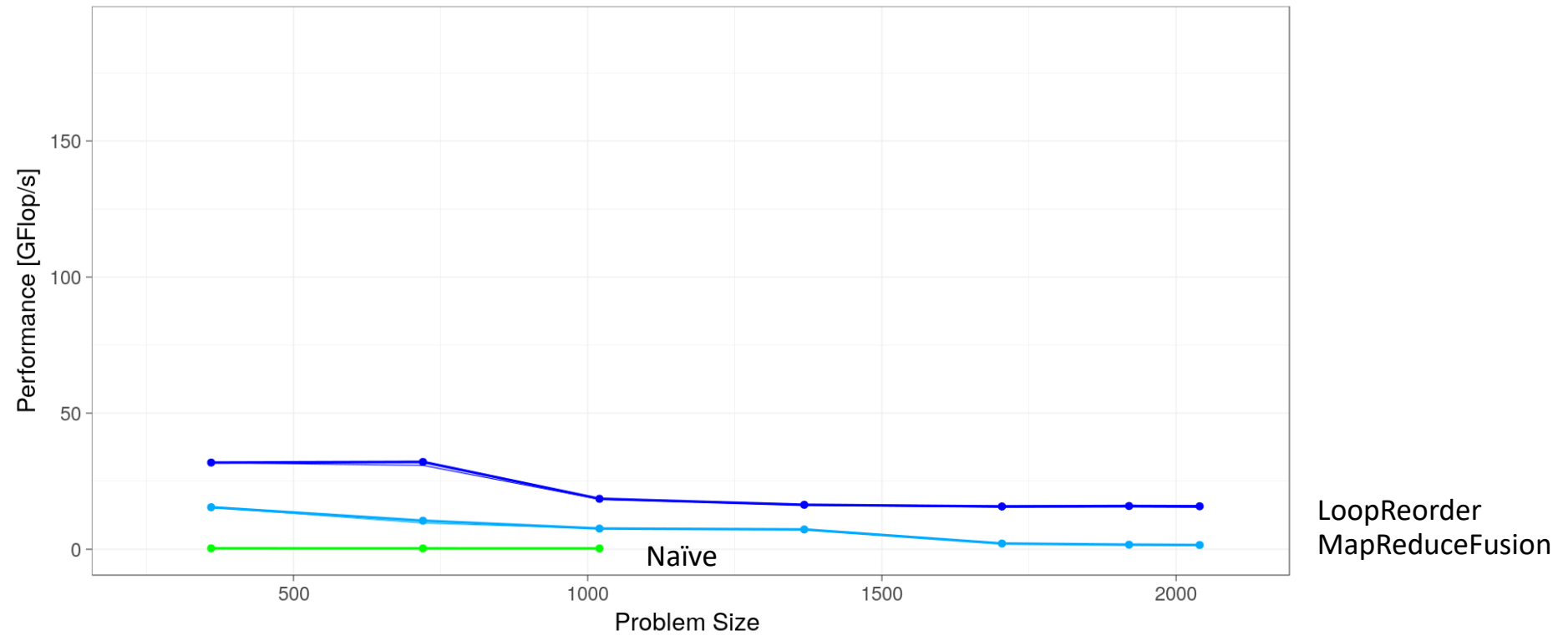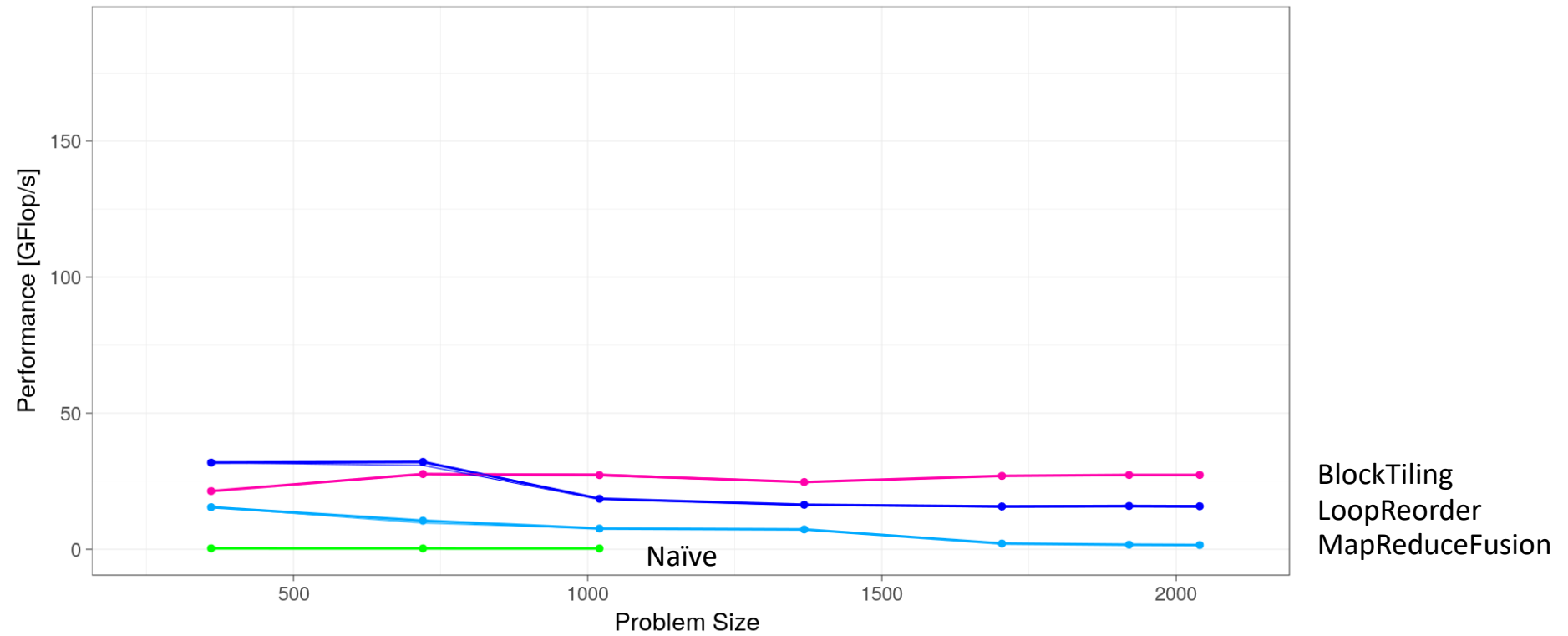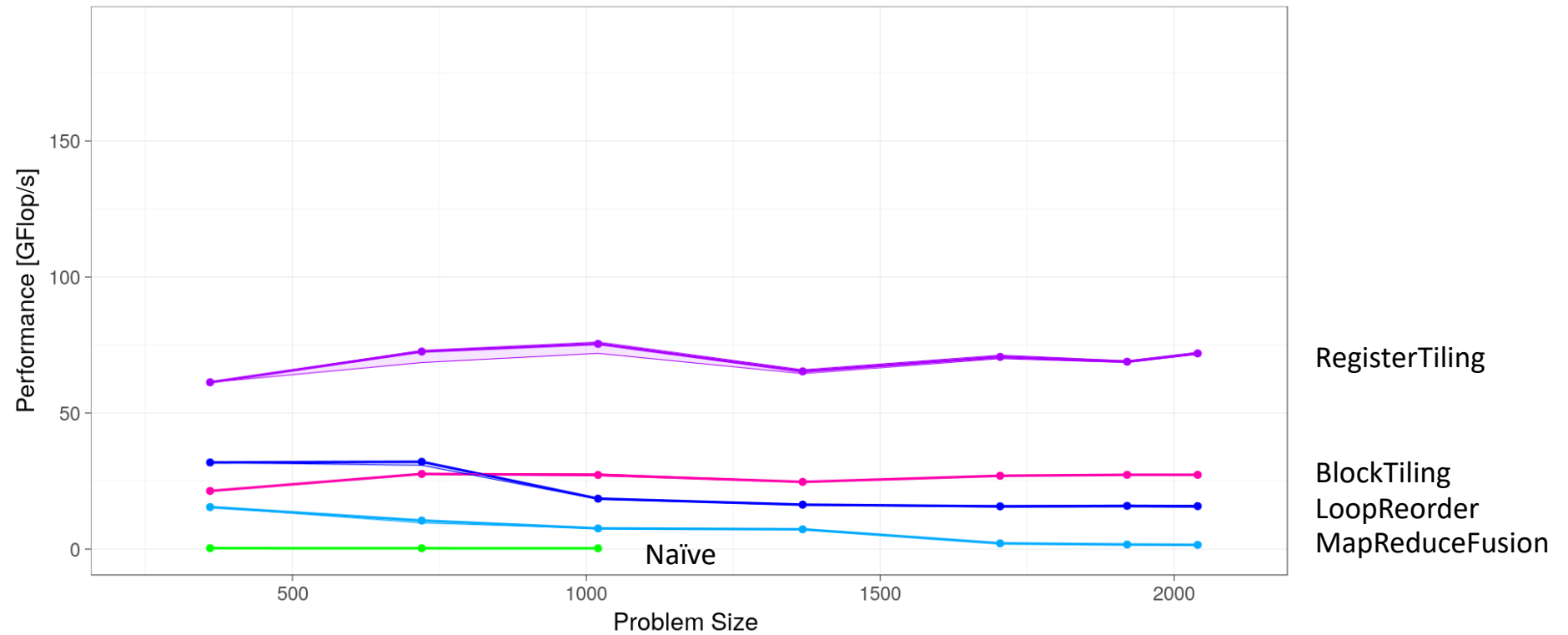


State s0

# MapReduceFusion Transformation

# Performance

# Performance

# Performance

# Performance

# Performance

# Performance

# Performance

# Performance



Intel MKL

25% difference

DaCe

With tuning: 98.6% of MKL
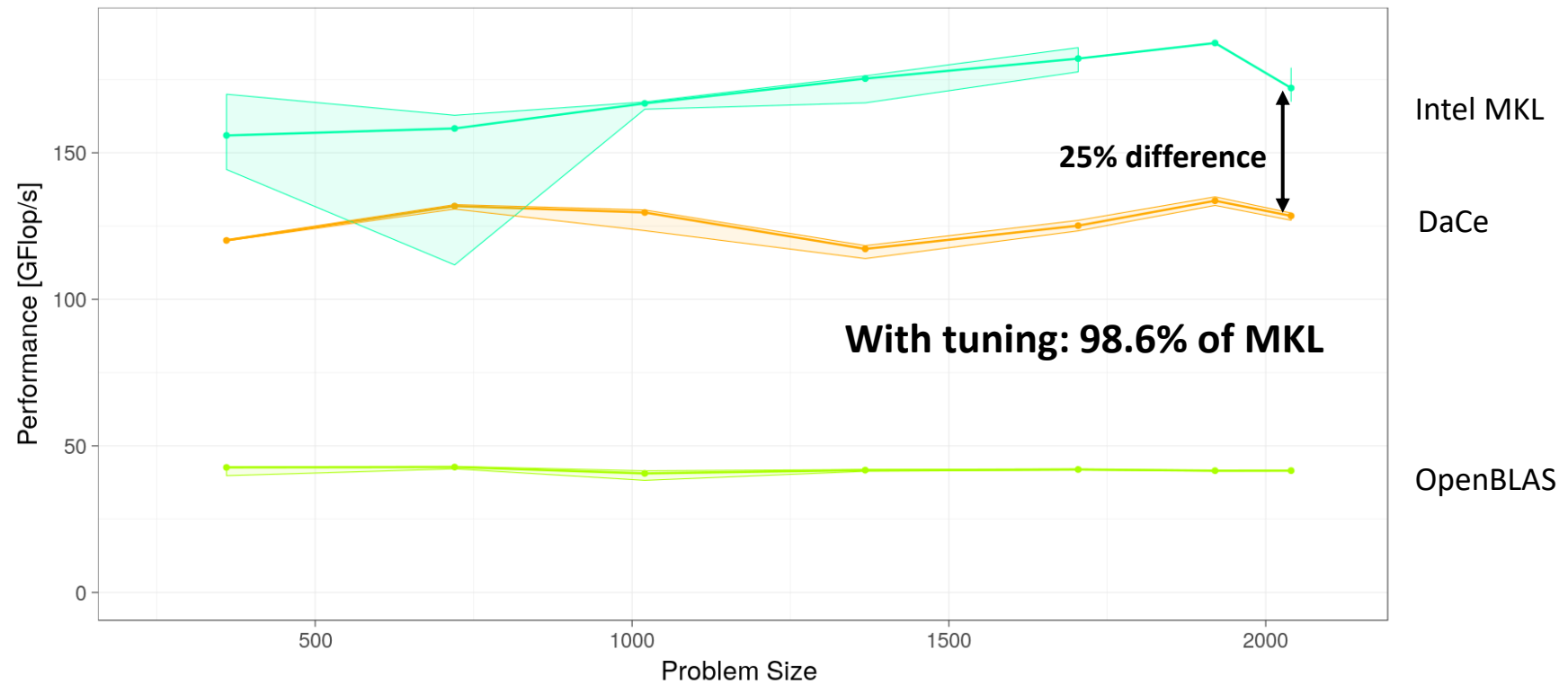
OpenBLAS

**Intel Xeon E5-2650 v4**    **NVIDIA Tesla P100**    **Xilinx VU9P**

**SDFG**

**General Compilers**

GCC 8, Clang 6, icc 18, NVCC 9.2, SDAccel

**Polyhedral Optimizers**

Polly 6, Pluto 0.11.4, PPCG 0.8

**Frameworks & Libraries**

HPX, Halide, Intel MKL, CUBLAS, CUSPARSE, CUTLASS, CUB

# Related work

**Separation of Concerns**

CHiLL [16], Halide [45], Tiramisu [9], SPIRAL [28]

**Transformable Representations**

Halide [45], HPVM [37], CHiLL [16], Lift [48]

**Data Access Decoupling**

**General**

Sequoia [26], Chapel [15], Halide [45], Legion [10], MAPS [12], SPIRAL [28], TensorFlow [4]

**Polyhedral**

PENCIL [8], Pluto [13], Polly [30], Tiramisu [9]

**Streaming**

OpenSPL [11], StreamIT [50]

**Graphs**

Galois [41], Gluon [21]

**Multi-Target IRs**

**Imperative Programming**

LLVM [38], OpenMP [20], OpenCL [31], Lift [48], Halide [45], SPIRAL [28], HPVM [37], Tiramisu [9]

**Spatial Programming**

FROST [47], OpenACC [39], OpenCL [19], SPIRAL [28]

23