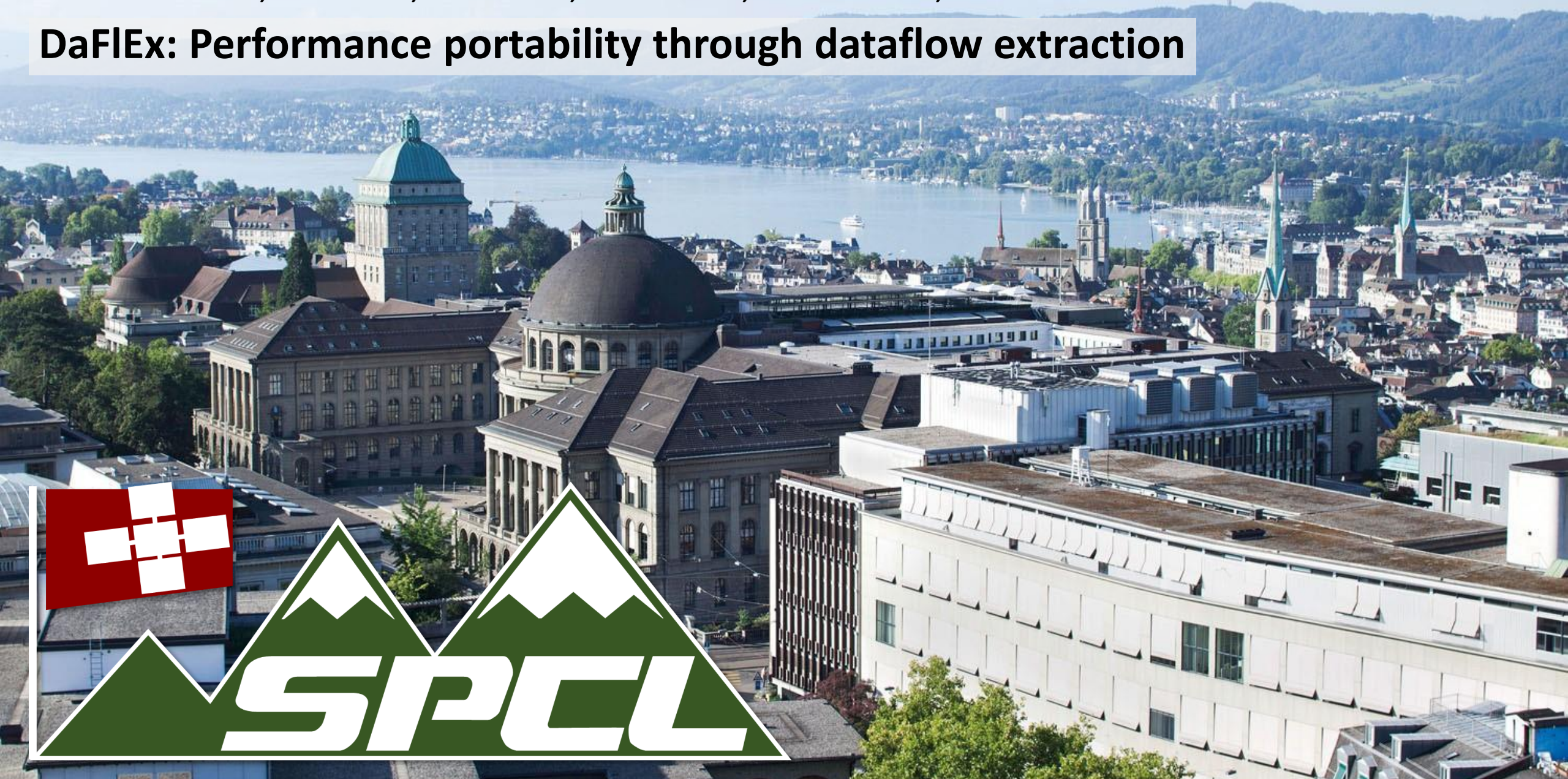


ALEXANDRU CALOTOIU, TAL BEN-NUN, NOAH HÜTTER, PAUL SCHEFFLER, TORSTEN HOEFER, LUCA BENINI

DaFLEX: Performance portability through dataflow extraction



Outline

- Introduction to DaCe
- DaFlEx overview
- DaFlEx first steps and goals

DaCe Overview

Domain Scientist

Problem Formulation

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

Python

DSLs

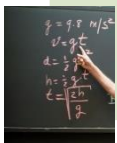
TensorFlow

MATLAB

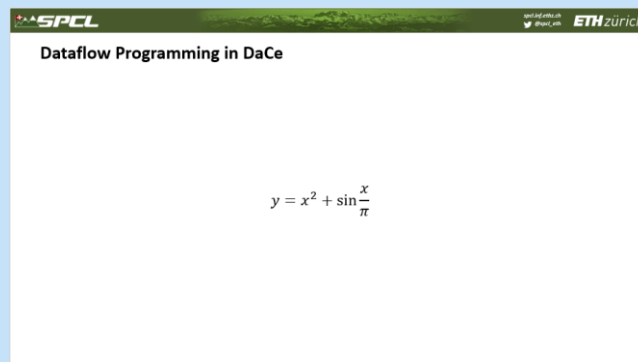
...



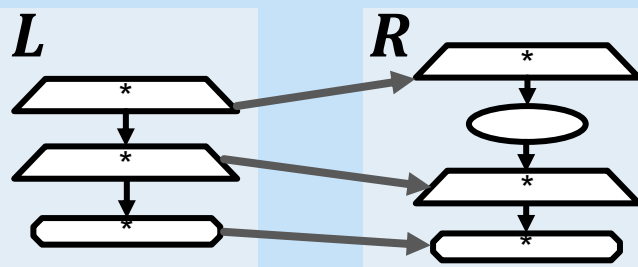
Scientific Frontend



Performance Engineer



Data-Centric Intermediate Representation (SDFG)



Graph Transformations



Transformed
Dataflow



Performance
Results



System

Hardware
Information

Compiler

Runtime

CPU Binary

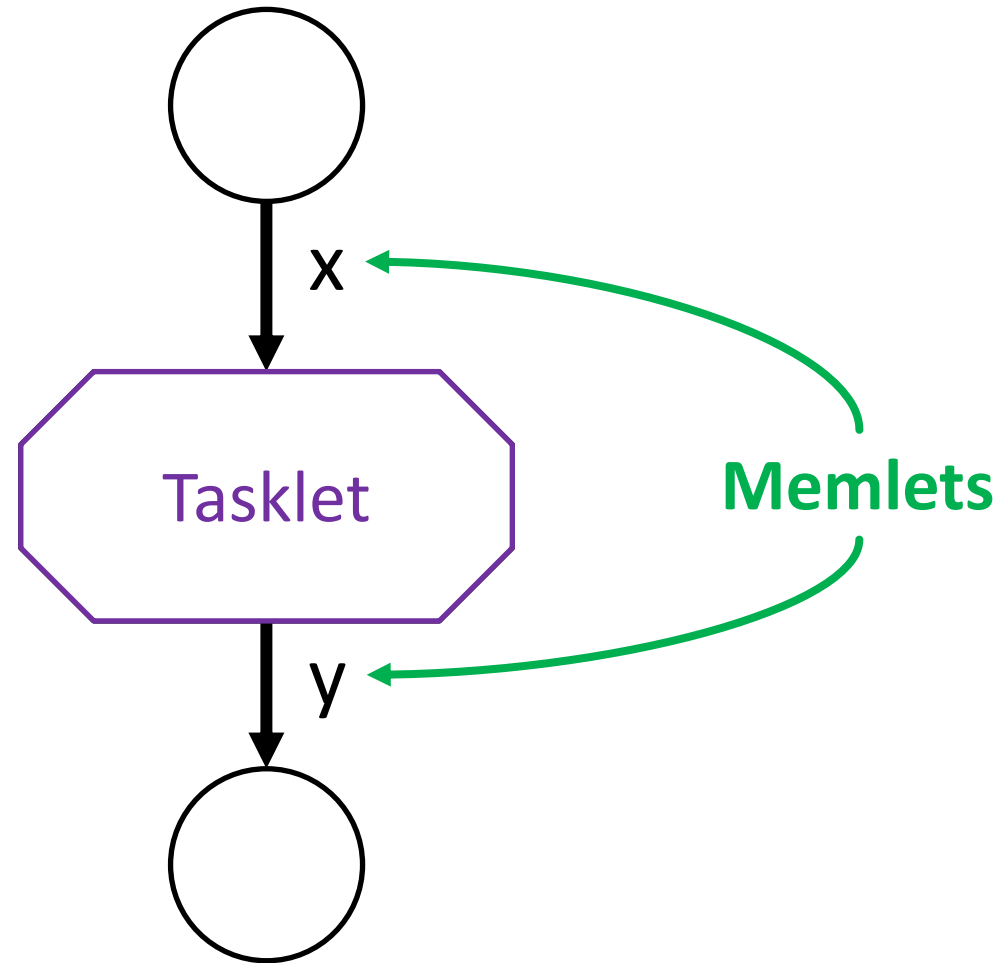
GPU Binary

FPGA Modules

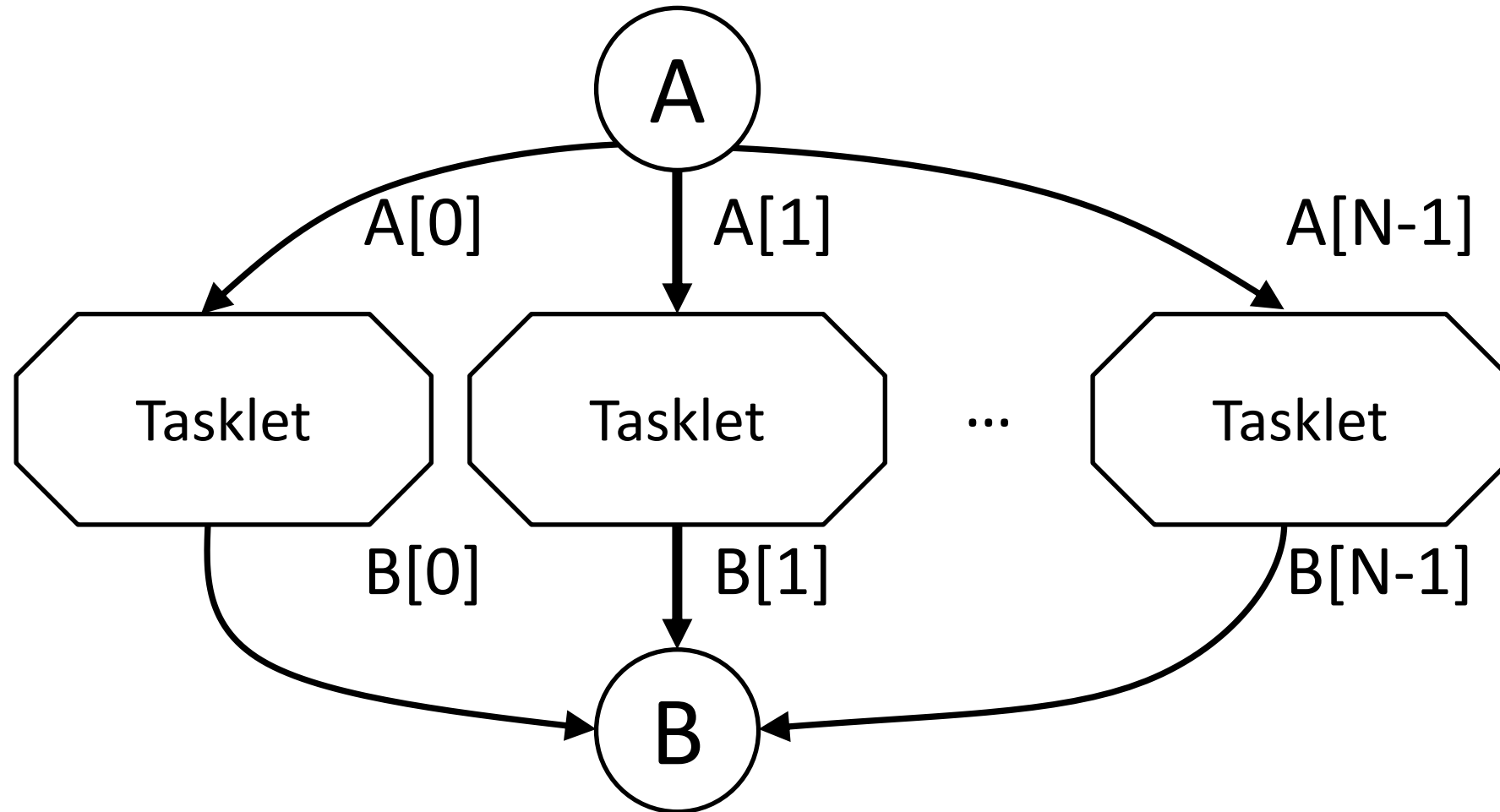
Dataflow Programming in DaCe

$$y = x^2 + \sin \frac{x}{\pi}$$

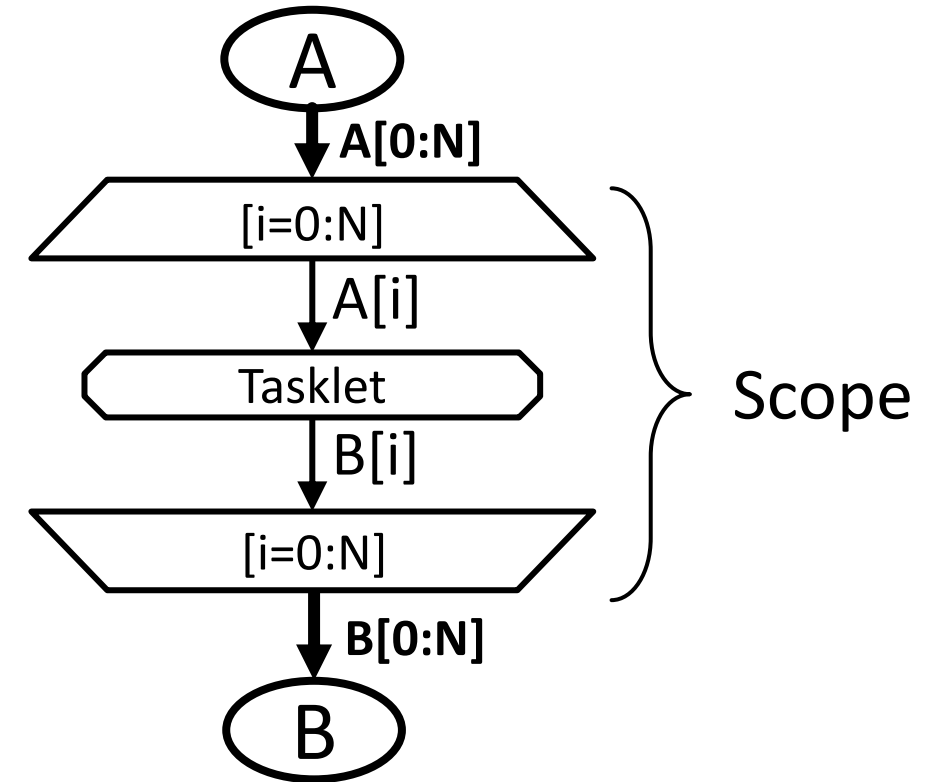
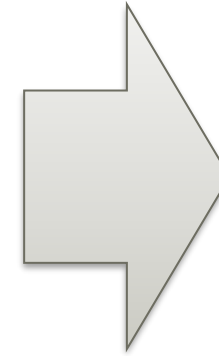
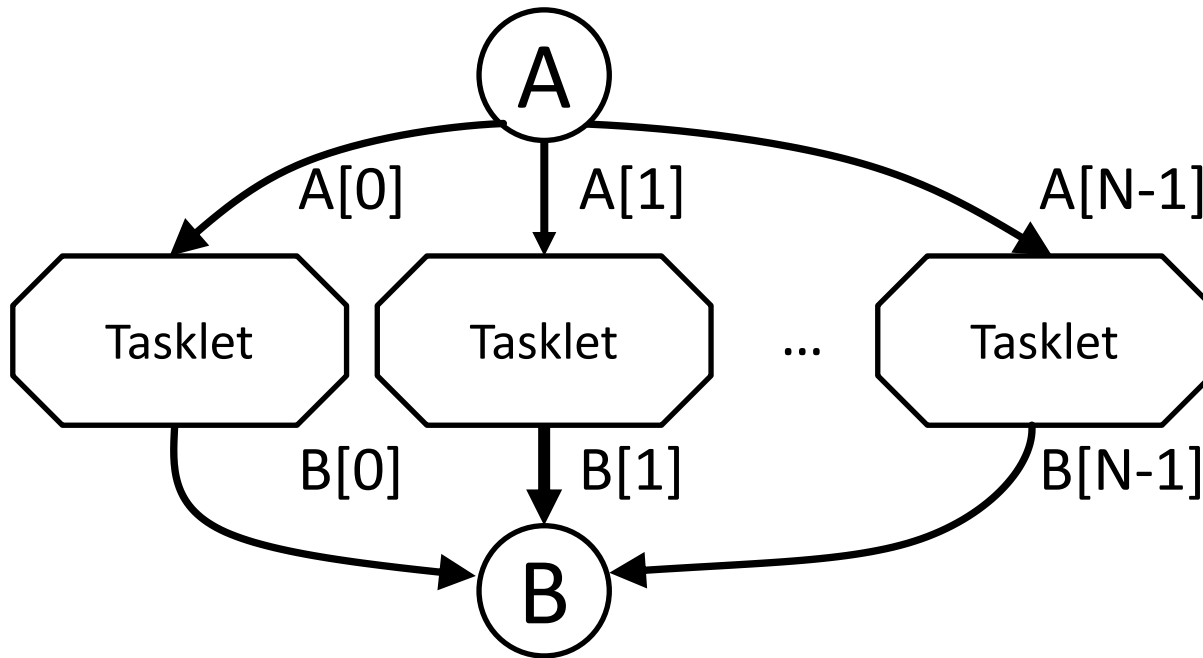
Dataflow Programming in DaCe



Parallel Dataflow Programming

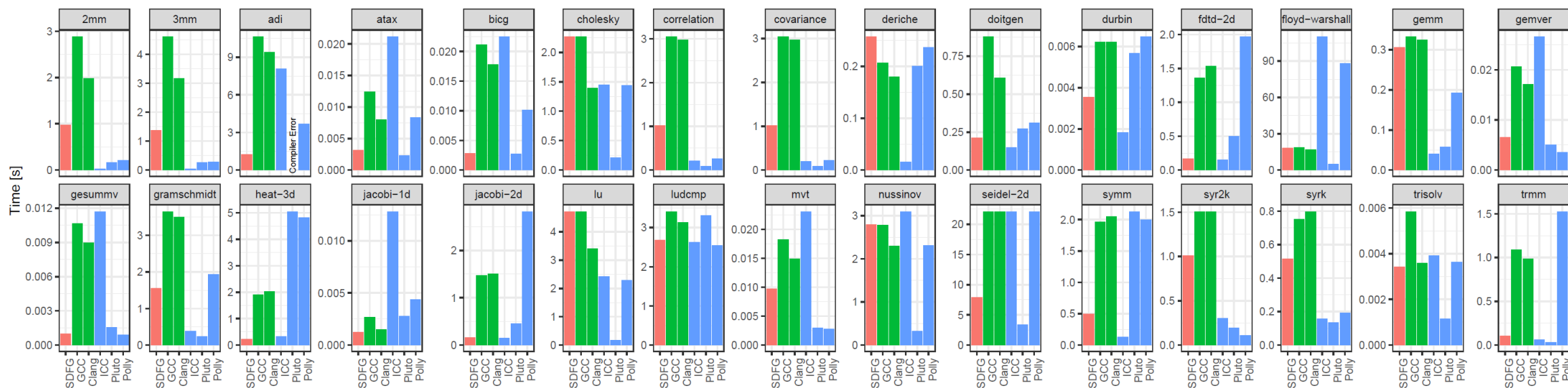


Parallel Dataflow Programming



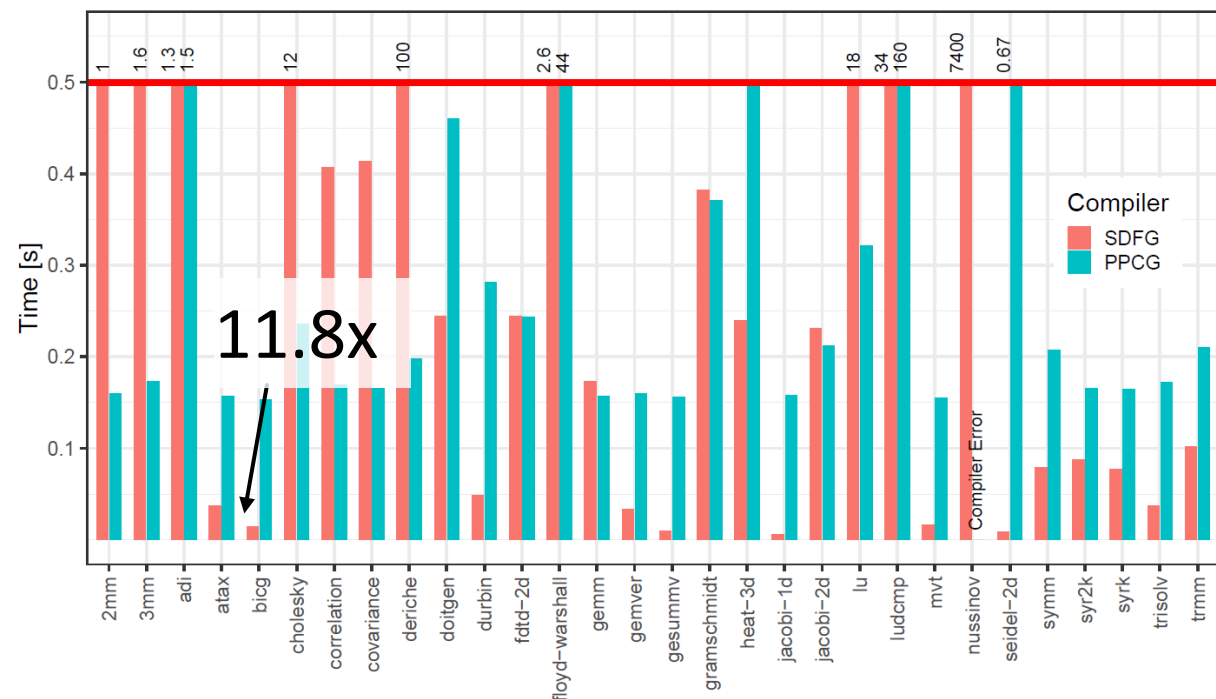
Performance Evaluation: Polybench (CPU)

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



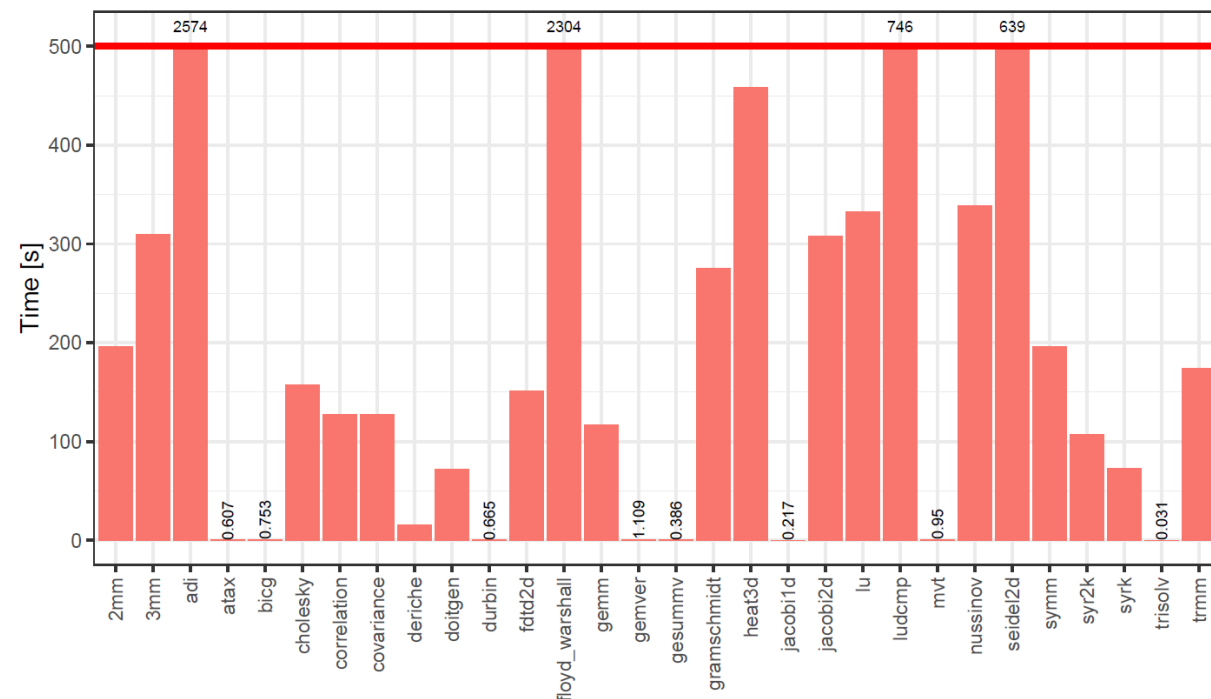
Performance Evaluation: Polybench (GPU, FPGA)

- Automatically transformed from CPU code



GPU

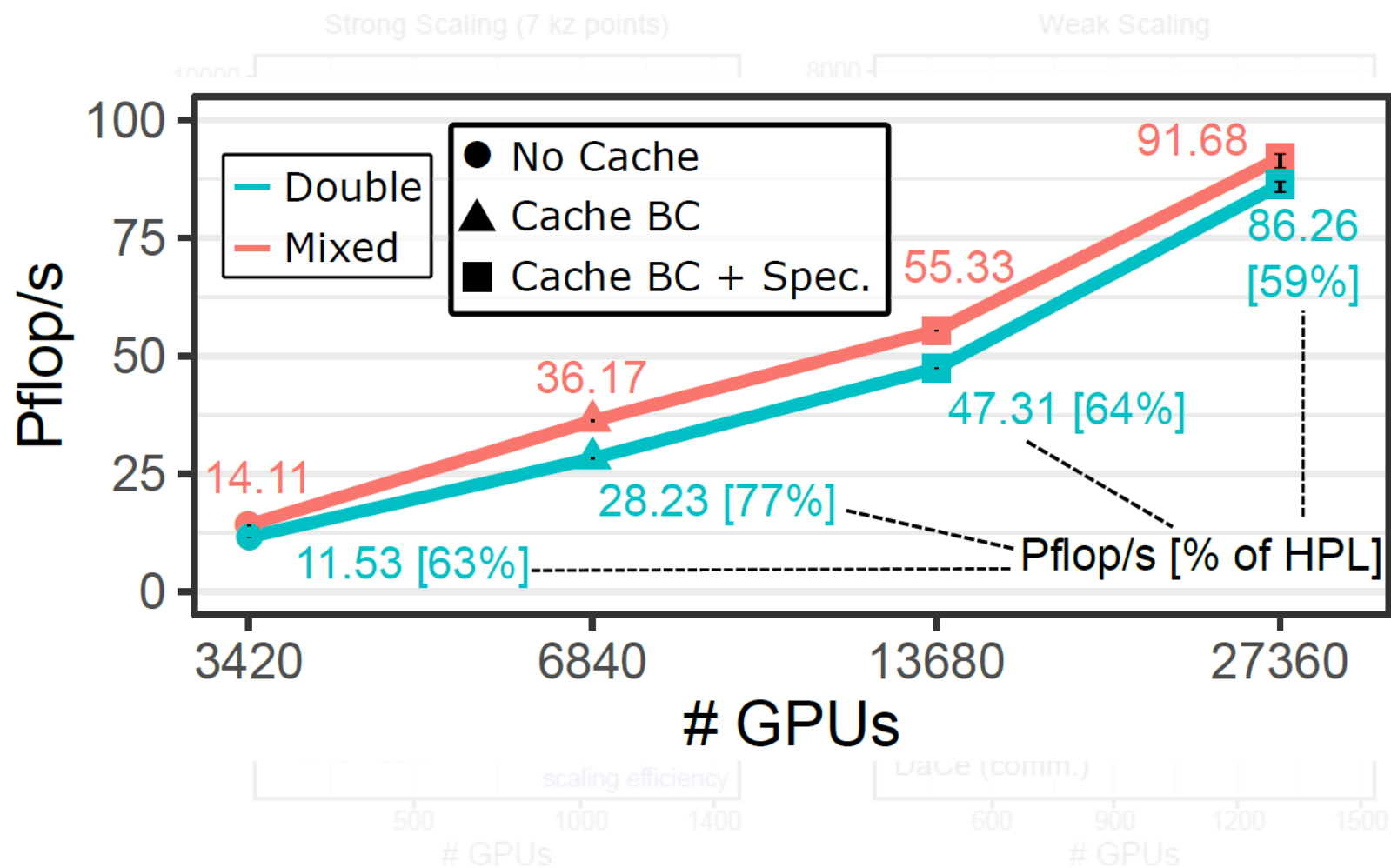
(1.12x geomean speedup)



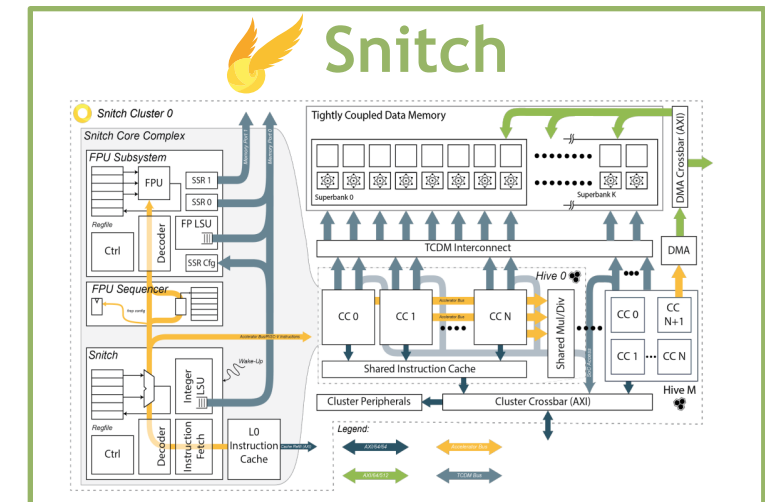
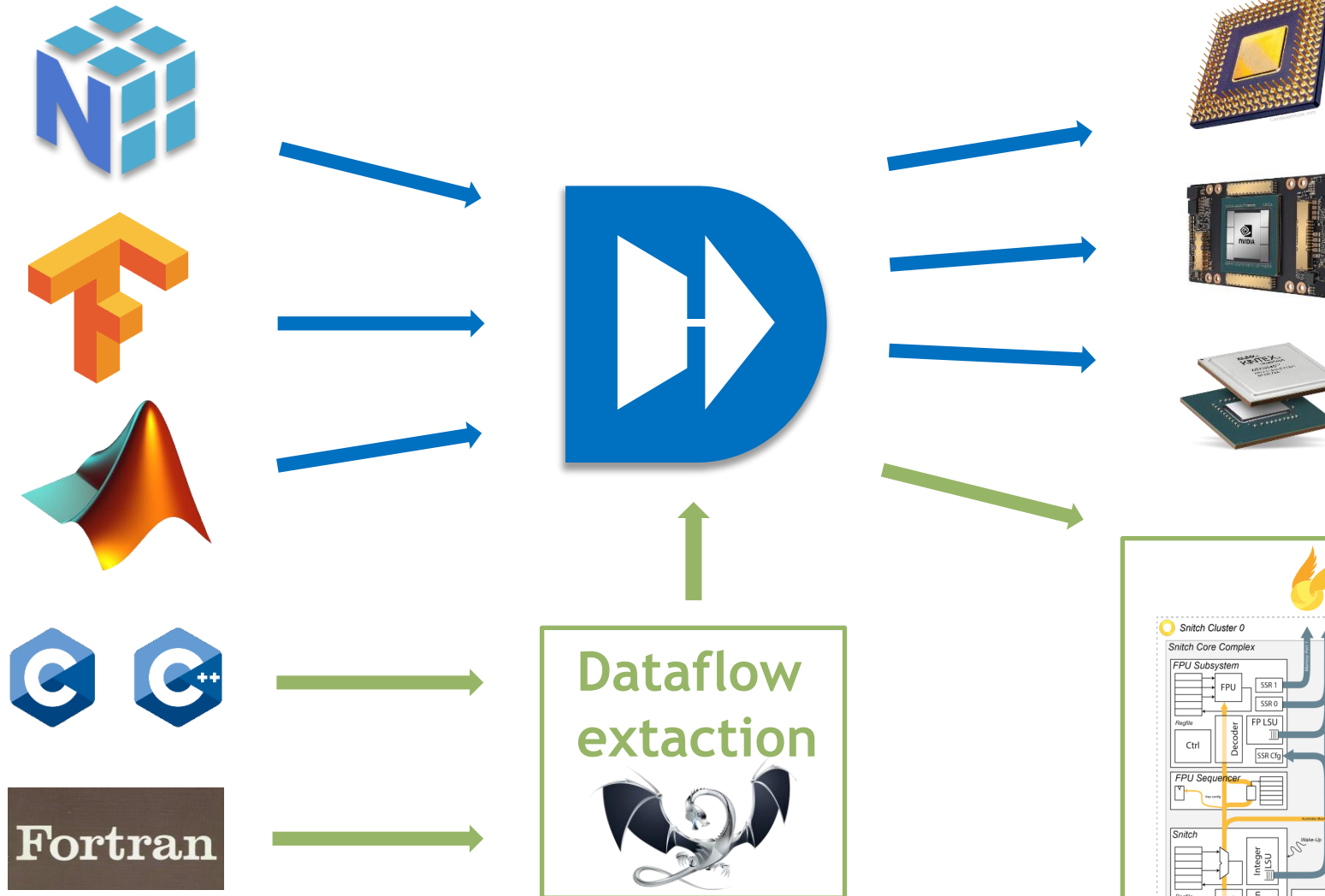
FPGA

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

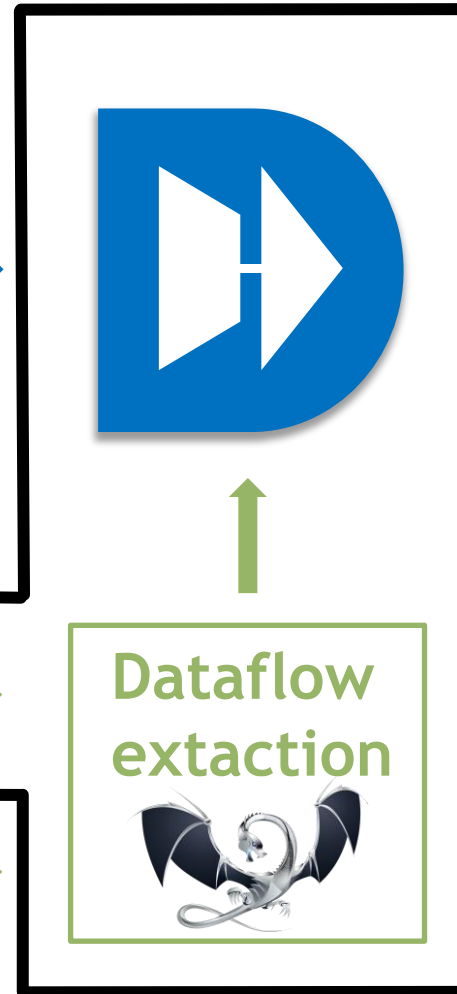
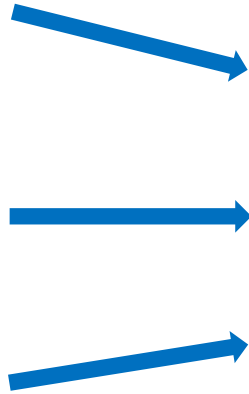
Performance Evaluation: OMEN



DaFlEx – project overview



DaFlEx – first goals



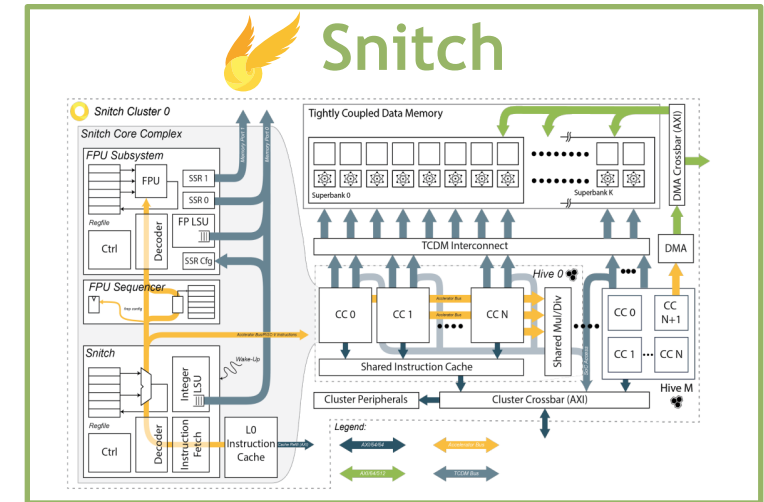
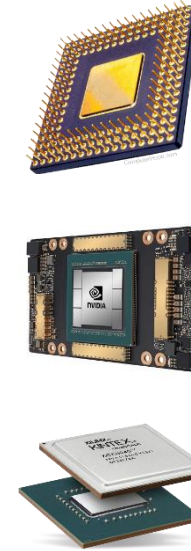
C2DaCe



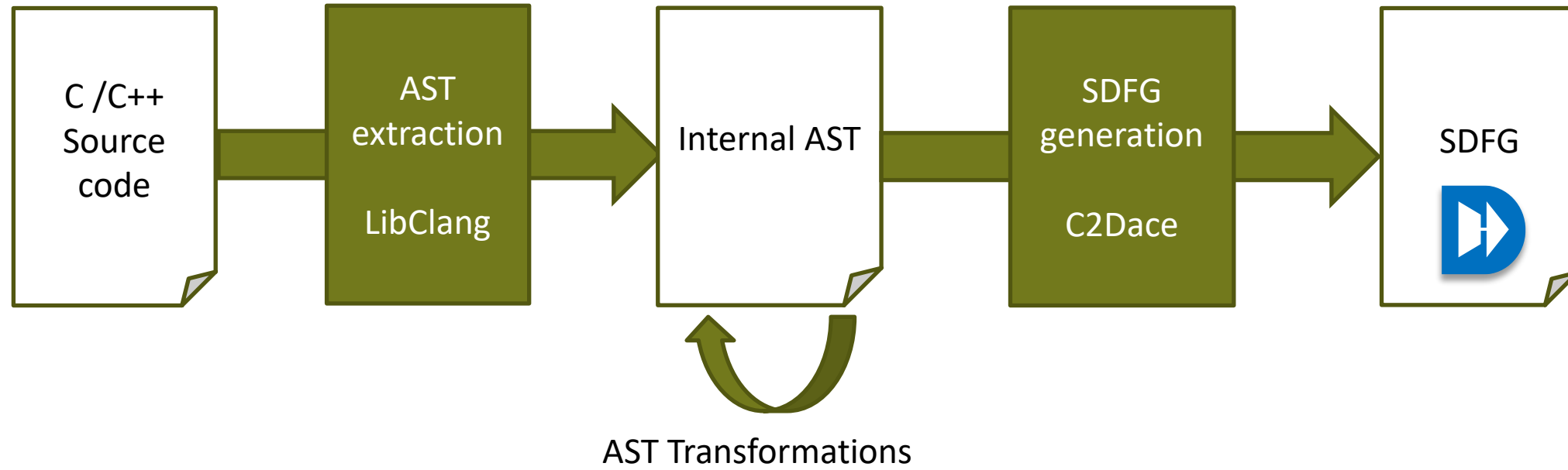
Fortran



Dataflow
extaction

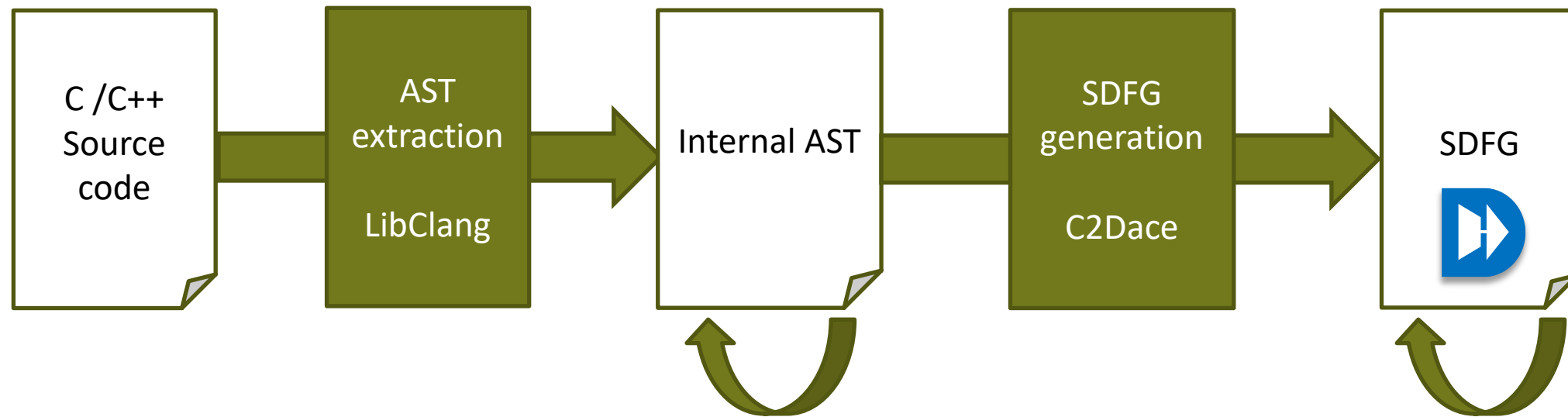


Dataflow extraction pipeline for C/C++



- Array index extraction
- Function call isolation
- Separating variable declaration and initialization

Dataflow extaction pipeline for C/C++



Candidates:

- Polybench
- LULESH
- C HPC codes?

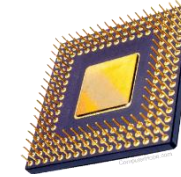
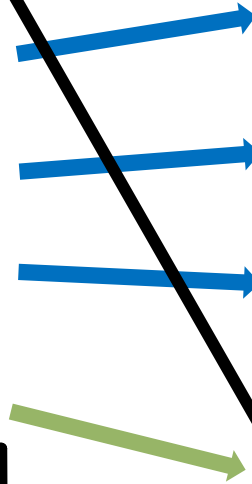
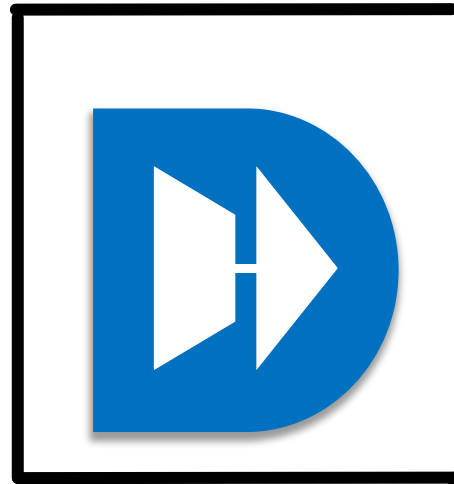
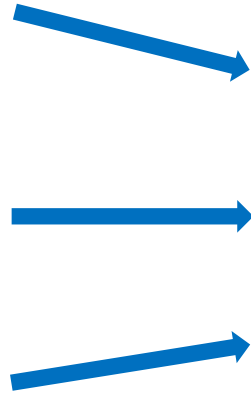
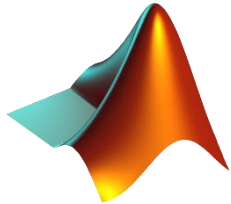
SDFG Transformations

- Symbol promotion
- Create WCR accesses where necessary
- **Loop to Map**

C2DACE

C Language	SDFG Equivalent	C Language	SDFG Equivalent
Declarations and Types		Statements	
Primitive data type	Scalar data container	Compound (blocks)	Nested SDFG
Type qualifier (const,volatile)	Implicit, hints provided are inherent to the SDFG representation	Branching (if,switch)	Branch conditions on state transition edges
Storage class (static,auto,register, etc.)		Iteration (for, while,do..while)	Nested SDFG for compound state-ment, with states and state transi-tions for loop logic
Pointer		Function flow	Edge to SDFG exit state, using helpervariable to
Array	Array data container	C Language	
Structure	Data container	SDFG Equivalent	
Union	Data container	Parallelism	
		Parametric Map scope	
C Language	SDFG Equivalent	Functions	
Expressions and Assignments		Function calls (with source)	Nested SDFG for content, memlets reduce shape of inputs and outputs
Operators (e.g., Unary,Binary)	Tasklet with incoming and outgoingmemlets for read/written operands	External/Library calls	Tasklet with library state
Compound assignments	Tasklet	Recursion	Unsupported
Array expression	Memlet	Function pointers	No equivalent, unsupported

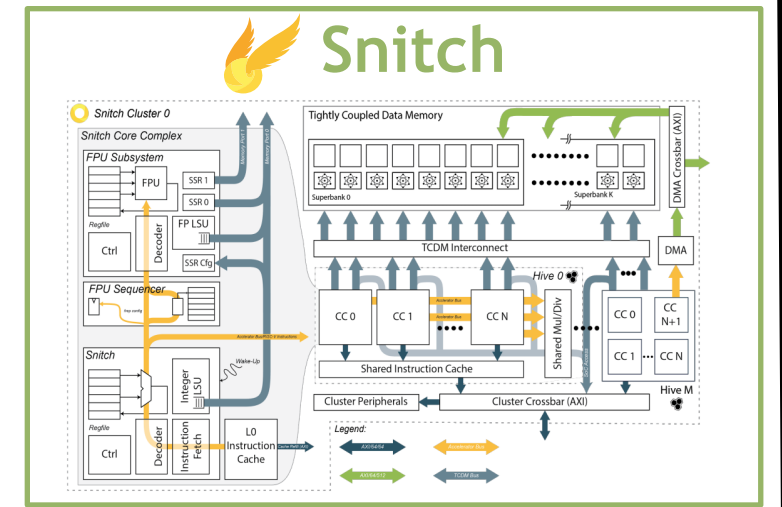
DaFlEx – first goals



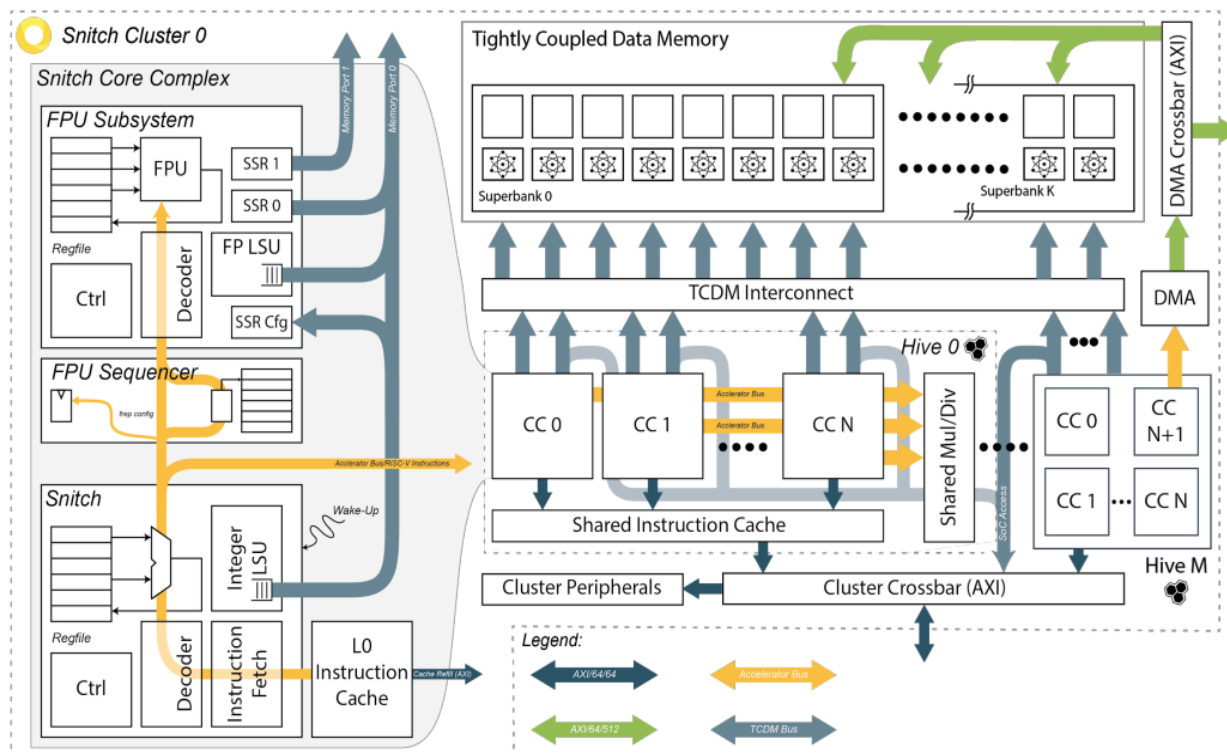
Dataflow
extaction



DaCe
for
Snitch



Snitch



RV32IMAFD Xfrep Xssr Xsdma

32 bit RISC-V
ISA

Multiplications
Divisions

Floating-Point

Atomics

Custom Extensions

- Streaming **Semantic Registers**
Register-memory interface
- **DMA Engine**
Asynchronous Data Movement
- Floating Point **Repetition**
Hardware loops for the FPU

DaCe code generation for the Snitch

Goal – create a new codegeneration backend to leverage Snitch’s custom extensions for data movement, specifically:

TCDM/DMA Double buffering from host memory to TCDM

SSR Hot-loop without load/store overhead

Frep Hardware loops for the FPU

Concept – from an SDFG emit the SSR configuration and express parallelism using OpenMP

```

void __program_daxpy_internal( ... ) {
    builtin_ssr_setup_1d_r(0, ... );
    builtin_ssr_setup_1d_r(1, ... );
    builtin_ssr_enable();
    #pragma omp parallel for
    for (auto i = 0; i < N; i += 1) {
        double in_A = A;
        double in_X = builtin_ssr_pop(0);
        double in_Y = builtin_ssr_pop(1);
        double out;
        // Tasklet code (multiplication)
        out = ((in_A * in_X) + in_Y);
        builtin_ssr_push(2, out);
    }
    builtin_ssr_disable();
}
  
```


Thank you!

stateFOR449

