

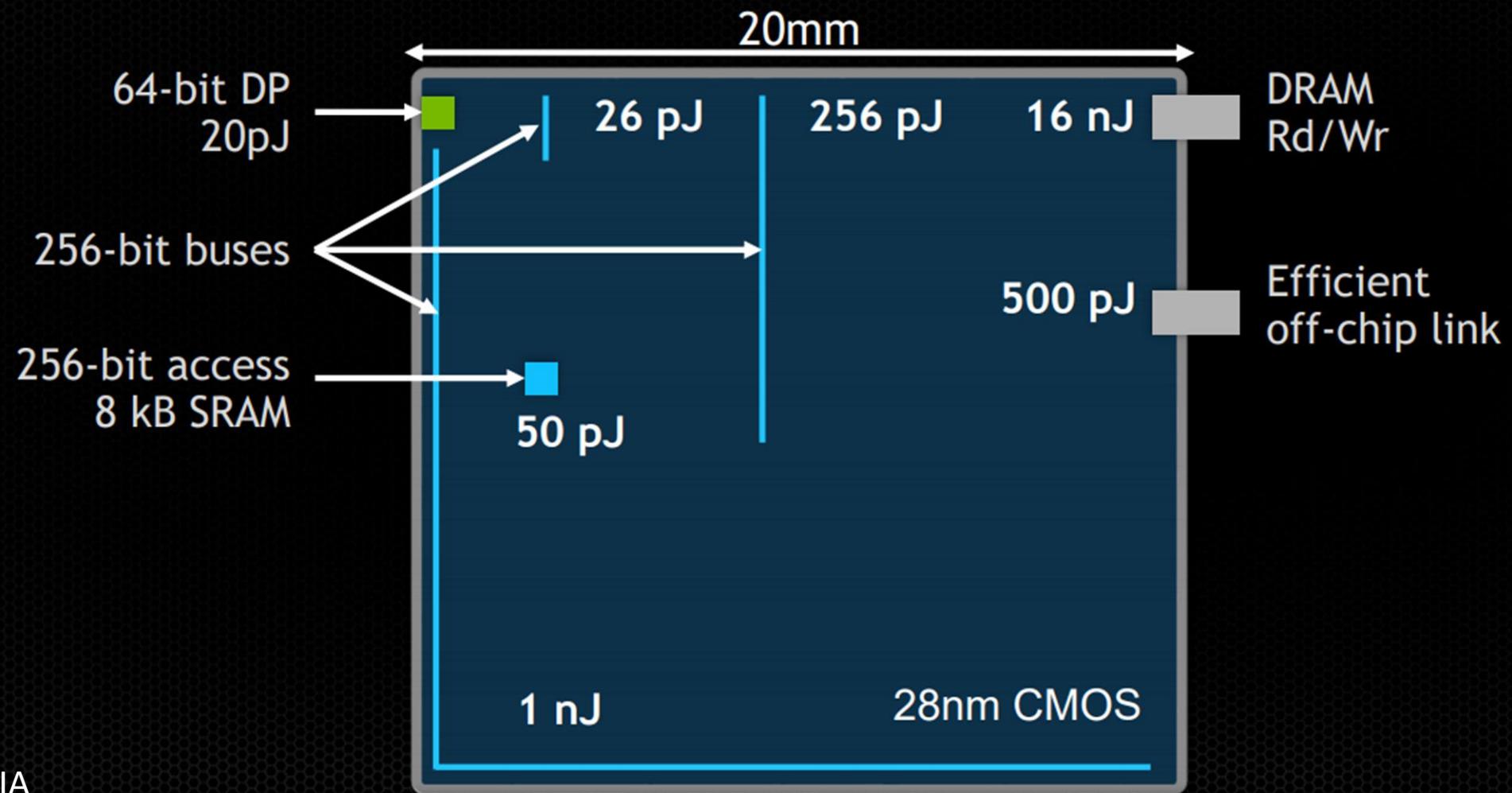
# Data-Centric Python - Productivity, portability and all with high performance!

Torsten Hoefler, keynote at the Russian Supercomputing Days (virtual)

Tal Ben-Nun, Alexandros Ziogas, Johannes de Fine Licht, Tiziano de Matteis, Timo Schneider, Andreas Kuster, Manuel Burger, Philip Schaad, Dominic Hofer and the whole DAPP team @ SPCL



# Communication Dominates Arithmetic



## Trends in Data Locality Abstractions for HPC Systems

Didem Unat, Anshu Dubey, Torsten Hoefer, John Shalf, Mark Abraham, Mauro Bianco, Bradford L. Chamberlain, Romain Cledat, H. Carter Edwards, Hal Finkel, Karl Fuerlinger, Frank Hannig, Senior Member, IEEE, Emmanuel Jeannot, Amir Kamil, Jeff Keesler, Paul H J Kelly, Vitus Leung, Hatem Ltaief, Naoya Maruyama, Chris J. Newburn, and Miquel Pericás

**Abstract**—The cost of data movement has always been an important concern in high performance computing (HPC) systems. It has now become the dominant factor in terms of both energy consumption and performance. Support for expression of data locality has been explored in the past, but those efforts have had only modest success in being adopted in HPC applications for various reasons. However, with the increasing complexity of the memory hierarchy and higher parallelism in emerging HPC systems, locality management has acquired a new urgency. Developers can no longer limit themselves to low-level solutions and ignore the potential for productivity and performance portability obtained by using locality abstractions. Fortunately, the trend emerging in recent literature on the topic alleviates many of the concerns that got in the way of their adoption by application developers. Data locality abstractions are available in the forms of libraries, data structures, languages and runtime systems; a common theme is increasing productivity without sacrificing performance. This paper examines these trends and identifies commonalities that can combine various locality concepts to develop a comprehensive approach to expressing and managing data locality on future large-scale high-performance computing systems.

**Index Terms**—Data locality, programming abstractions, high-performance computing, data layout, locality-aware runtimes

### 1 INTRODUCTION

THE computing industry has entered a period of technology transition as we strive for the next 1,000 x performance improvement over the previous generation of petaflops-scale computing platforms. Over the past 30 years, we have come to expect a 1,000 x increase in HPC system

• D. Unat is with the Department of Computer Engineering, Koç University, Istanbul 34450, Turkey. E-mail: dunat@ku.edu.tr.

• A. Dubey is with Argonne National Laboratory, Lemont, IL 60439. E-mail: adubey@anl.gov.

• T. Hoefer is with ETH Zürich, Zürich 8092, Switzerland. E-mail: thof@inf.ethz.ch.

• J. Shalf is with Lawrence Berkeley National Laboratory, Berkeley, CA 94720. E-mail: jshalf@lbl.gov.

• M. Abraham is with KTH Royal Institute of Technology, Solna 17121, Sweden. E-mail: mjab@kth.se.

• M. Bianco is with Swiss National Supercomputer Centre, Lugano 6900, Switzerland. E-mail: mauro.bianco@csccs.ch.

• B. Chamberlain is with Cray Inc., Seattle, WA 98164. E-mail: brad@cray.com.

• R. Cledat is with Intel Cooperation, Santa Clara, CA 95050. E-mail: rcedat@intel.com.

• C. Edwards is with Sandia National Laboratories, Albuquerque, NM 87185. E-mail: hcedars@sandia.gov.

• H. Finkel is with Argonne National Laboratory, Argonne, IL 60439. E-mail: hfinkel@anl.gov.

• K. Fuerlinger is with Ludwig-Maximilians-Universität München, Munich D-80538, Germany. E-mail: Karl.Fuerlinger@nm.iifi.lmu.de.

• F. Hannig is with University of Erlangen-Nuremberg, Erlangen 91058, Germany. E-mail: frank.hannig@fau.de.

Manuscript received 2 June 2016; revised 12 Apr. 2017; accepted 14 Apr. 2017. Date of publication 10 May 2017; date of current version 13 Sept. 2017. (Corresponding author: Didem Unat.)

Recommended for acceptance by U.V. Catalyurek.

For information on obtaining reprints of this article, please send e-mail to: reprints@iee.org, and reference the Digital Object Identifier below. Digital Object Identifier no. 10.1109/TPDS.2017.2703149

• V. Leung is with INRIA Bordeaux Sud-Ouest, Talence 33405, France. E-mail: emmanuel.jeannot@inria.fr.

• A. Kamil is with University of Michigan, MI 48109, and also with Lawrence Berkeley National Laboratory, Berkeley, CA 94720. E-mail: akamil@umich.edu.

• J. Keesler is with Lawrence Livermore National Laboratory, Livermore, CA 94550. E-mail: jkeeler1@llnl.gov.

• P. Kelly is with Imperial College London, London, United Kingdom. E-mail: p.kelly@imperial.ac.uk.

• V. Leung is with Sandia National Laboratories, Albuquerque, NM 87185. E-mail: vleung@sandia.gov.

• H. Ltaief is with King Abdullah University of Science and Technology, Thuwal 23955, Kingdom of Saudi Arabia. E-mail: ltaief\_hatem@yahoo.fr.

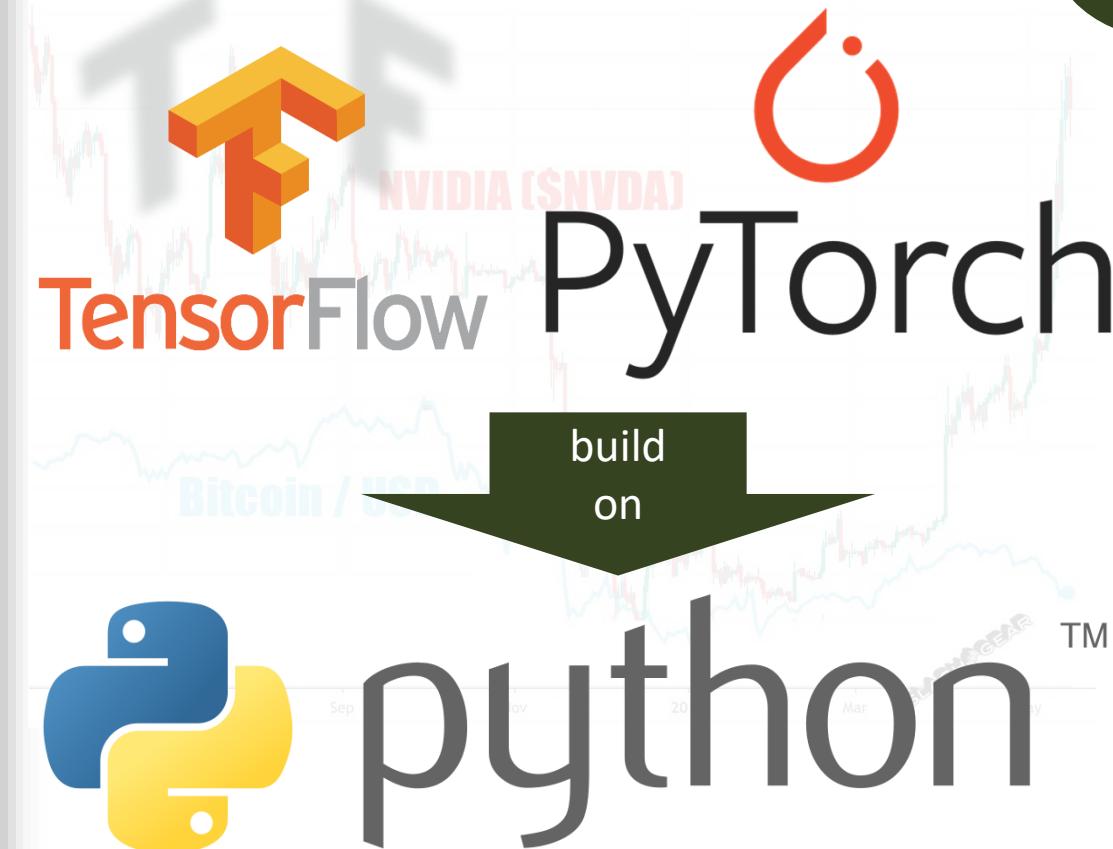
• N. Maruyama is with RIKEN, Kobe, Hyogo 650-0047, Japan. E-mail: nmaruyama@riken.jp.

• C. J. Newburn is with Nvidia Corporation, Santa Clara, CA 95050. E-mail: chris.newburn@intel.com.

• M. Pericás is with Chalmers University of Technology 41296, Göteborg, Sweden. E-mail: miquelp@chalmers.se.

How do we bring these ideas into productive performance computing?

What is the fastest growing high-performance user-base to adopt GPUs even faster than HPC?



# Comparing fast python implementations

51 kernels from 9 domains

Learning (6) 

LinAlg (12) 

Chemistry (4) 

Signals (3) 

Physics (9) 

Graphs (2) 

Weather (2) 

Solver (10) 

Other (3) 

Yours?

Metrics



Performance

Meet  
NPBench



Productivity

Frameworks



NumPy baseline



Pythran



Numba



CuPy



DaCe

Python is the  
language of  
science and  
engineering

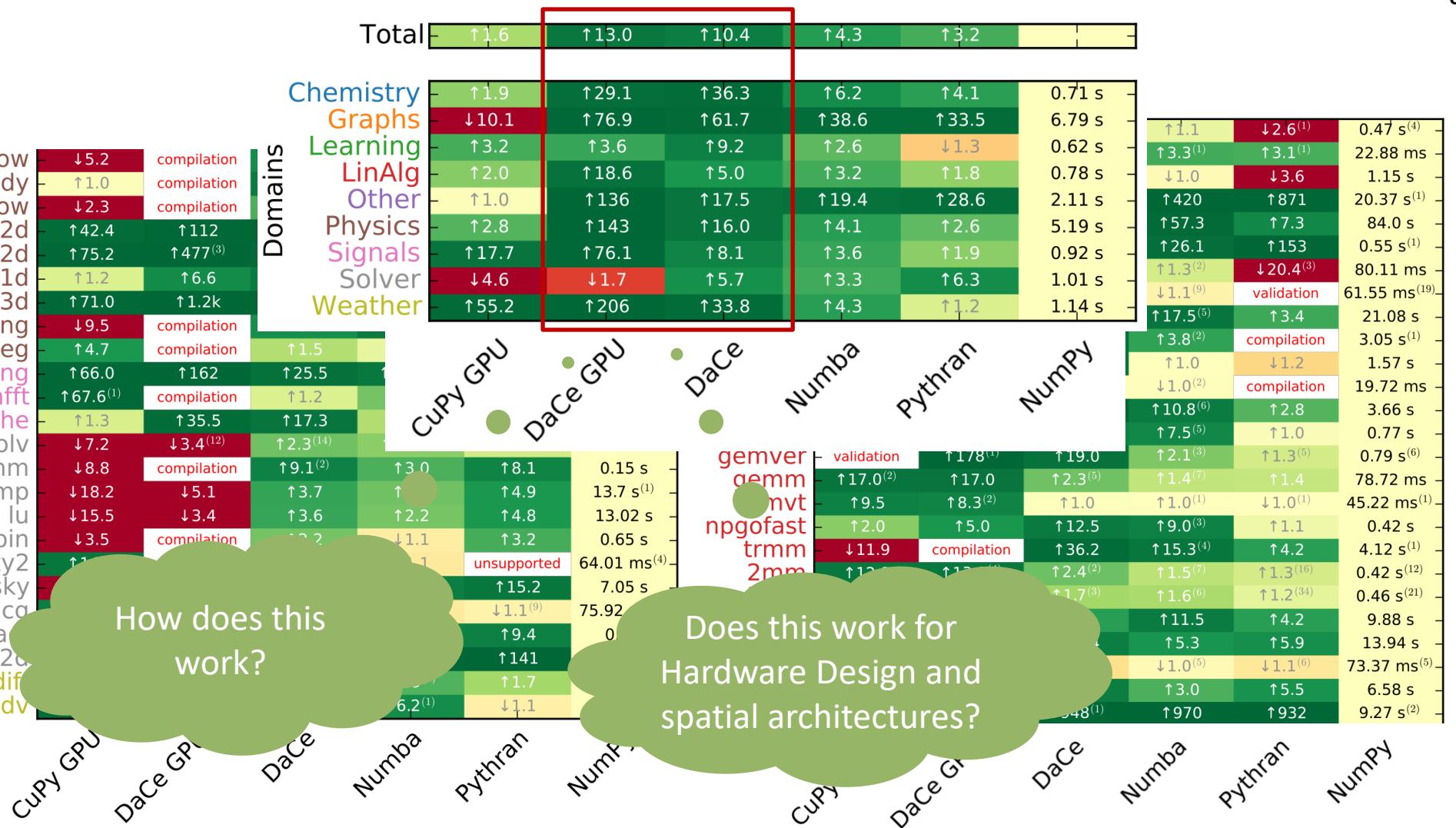
We need a fair  
comparison  
between  
frameworks



# NPBench results

Machine with two 16-core  
Intel Xeon Gold 6130 processors  
and an Nvidia V100 GPU with  
32GB of memory

chanflow	$\downarrow 5.2$	compilation
nbody	$\uparrow 1.0$	compilation
cavtflow	$\downarrow 2.3$	compilation
fdtd_2d	$\uparrow 42.4$	$\uparrow 112$
jacobi2d	$\uparrow 75.2$	$\uparrow 477^{(3)}$
jacobi1d	$\uparrow 1.2$	$\uparrow 6.6$
heat3d	$\uparrow 71.0$	$\uparrow 1.2k$
sselfeng	$\downarrow 9.5$	compilation
coninteg	$\uparrow 4.7$	compilation
clipping	$\uparrow 66.0$	$\uparrow 162$
sthamfft	$\uparrow 67.6^{(1)}$	compilation
deriche	$\uparrow 1.3$	$\uparrow 35.5$
trisolv	$\downarrow 7.2$	$\downarrow 3.4^{(12)}$
gramschm	$\downarrow 8.8$	compilation
ludcmp	$\downarrow 18.2$	$\downarrow 5.1$
lu	$\downarrow 15.5$	$\downarrow 3.4$
durbin	$\downarrow 3.5$	compilation
cholesky2	$\uparrow 1$	
cholesky		
bicg		
a		
seidel2d		
hdif		
vadv		



How does this work?

Does this work for  
Hardware Design and  
spatial architectures?

# Data-Centric Python



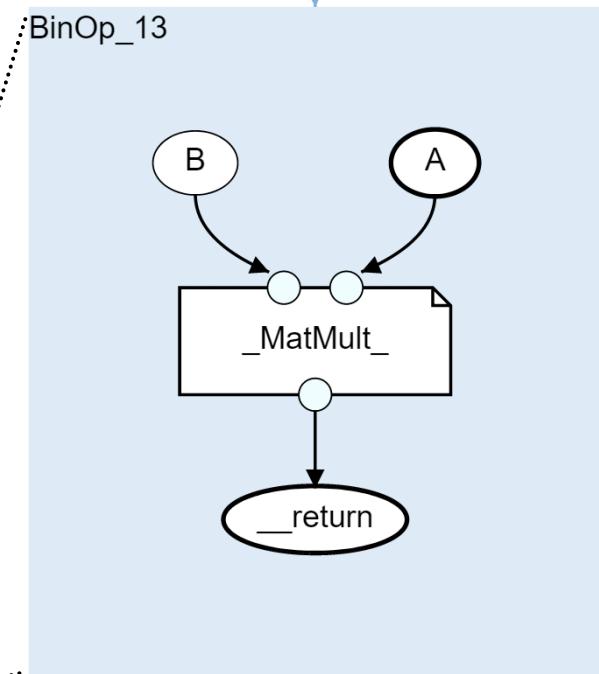
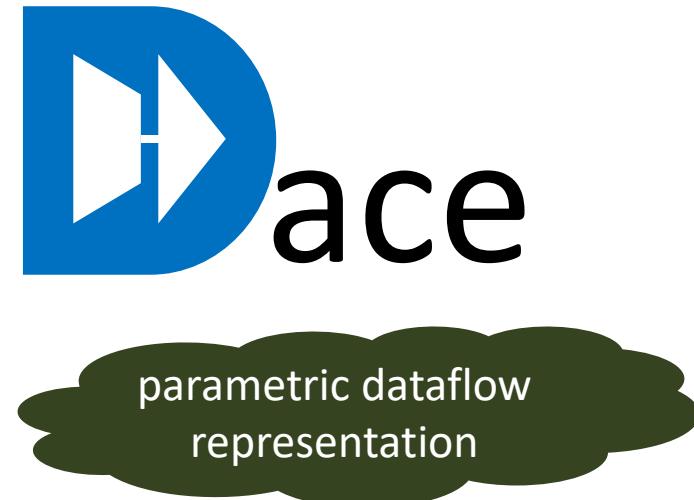
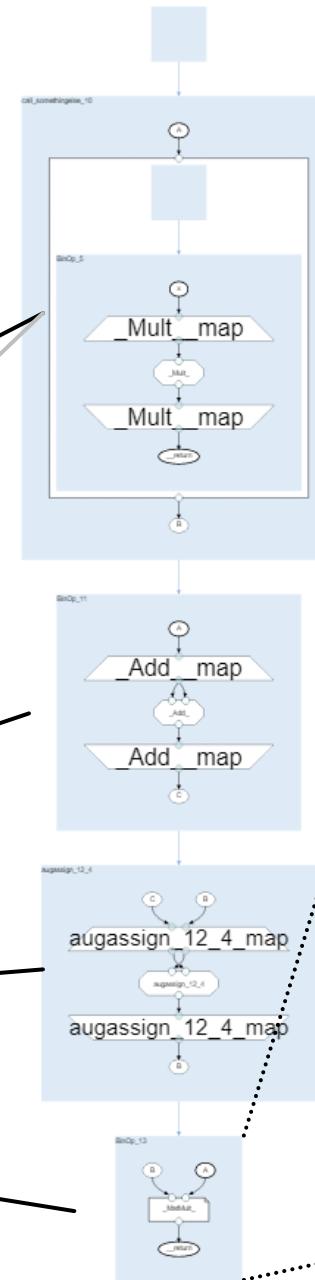
# python

```
import dace

@dace.program
def somethingelse(x):
    return x * 5

@dace.program
def example(A: dace.float64[20, 20]):
    B = somethingelse(A)
    C = A + A
    B += C
    return np.dot(B, A)
```

imperative code



# Data-Centric Python Vision – Performance Portability

10s of SLOC

main Scientist



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



DSLs



MATLAB



Applied Scientist



translate DSL into parametric dataflow graphs

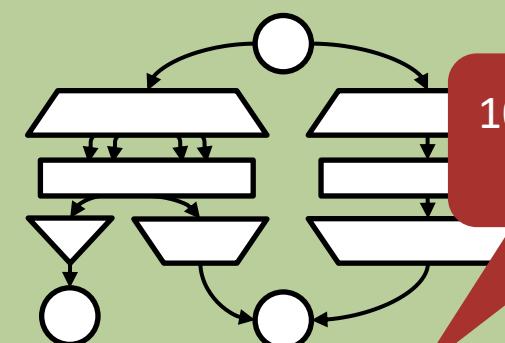
SDFG Builder API

Multi-Level Library Nodes

Performance Engineer



100s of reusable SLOC



Parametric Dataflow Graphs (SDFG)

SPCL

Dataflow Programming in DaCe

Transformed Dataflow

Performance Results

Graph Transformations (API, Interactive)



1000s of auto-generated SLOC

Specialized Code Generation

Runtime

CPU Code

GPU Code

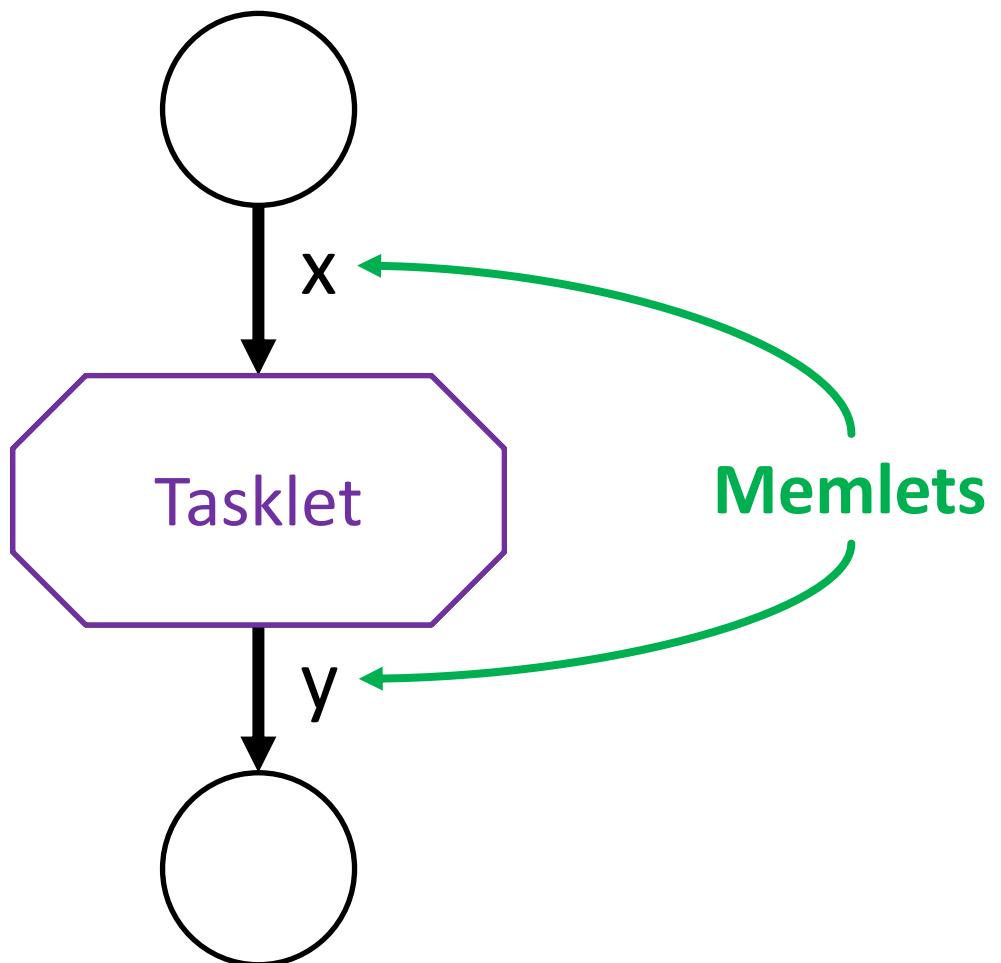
FPGA Code

C++ code generation/runtime

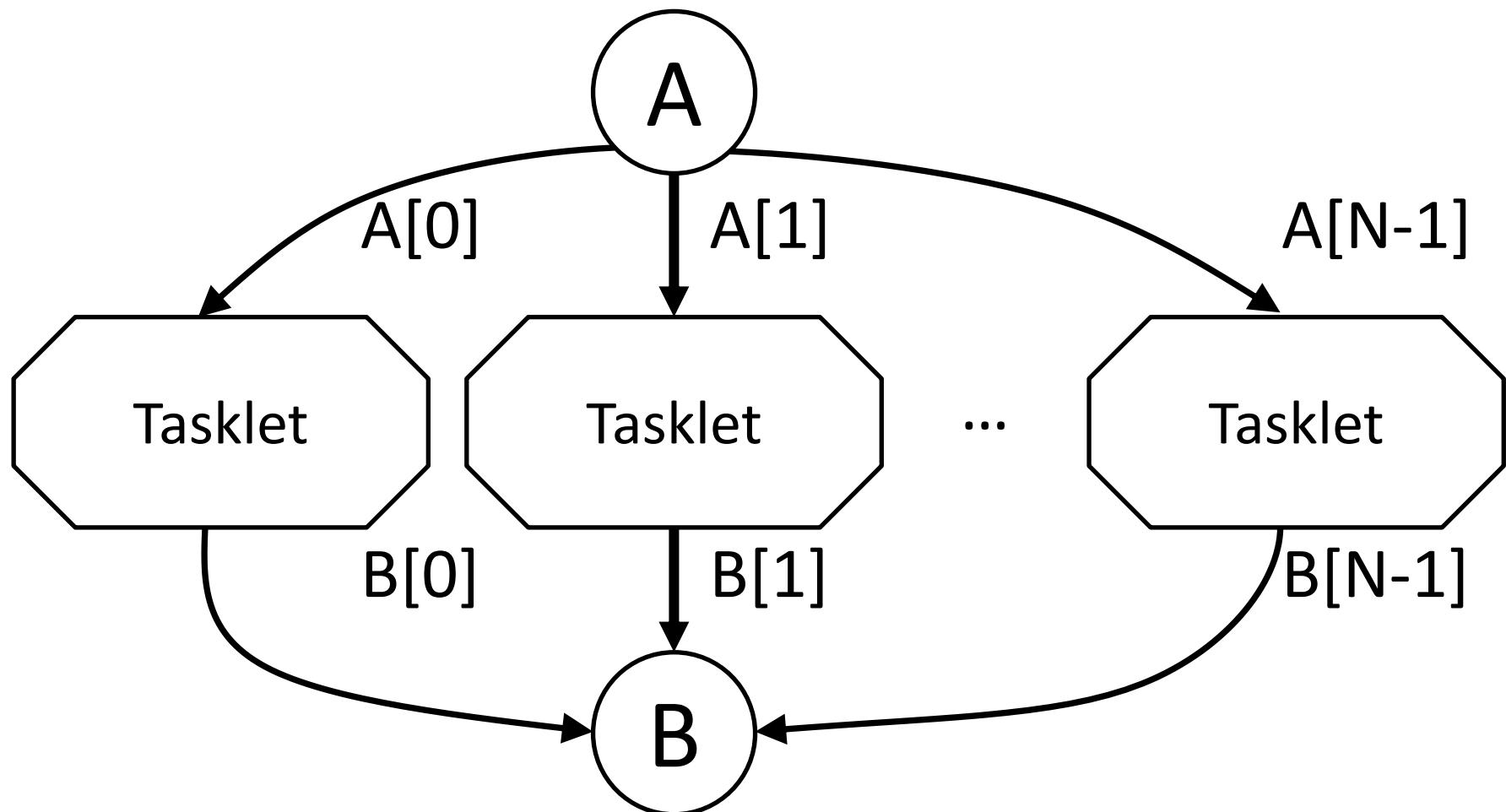
# Dataflow Programming in DaCe

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

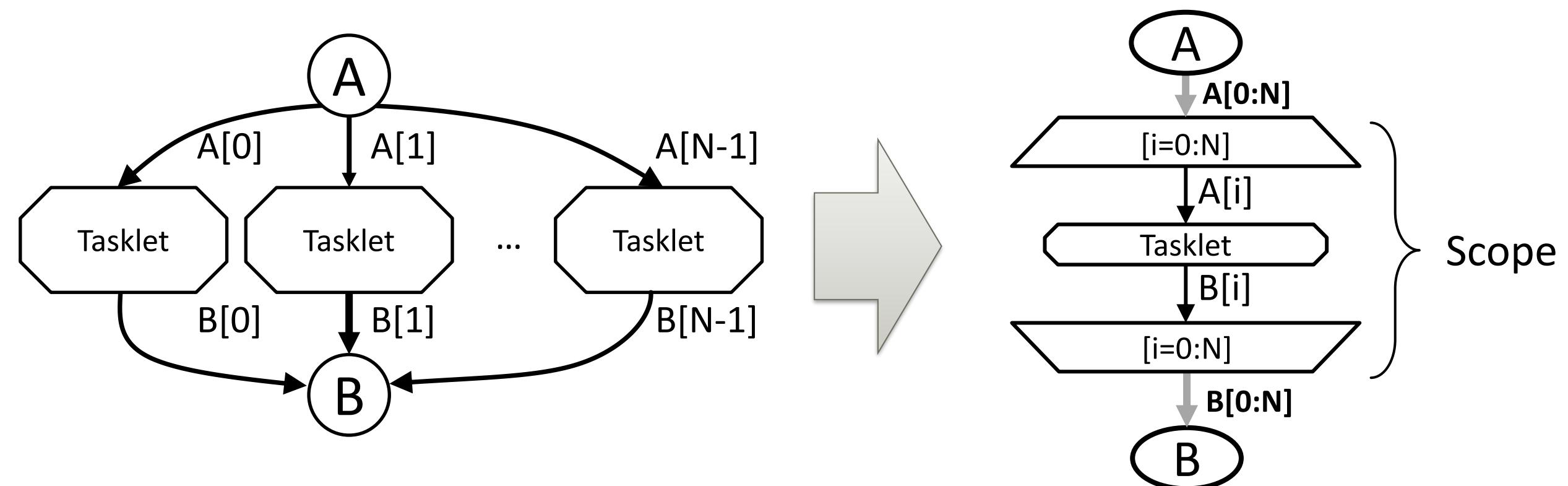
# Dataflow Programming in DaCe



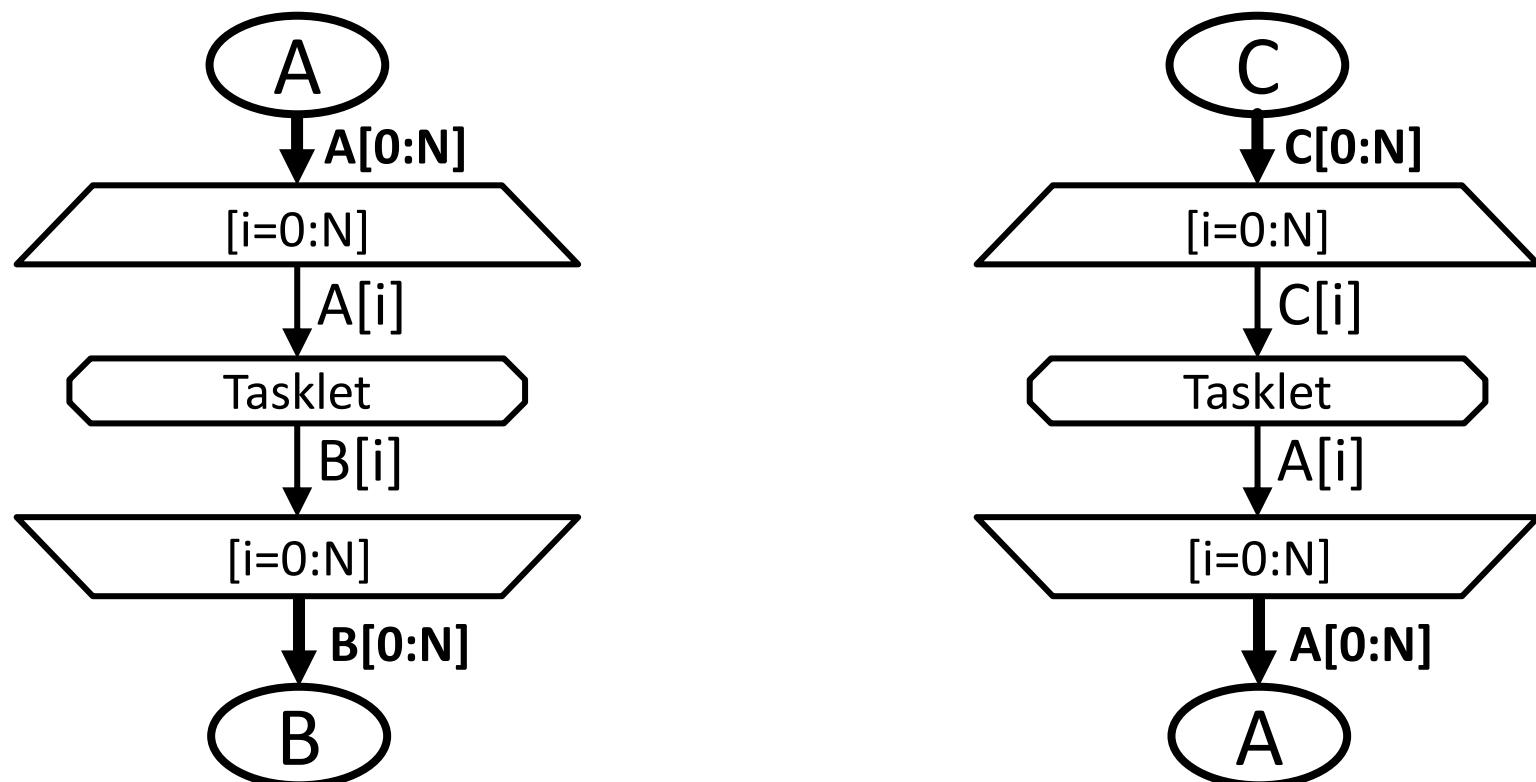
# Parallel Dataflow Programming



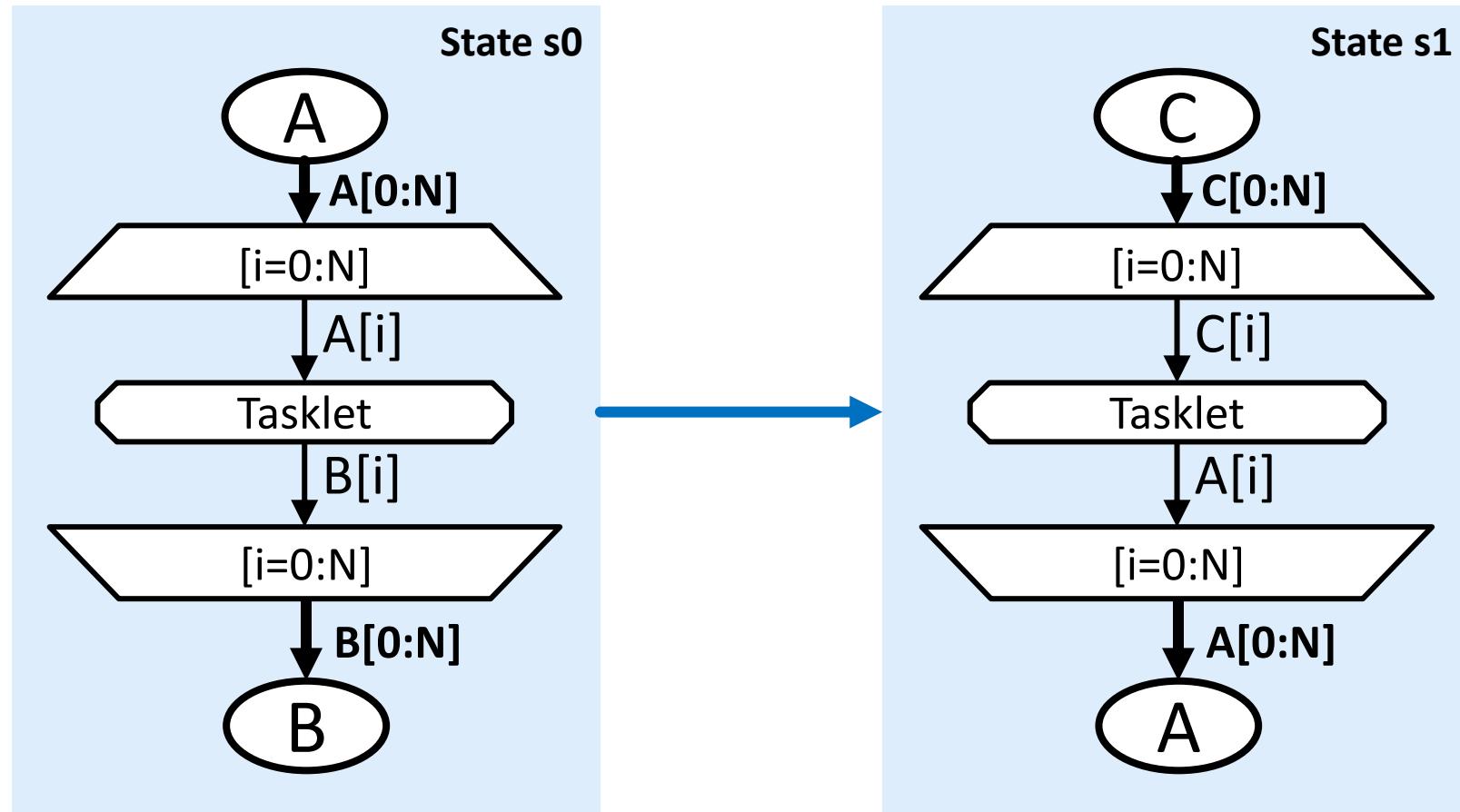
# Parallel Dataflow Programming



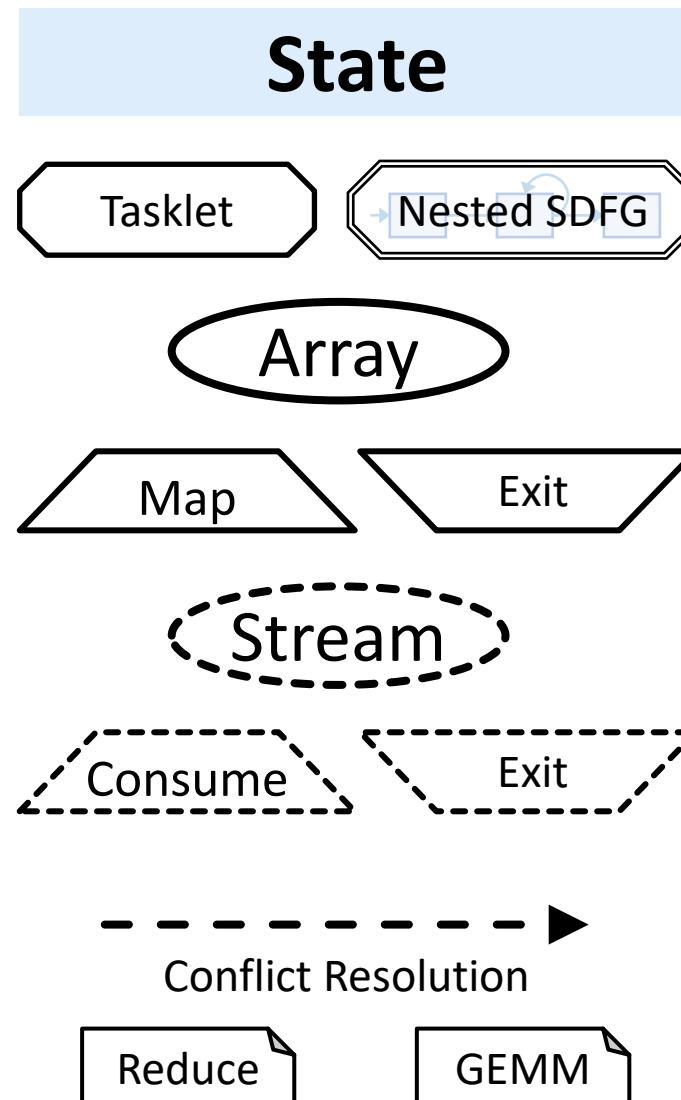
# Stateful Parallel Dataflow Programming



# Stateful Parallel Dataflow Programming



# Meet the Nodes



State machine element

Fine-grained computational block

N-dimensional data container

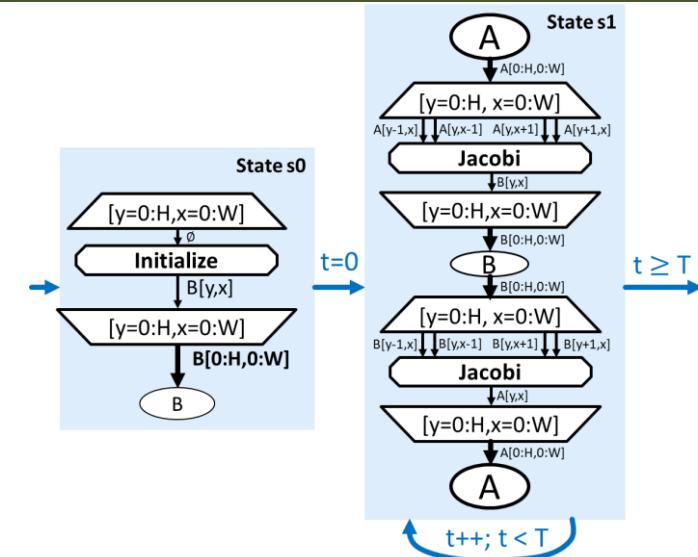
Parametric graph abstraction for parallelism

Streaming data container

Dynamic mapping of computations on streams

Defines behavior during conflicting writes

Customizable computation with multiple implementations

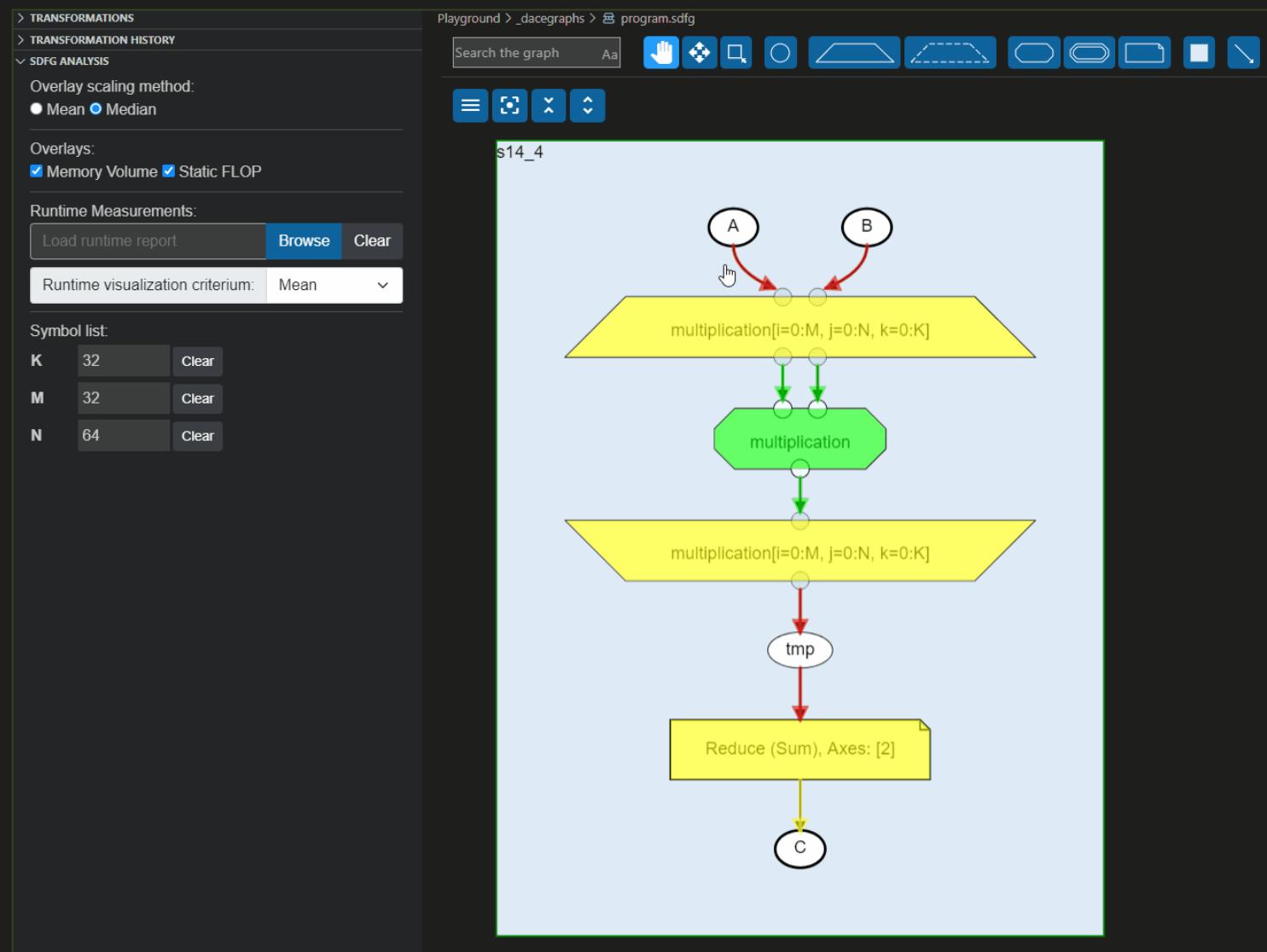


# Programmer/Performance Engineer view: Visual Studio Code Integration

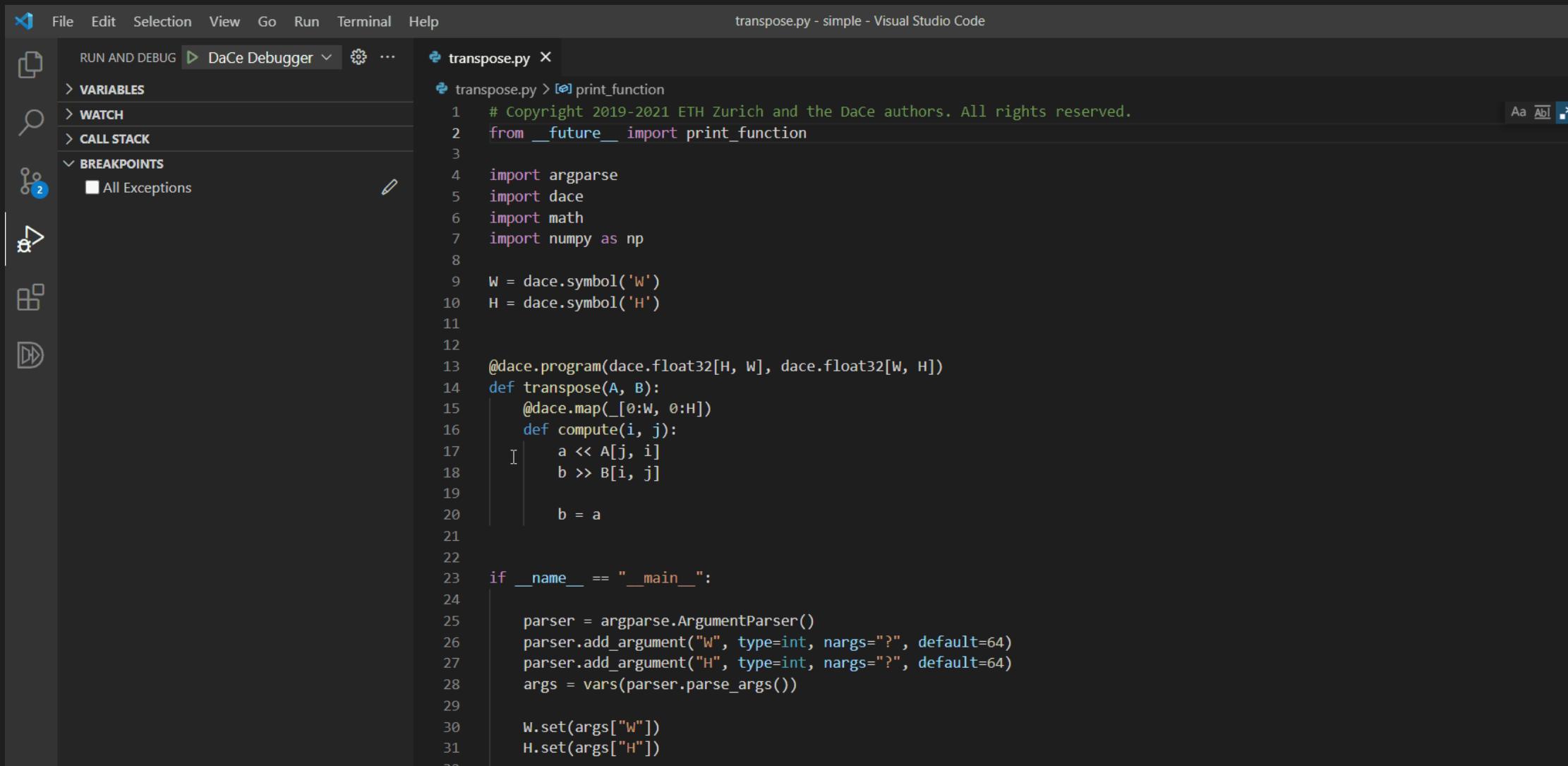
The screenshot illustrates the SPCL Visual Studio Code extension's integration for programmer and performance engineers. It features three main panes:

- Left Pane:** A code editor showing Python code for a GEMM operation using DACE (Data Access Configuration Environment). The code defines symbols for M, N, K, initializes arrays A, B, and C, and performs the computation.
- Middle Pane:** A Stateful Dataflow Multigraph (SDFG) visualization for the GEMM operation. The graph shows data flow from inputs A and B through a multiplication stage to produce intermediate results and finally output C. The graph is labeled "s14\_4".
- Right Pane:** A configuration interface for the SDFG gemm transformation. It includes sections for General, arg\_names, constants\_prop, exit\_code, global\_code, init\_code, instrument, openmp\_sections, symbols, and Uncategorized. The "symbols" section lists K, M, and N as int32 types.

# Programmer/Performance Engineer view: Analyzing data flows



# Programmer/Performance Engineer view: Debugging

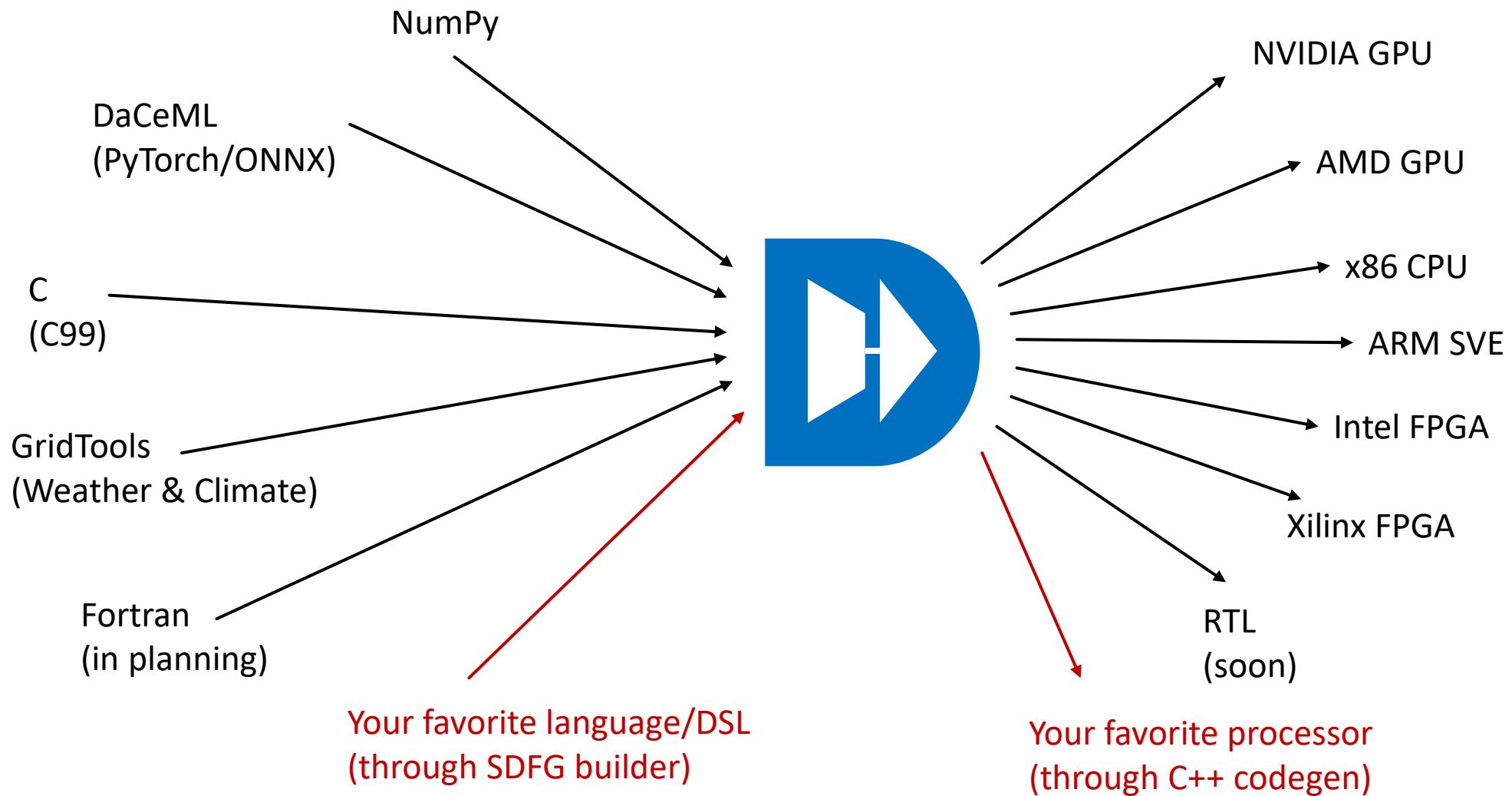


The screenshot shows the Visual Studio Code interface with the following details:

- File Bar:** File, Edit, Selection, View, Go, Run, Terminal, Help.
- Title Bar:** transpose.py - simple - Visual Studio Code.
- Left Sidebar:** RUN AND DEBUG (DaCe Debugger selected), VARIABLES, WATCH, CALL STACK, BREAKPOINTS (2 items). A tooltip for print\_function is shown above the code area.
- Code Area:** The transpose.py script is displayed. It includes imports for argparse, dace, math, and numpy, along with Dace annotations for symbol declarations and a map loop. It defines a transpose function that performs matrix transposition using a map loop and computes element-wise operations. The script also includes a main block for argument parsing and setting symbols W and H.

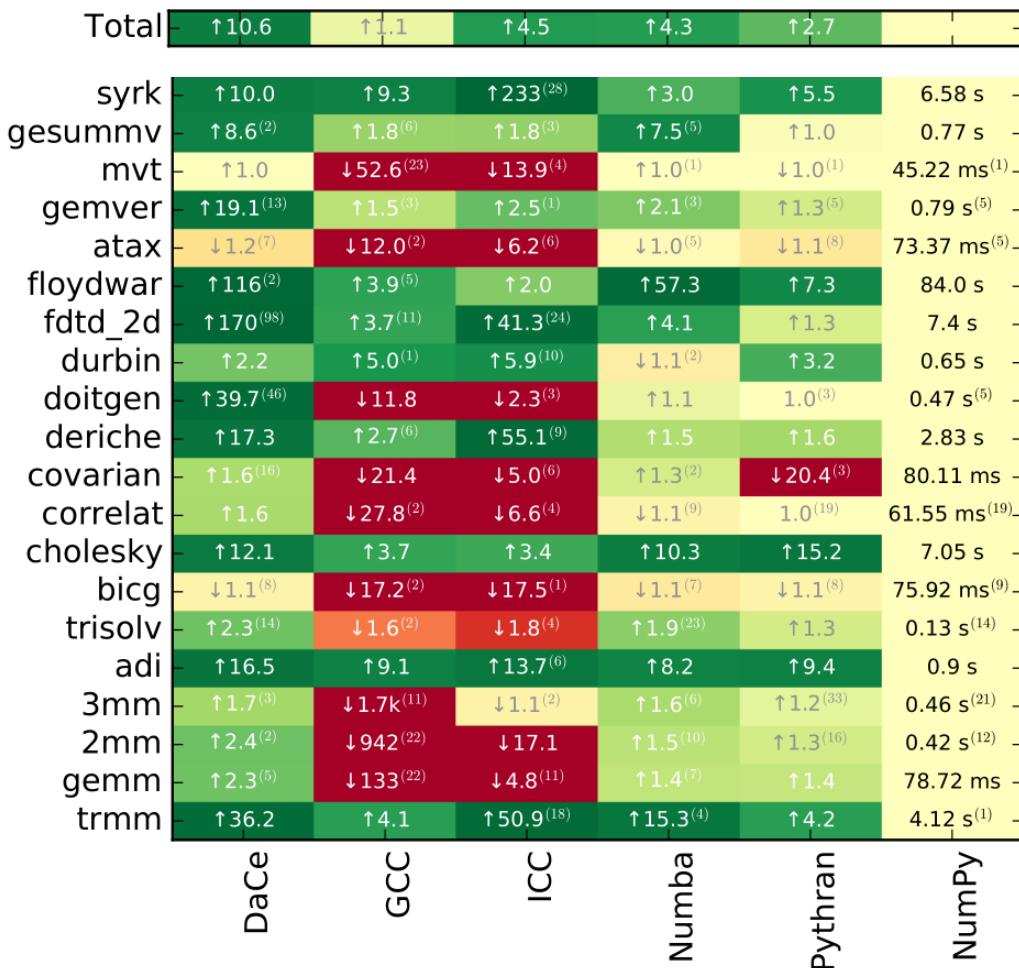
```
1 # Copyright 2019-2021 ETH Zurich and the DaCe authors. All rights reserved.
2 from __future__ import print_function
3
4 import argparse
5 import dace
6 import math
7 import numpy as np
8
9 W = dace.symbol('W')
10 H = dace.symbol('H')
11
12
13 @dace.program(dace.float32[H, W], dace.float32[W, H])
14 def transpose(A, B):
15     @dace.map(_[0:W, 0:H])
16     def compute(i, j):
17         a << A[j, i]
18         b >> B[i, j]
19
20         b = a
21
22
23 if __name__ == "__main__":
24
25     parser = argparse.ArgumentParser()
26     parser.add_argument("W", type=int, nargs="?", default=64)
27     parser.add_argument("H", type=int, nargs="?", default=64)
28     args = vars(parser.parse_args())
29
30     W.set(args["W"])
31     H.set(args["H"])
```

# DaCe is a versatile platform



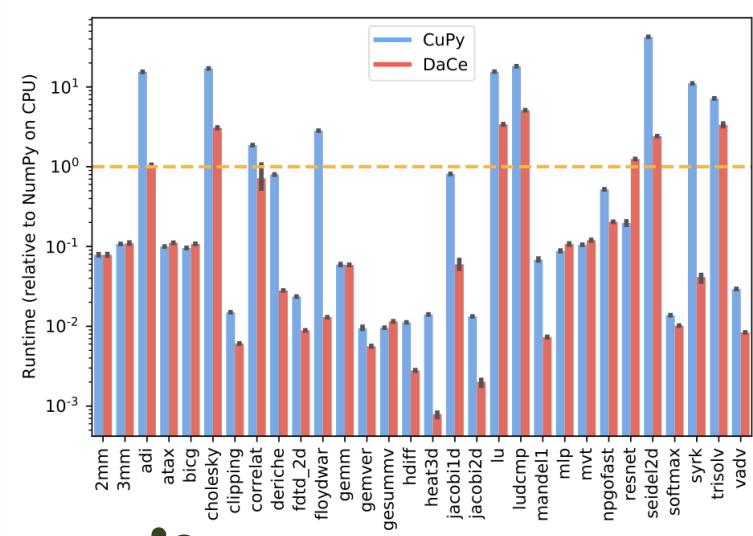
# Mapping NumPy to CPU

**Key principle: keep dataflow of vectorized code around and use it for mapping**



• •

NumPy can even beat the C versions of PolyBench codes!



• •

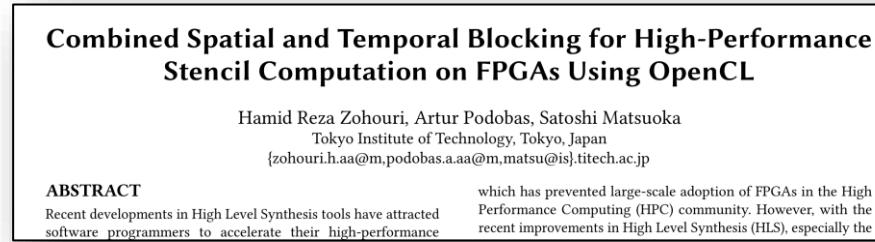
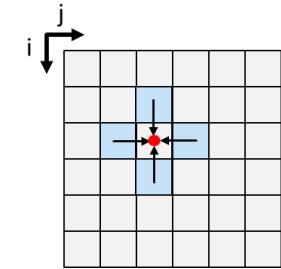
DaCe outperforms CuPy in most cases

## PolyBench/C

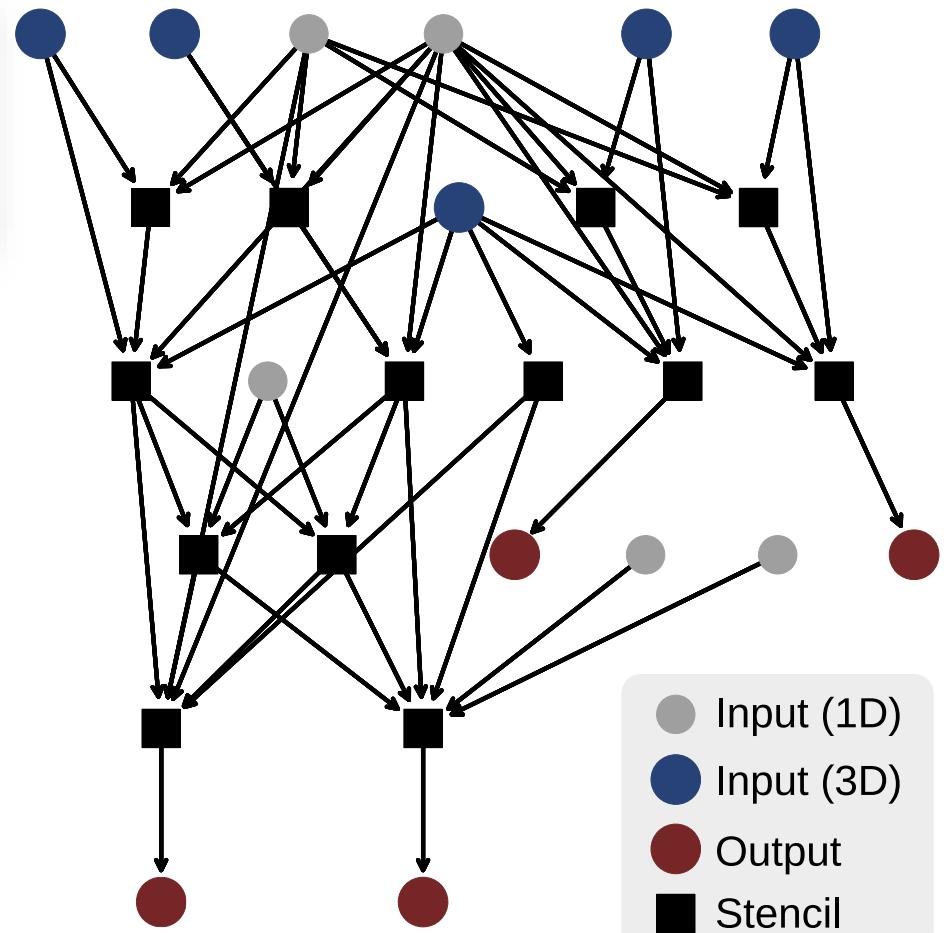
the Polyhedral Benchmark suite

# Mapping Weather/Climate Stencils (in a Python DSL) to FPGAs

**Key principle: spatial layout and pipelining using streams and delay buffers**



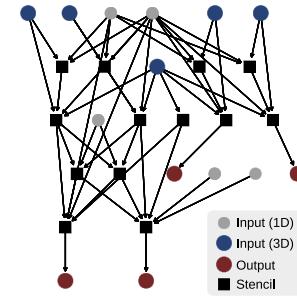
	<b>Performance</b>	<b>ALM</b>	<b>FF</b>	<b>M20K</b>	<b>DSP</b>
Total Avail.		103 M 692 K	3.7 M 2.8 M	11.7 K 8.9 K	5760 4468
Jacobi 3D (Ours)	265 GOP/s	233 K 33.6%	534 K 19.3%	1495 16.7%	784 17.6%
Jacobi 3D W=8 (Ours)	921 GOP/s	437 K 63.1%	1207 K 43.6%	2285 25.5%	3072 68.8%
Diffusion 2D W=8 (Ours)	1,313 GOP/s	449 K 64.8%	1329 K 48.0%	2565 28.6%	2304 51.6%
Diffusion 3D W=8 (Ours)	1,152 GOP/s <b>23-44% faster!</b>	567 K 81.9%	1606 K 57.9%	5357 59.8%	3072 68.8%
Diffusion 2D (Zohouri et. al. [8])	913 GOP/s	471.4 K 68.0%	1173.6 K 42.3%	2204 24.6%	3844 86.0%
Diffusion 3D (Zohouri et. al. [8])	934 GOP/s	450.5 K 65.0%	1078.2 K 38.9%	8684 97.0%	3592 80.4%



# Mapping Weather Stencils (in a Python DSL) to FPGAs

**Key principle: spatial layout and pipelining using streams and delay buffers**

318 MHz at 48% DSP utilization



	Runtime	Performance	Peak BW.	%Roof.
Stratix 10	1,178 µs	145 GOp/s	77 GB/s	52%

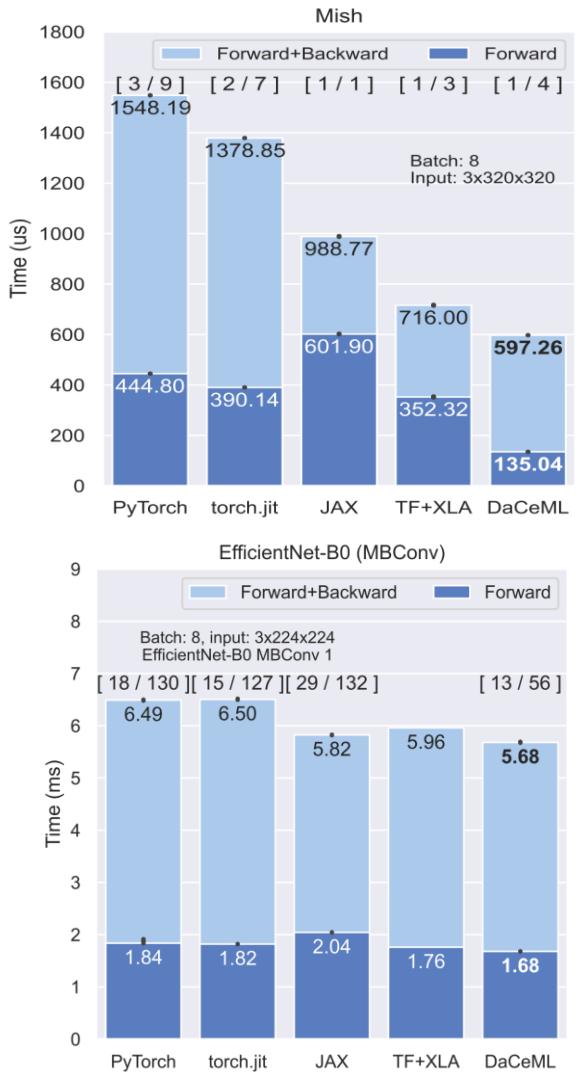
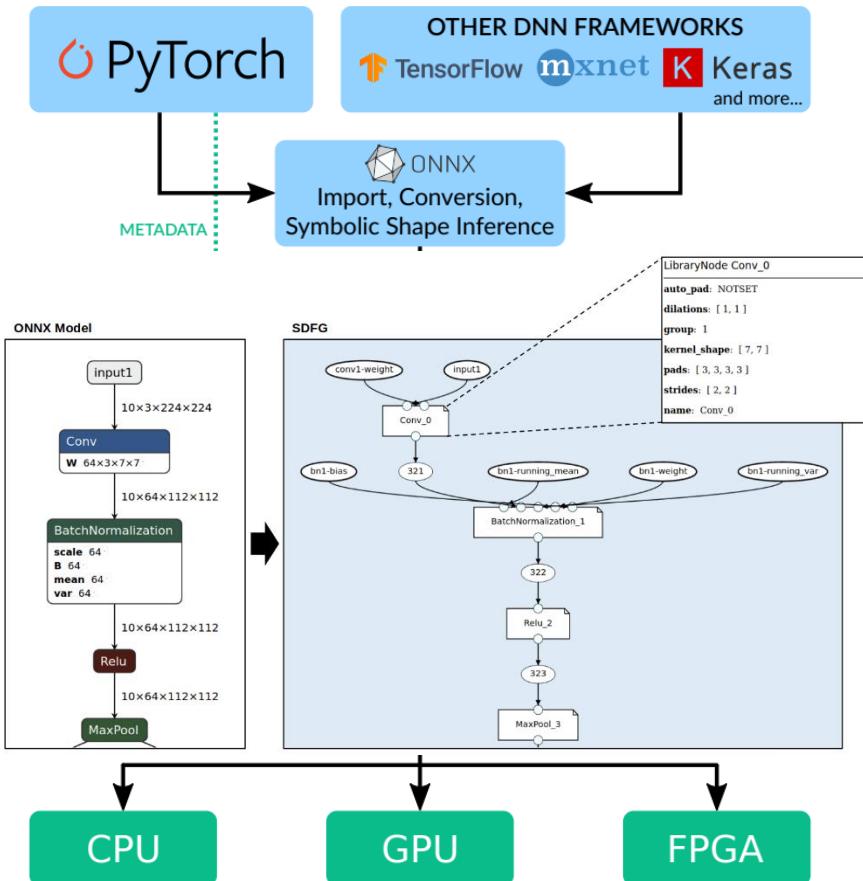
all memory bound

The Stratix 10 is held back by insufficient bandwidth.

FPGAs are good at deterministically **exploiting bandwidth**,  
but require a lot of **pipeline parallelism**.

# Mapping Deep Learning Codes to GPUs

**Key principle: minimize data movement and optimize data layout**



## Data Movement Is All You Need: A Case Study on Optimizing Transformers

Andrei Ivanov\*, Nikoli Dryden\*, Tal Ben-Nun, Shigang Li, Torsten Hoefer  
ETH Zürich  
firstname.lastname@inf.ethz.ch  
\* Equal contribution



challenges such as artificial general intelligence [27]. Thus, improving transformer performance has been in the focus of numerous research and industrial groups.

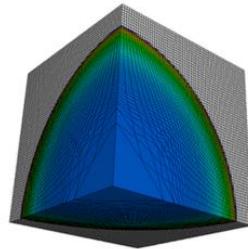
Significant attention has been given to optimizing transformers: local and fixed-window attention [28]–[32], more general structured sparsity [33], learned sparsity [34]–[36], and other algorithmic techniques [19], [37] improve the performance of transformers. Major hardware efforts, such as Tensor Cores and TPUs [38] have accelerated tensor operations like matrix-matrix multiplication (MMM), a core transformer operation. Despite this, existing implementations do not efficiently utilize GPUs. Even optimized implementations such as Megatron [18] report achieving only 30% of peak GPU flops.

We find that the **key bottleneck when training transformers is data movement**. Improvements in compute performance have reached the point that, due to Amdahl's Law and the acceleration of tensor contractions, training is now memory-bound.

Over a third (37%) of the runtime in a BERT training iteration is spent in memory-bound operators: While tensor contractions account for over 99% of the flop performed, they are only 61% of the runtime. By optimizing these, we show that the overhead of data movement can be reduced by up to 22.91%. Further, while MMM is highly tuned by BLAS libraries and hardware, we also find that **existing frameworks use suboptimal data layouts**. Using better layouts enables us to speed up MMM by up to 52%. Combining these insights requires moving beyond peephole-style optimizations and **globally optimizing data movement**, as selecting a single layout is insufficient. Overall, we achieve at least 1.30× performance improvements in training over general-purpose deep learning frameworks, and 1.08× over DeepSpeed [39], the state of the art manually-tuned implementation of BERT.

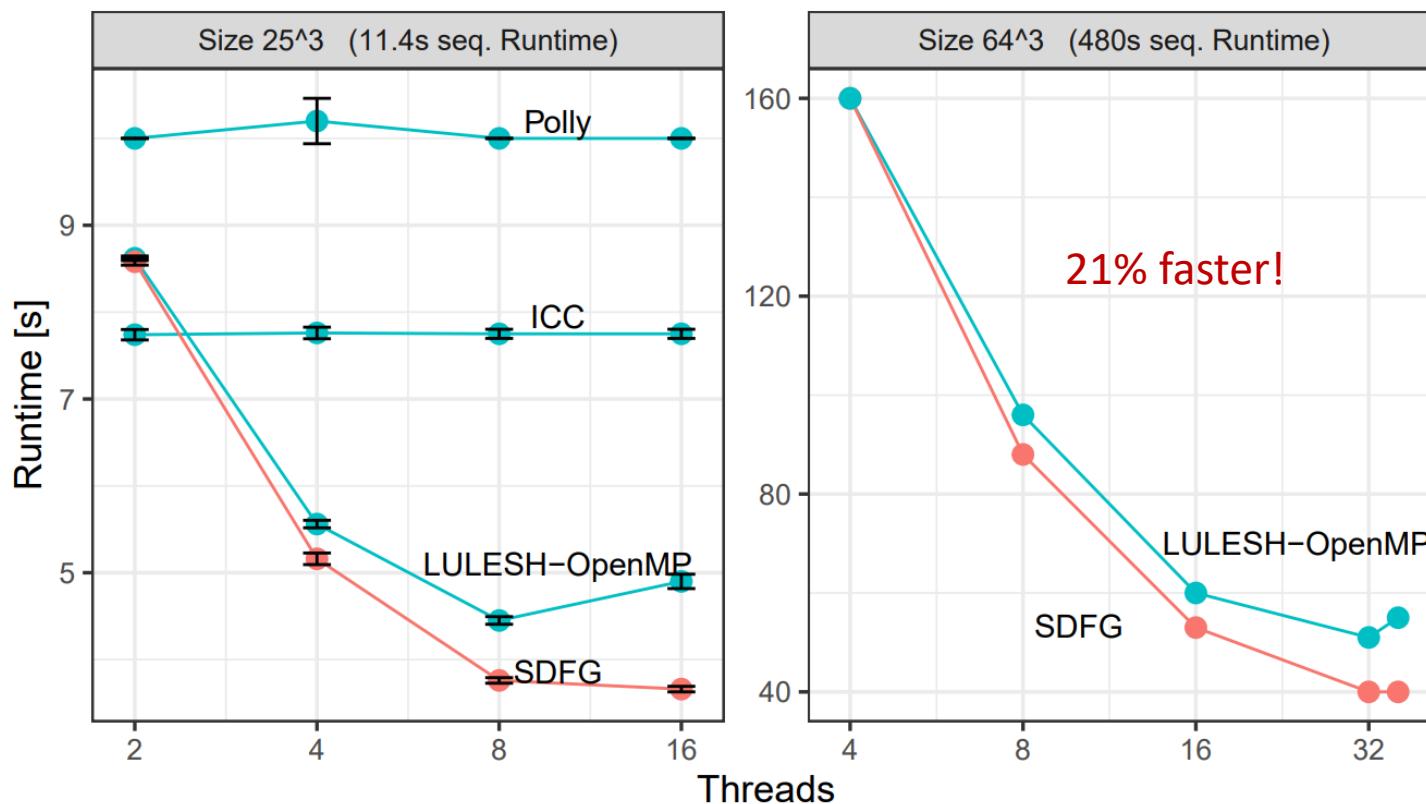
# Bonus: Mapping C/C++ (not just Python!) to CPUs

Key principle: find parallel regions using parametric dataflow



## LLNL's LULESH C++ benchmark

dual-socket 2×18 core Intel Xeon Gold 6154

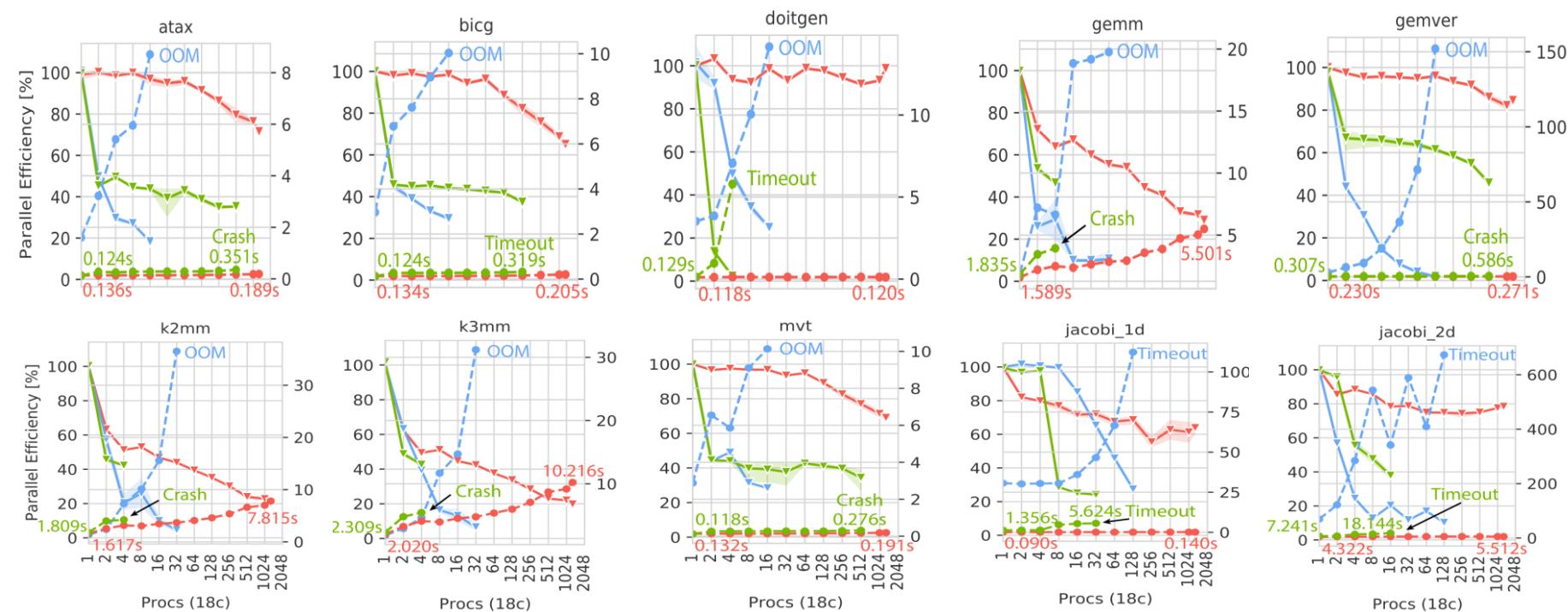
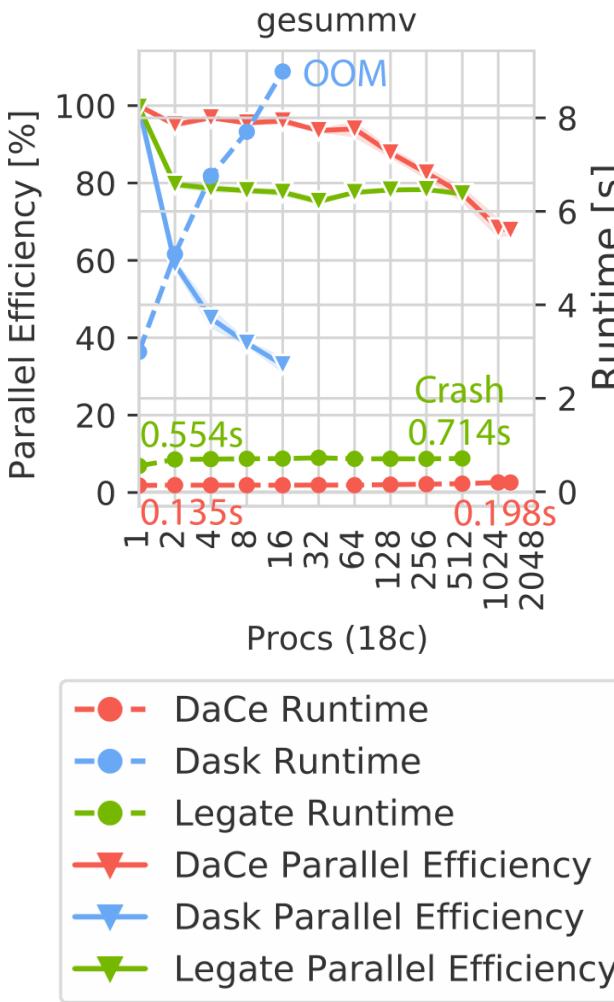


# Mapping NumPy to large-scale CPU clusters

**Key principle: communication-minimizing data mapping to nodes**



nv-legate/  
legate.numpy



**Using MPI library nodes  
(MPI-like interface)**

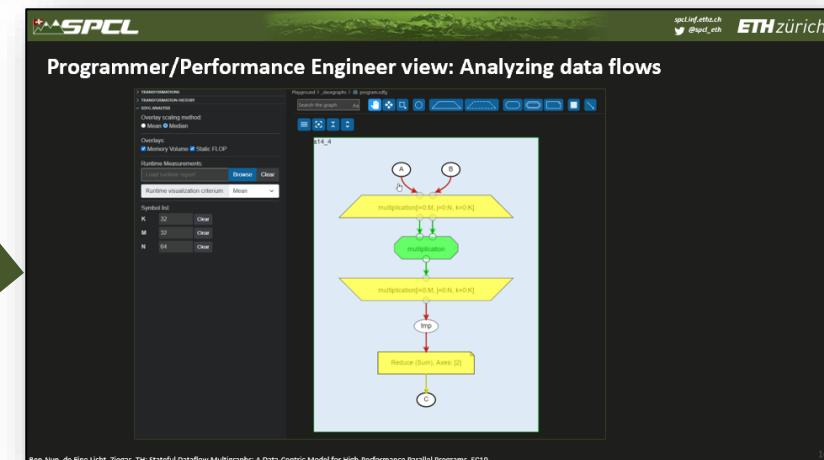
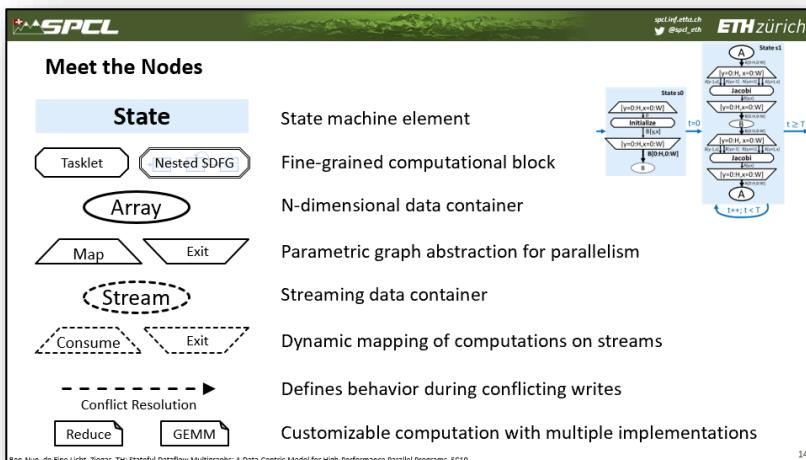
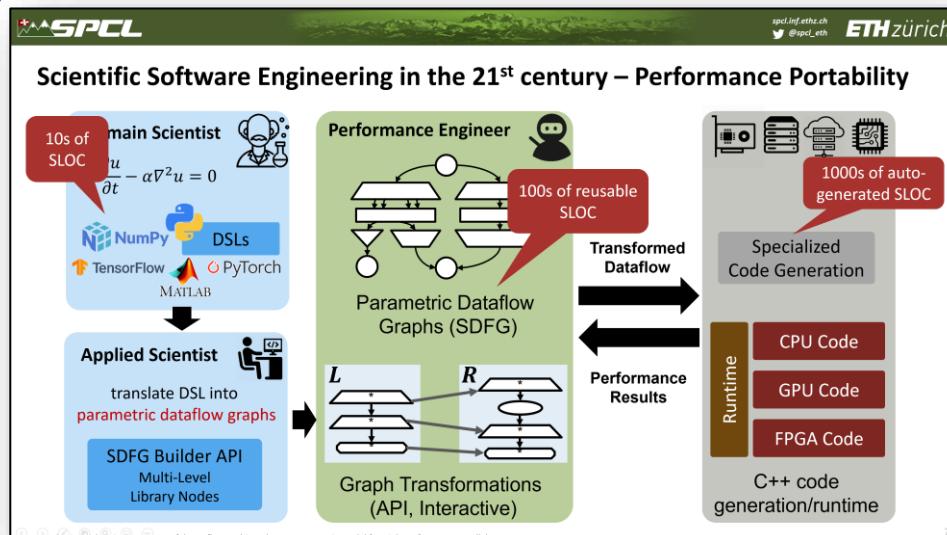
# Gordon Bell Prize 2019 on ORNL's Summit (Top-1 machine run with >21k GPUs)

- `pip install dace`
- **Gordon Bell Prize 2019**
  - Quantum Nano Transport simulation  
*Design of future micro-processors*
- **Now working on large-scale:**
  - Deep Learning (transformers)
  - Climate (COSMO, icon, fv3)
  - Green's functions solvers
  - ... your project?

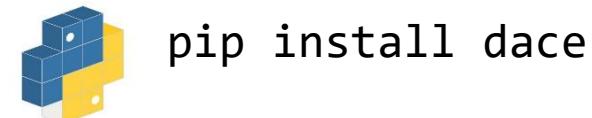


<http://spcl.inf.ethz.ch/DAPP>

# Overview and wrap-up



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



Open PhD and Postdoc positions:  
<https://spcl.inf.ethz.ch/Jobs/>



# SPCL is hiring PhD students and highly-qualified postdocs to reach new heights!

