

DaCe: Data Centric parallel programming for heterogeneous platforms

Tiziano De Matteis
Vrije Universiteit Amsterdam

Who am I?

Tiziano De Matteis

Assistant Professor at VU Amsterdam (joined March 2023)



Research Interests

I am interested in High-Performance Computing, Energy Efficiency, Systems.

Current focus: abstractions, methods, and systems to overcome the end of Moore's law.

Previously:

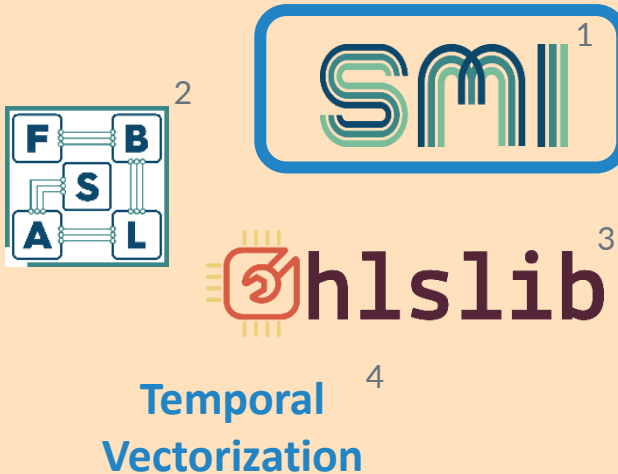
- ▷ PhD University of Pisa (on Parallel Streaming Analytics)
- ▷ PostDoc ETH Zurich (on FPGA for HPC)

Contributions to the FPGA-HPC community

Focus on methods and tools to productively program (with HLS) FPGAs for High-Performance Computing

Libraries/Tools

To increase programming productivity



Programming Model

To provide performance portability in Heter. Architectures (not only FPGA!)



Applications

To show the potentials of reconfigurable hardware in applications



¹ T. De Matteis, J. de Fine Licht, J. Beránek, T. Hoefer. "Streaming Message Interface: High-Performance Distributed Memory Programming on Reconfigurable Hardware". SC'19

² T. De Matteis, J. de Fine Licht, T. Hoefer. "FBLAS: Streaming Linear Algebra on FPGA". SC'20

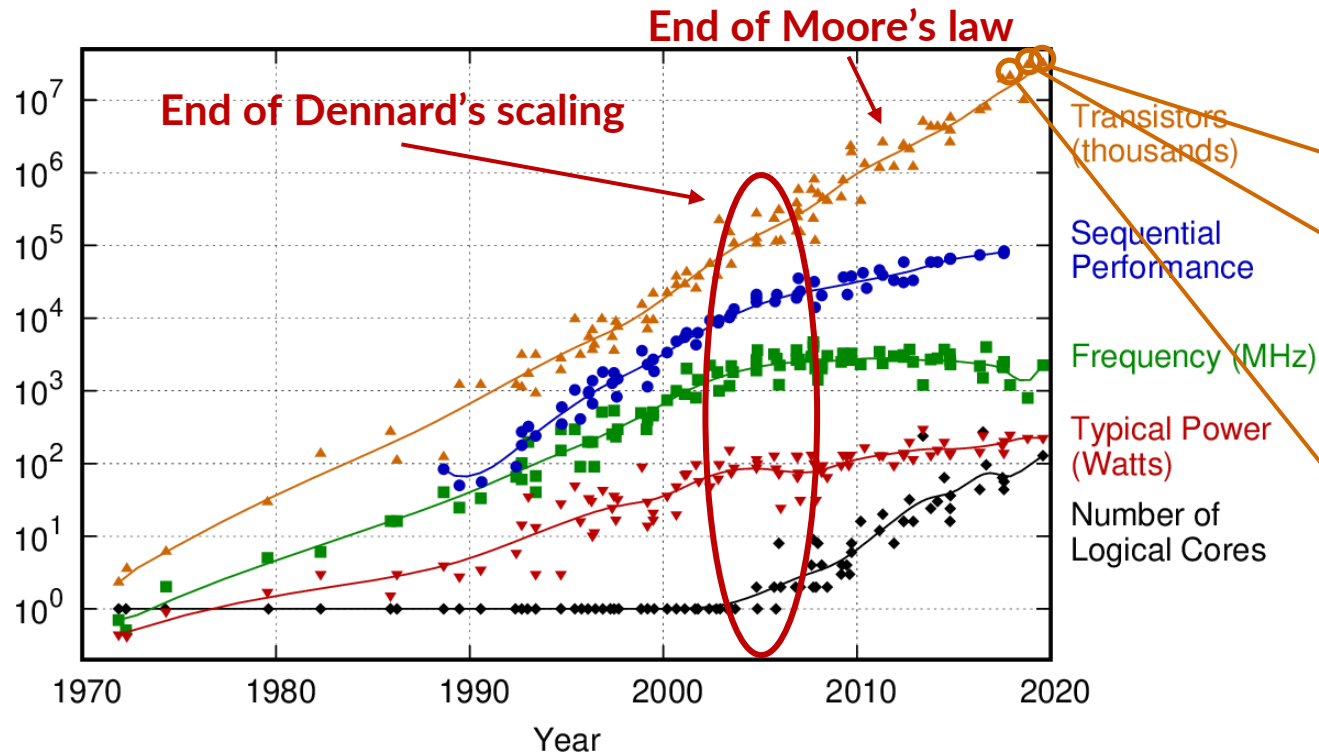
³ J. de Fine Licht, T. Hoefer. "hlslib: Software Engineering for Hardware Design". H2RC'19

⁴ C. Johnsen, T. De Matteis, T. Ben-Nun, J. de Fine Licht, T. Hoefer de Fine Licht, T. Hoefer. "Temporal Vectorization: A Compiler Approach to Automatic Multi-Pumping". ICCAD'22

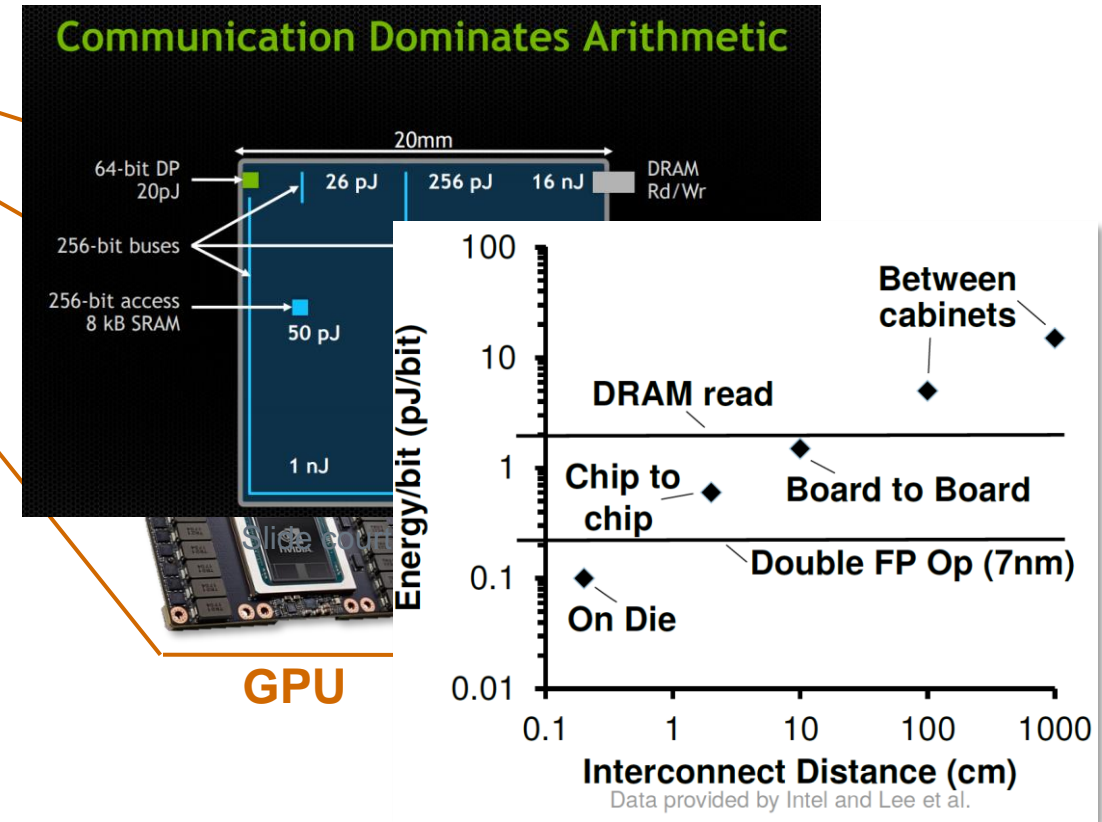
⁵ A. N. Ziogas, T. Schneider, T. Ben-Nun, A. Calotoiu, T. De Matteis, J. de Fine Licht, L. Lavarini, T. Hoefer. "Productivity, Portability, Performance: Data-Centric Python". SC '21.

⁶ J. de Fine Licht, A. Kuster, T. De Matteis, T. Ben-Nun, D. Hofer, T. Hoefer. "StencilFlow: Mapping Large Stencil Programs to Distributed Spatial Computing Systems". CGO '21

Changes in the Computing Landscape



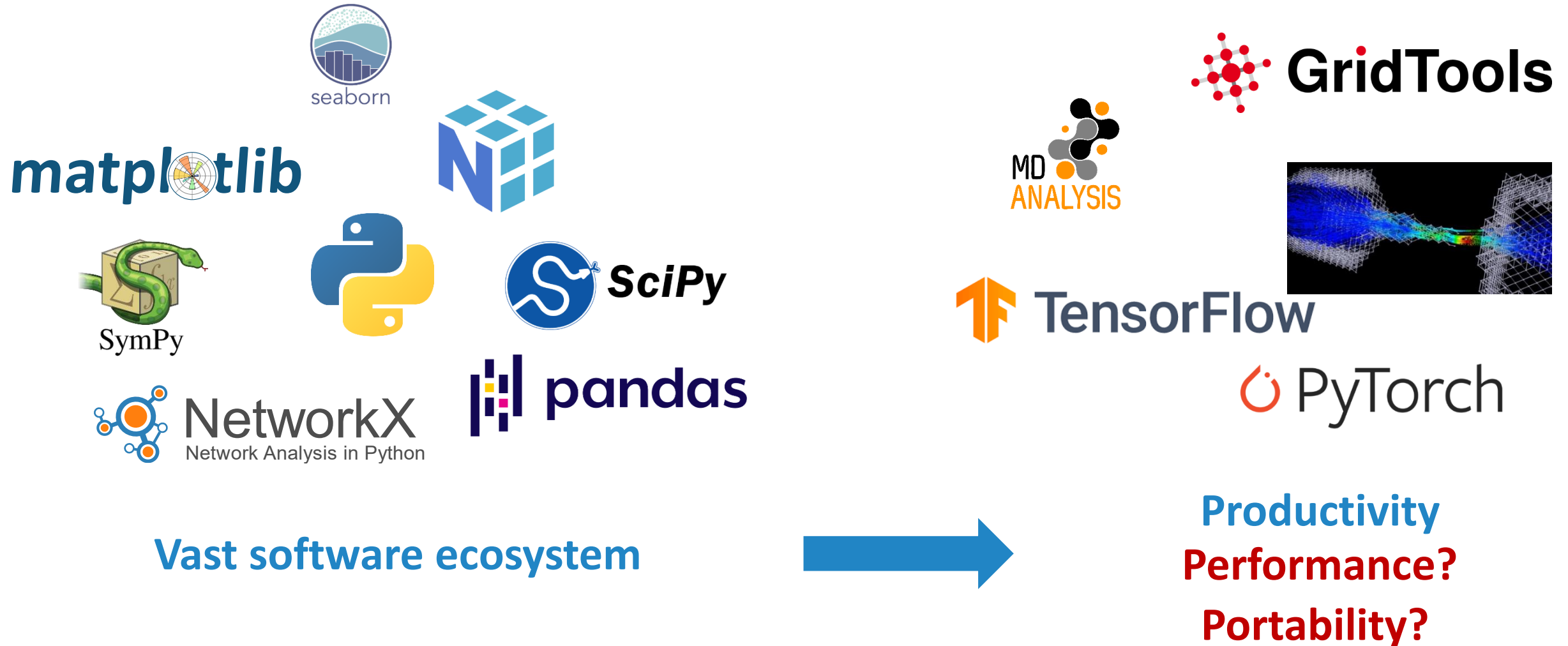
Original data up to the year 2010 collected and plotted by M. Horowitz, F. Labonte, O. Shacham, K. Olukotun, L. Hammond, and C. Batten
New plot and data collected for 2010-2017 by K. Rupp



The holy grail of Performance
Portability!

High-Performance Optimization -> Data
Movement Reduction

Changes in the Software Landscape – The Scientific Languages



DaCe Overview

Domain Scientist

Problem Formulation

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

Python

DSLs

PyTorch

C

...

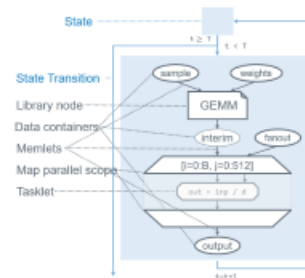
Scientific Frontend



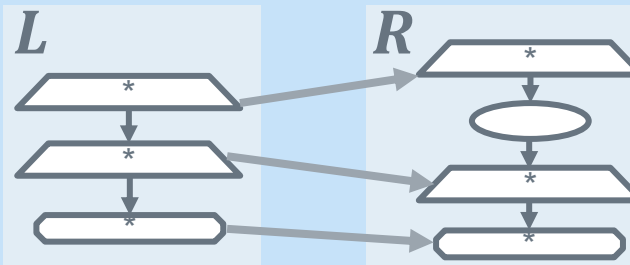
Performance Engineer

Data-Centric Intermediate Representation

1. Separate data containers from computation
2. Coarsening: multi-level view of data movement
3. Data movement as a first-class citizen (dependencies → program order)
4. Control dependencies only when dataflow is not implied



Data-Centric Intermediate Representation (SDFG)



Graph Transformations



Transformed
Dataflow



Performance
Results



System

Hardware
Information

Compiler

Runtime

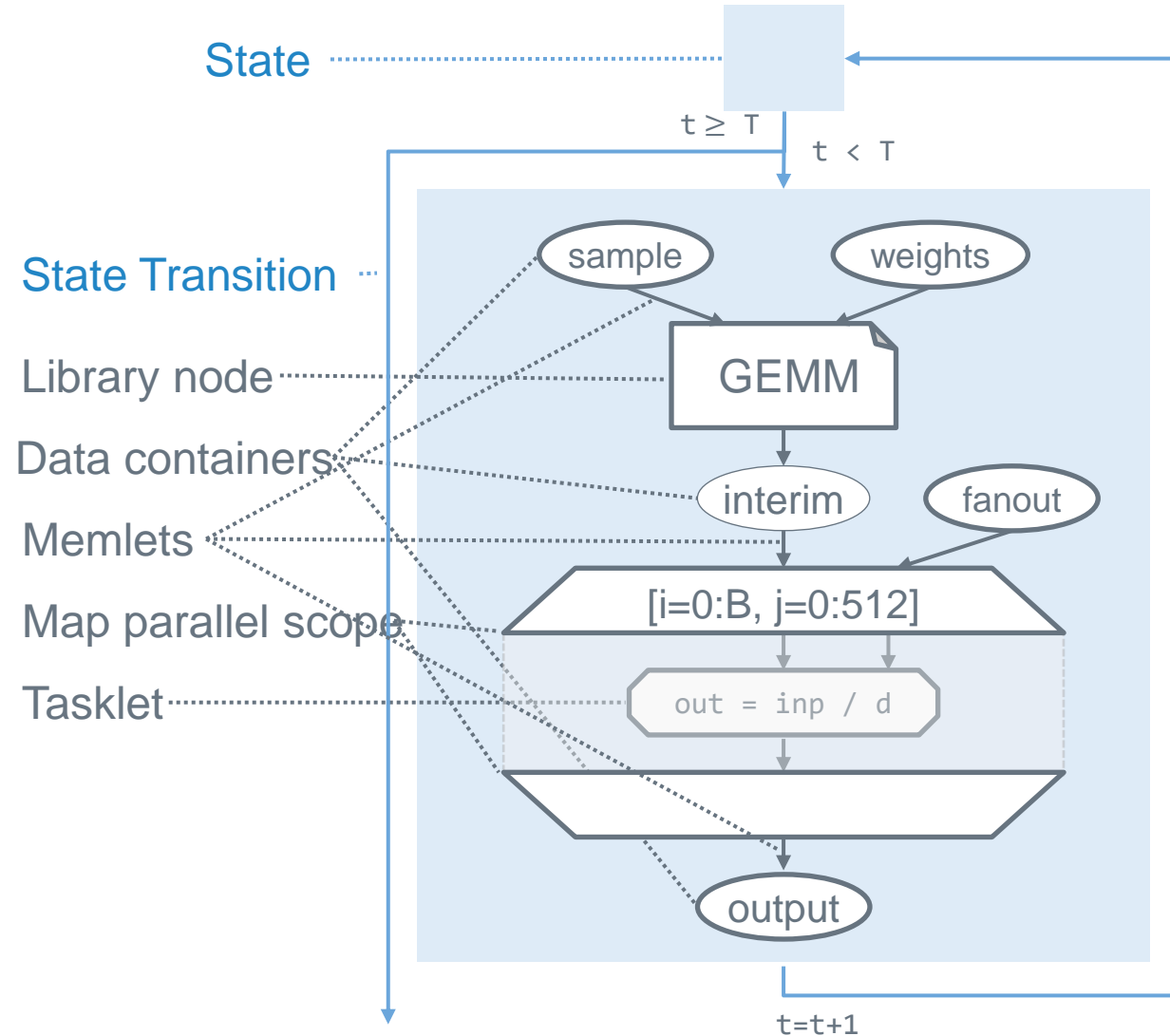
CPU Binary

GPU Binary

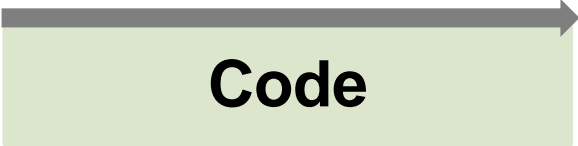
FPGA Modules

Data-Centric Intermediate Representation

1. Separate data containers from computation
2. Coarsening: multi-level view of data movement
3. Data movement as a first-class citizen (dependencies \rightarrow program order)
4. Control dependencies *only when dataflow is not implied*



From source code to SDFG



Code

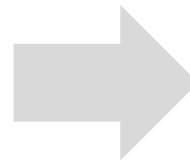
```
glob_b = ...  
class ClassA:  
    def __init__(self, arr):  
        self.q = arr  
  
    @dace.method  
    def __call__(self, a):  
        return a * self.q + glob_b
```



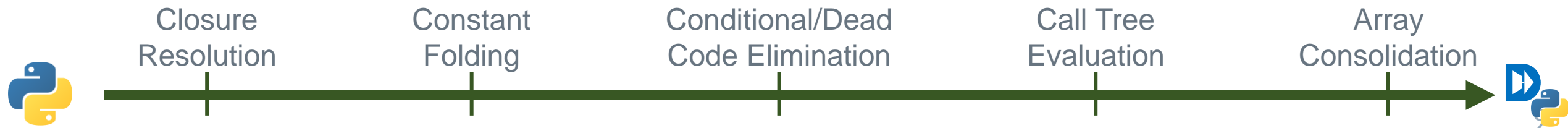
From source code to SDFG



```
glob_b = ...  
class ClassA:  
    def __init__(self, arr):  
        self.q = arr  
  
    @dace.method  
    def __call__(self, a):  
        return a * self.q + glob_b
```



```
@dace.program  
def ClassA__call__(a, __g_self_q, __g_glob_b):  
    return a * __g_self_q + __g_glob_b
```

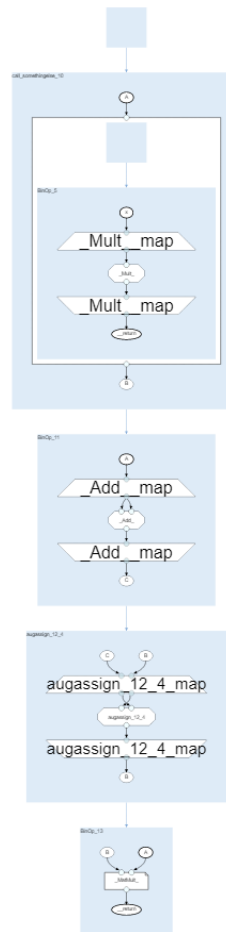


From source code to SDFG

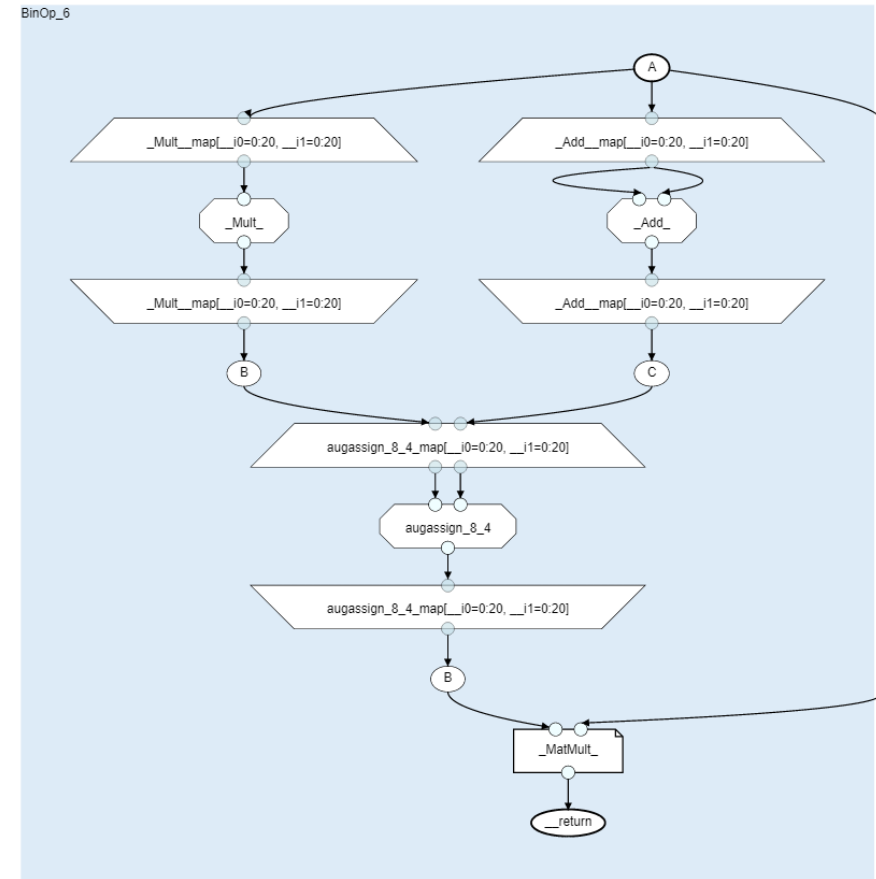
Code

Preprocess

Simplify



Control flow
coarsening passes



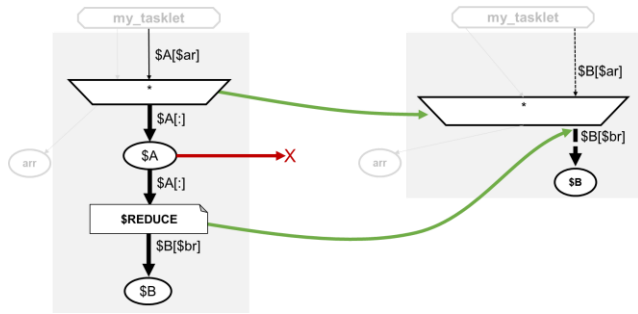
From source code to SDFG

Code

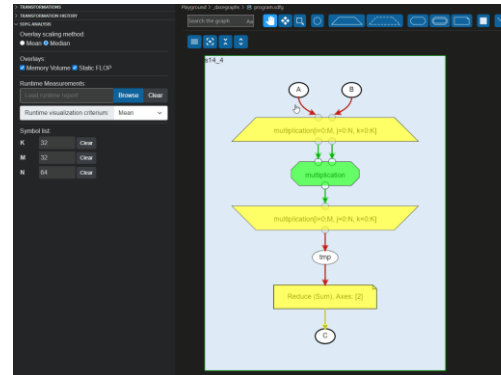
Preprocess

Simplify

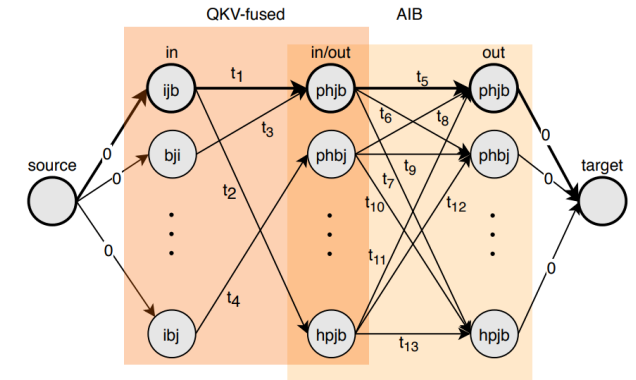
Optimize



Graph Rewriting Transformations



Interactive Transformation and Instrumentation



Local and Global Tuning Interface

```
class ClassA:
    def __init__(self, arr):
        self.q = arr

    @dace.method(auto_optimize=True, device=dace.DeviceType.GPU)
    def __call__(self, a):
        return a * self.q + glob_b
```



Automatic Optimization Heuristics

Visual Studio Code Integration

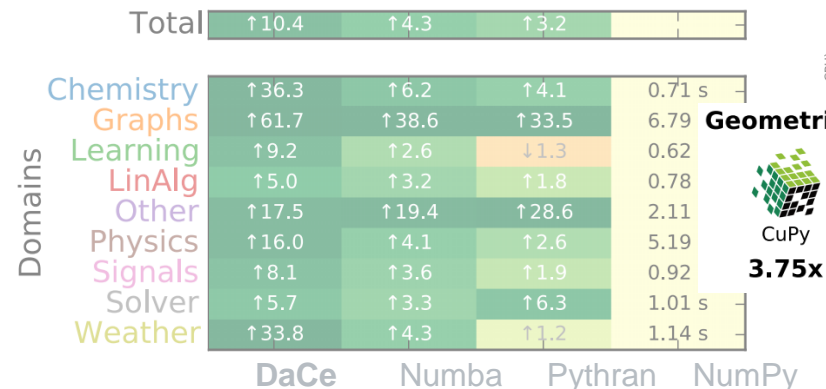
The screenshot displays the Visual Studio Code interface with the following components:

- Left Sidebar:** Contains two panels. The top panel, 'TRANSFORMATIONS', lists various optimization techniques under 'Selection' and 'Viewport' categories, including 'FPGATransformState', 'MapTiling', 'StripMining', etc. The bottom panel, 'TRANSFORMATION HISTORY', shows 'No previously applied transformations'.
- Central Editor:** Displays a Python file named 'gemm.py'. The code defines a Dace program for a matrix multiplication (gemm) using Dace symbols and a nested loop structure. It includes comments and uses Dace's `@dace.program` decorator.
- Right Sidebar:** Shows the 'program.sdfg' graph. The top part is a visual graph with nodes for inputs A and B, multiplication operations, a temporary variable 'tmp', and a 'Reduce (Sum)' operation. The bottom part is a configuration panel for the 'SDFG gemm' with fields for 'arg_names', 'constants_prop', 'exit_code', 'global_code', 'init_code', 'instrument', 'openmp_sections', and 'symbols'.

DaCe Performance



<https://github.com/spcl/npbench>



Geometric Mean of Improvements over Other Frameworks



3.75x



2.47x



10.6x



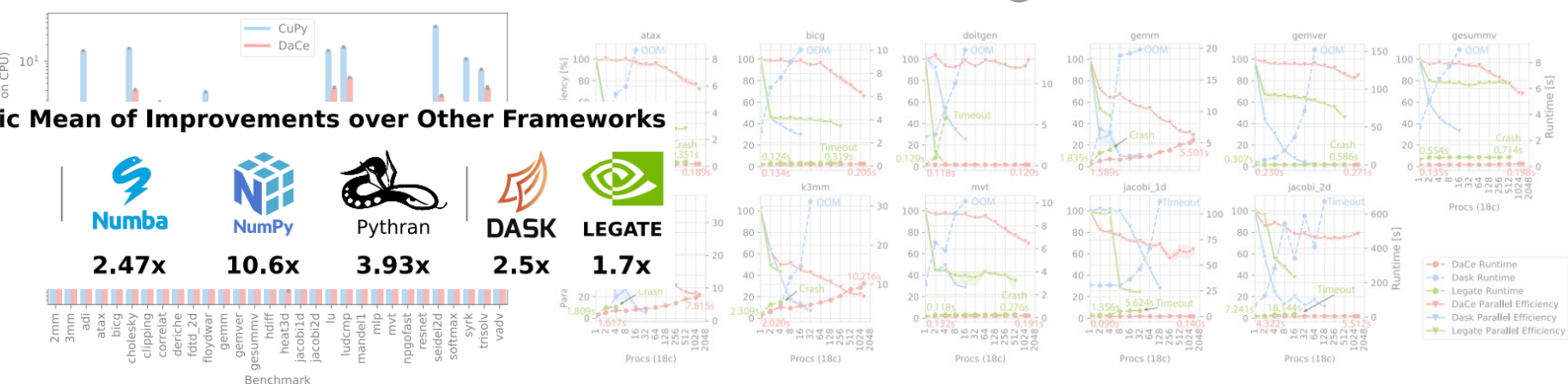
3.93x



2.5x



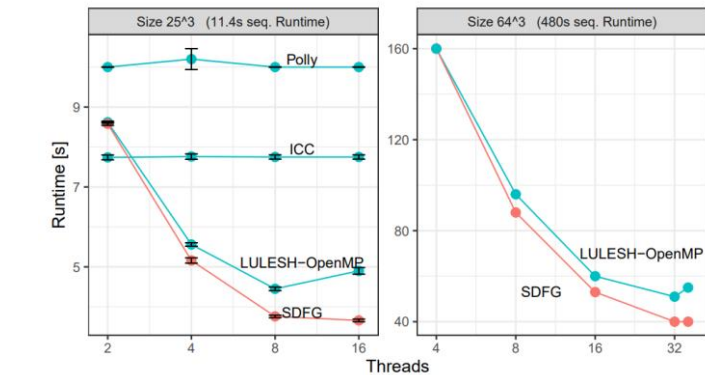
1.7x



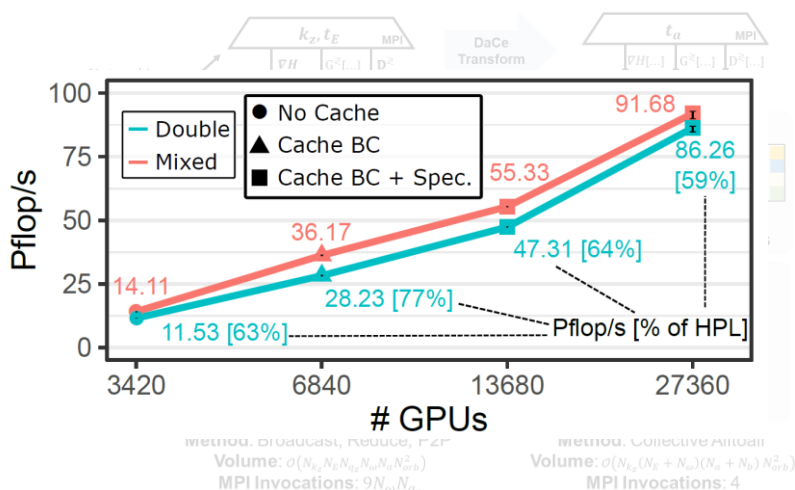
NumPyBench

NumPyBench GPU

NumPyBench Distributed



Unstructured Hydrodynamics (LULESH)



Quantum Transport Simulation

DaCe and FPGA

Let's use DaCe to compute $y = A^T Ax$ (ATAX)

```
import dace
```

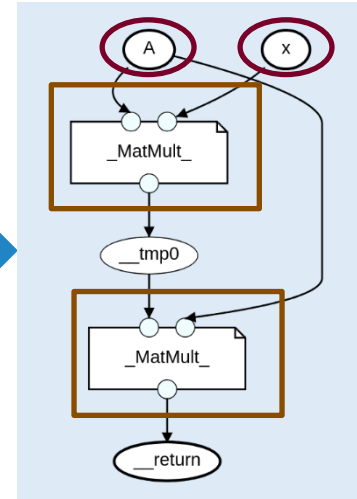
```
M, N = 24, 24
```

```
@dace.program
```

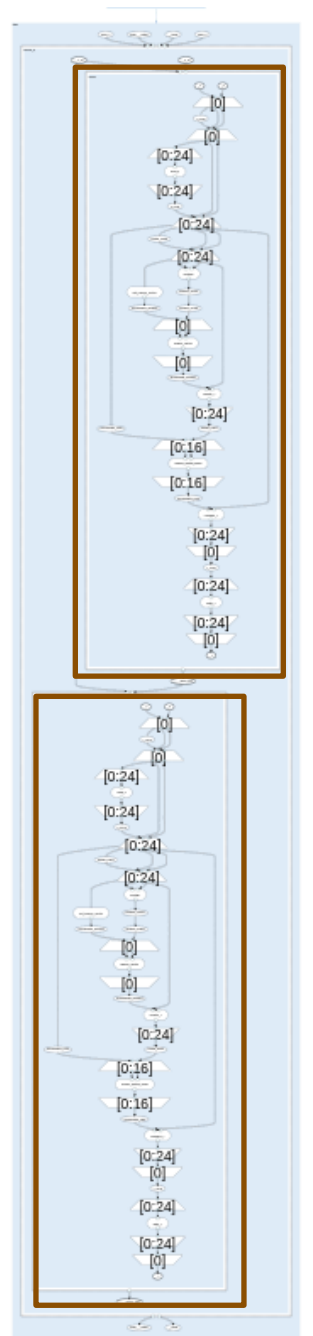
```
def atax(A: dace.float32[M, N],  
        x: dace.float32[N]):
```

```
    return (A @ x)@A
```

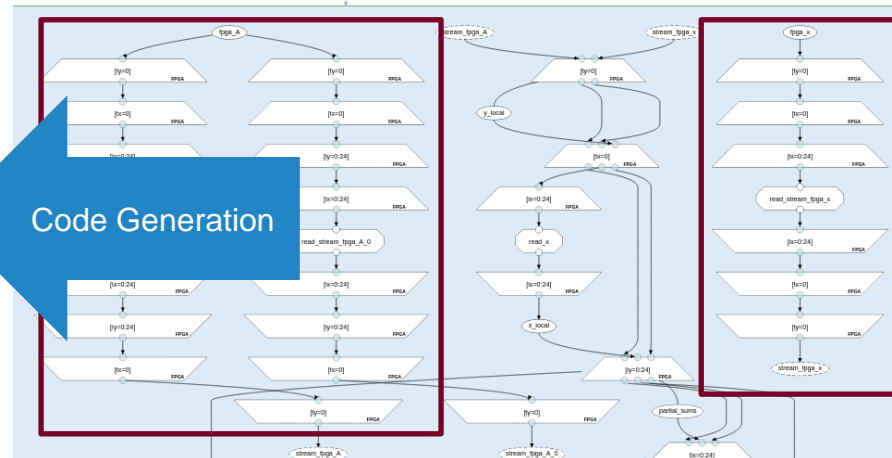
SDFG Parsing and
Simplification



Specialization



Code Generation



Streaming
Memory

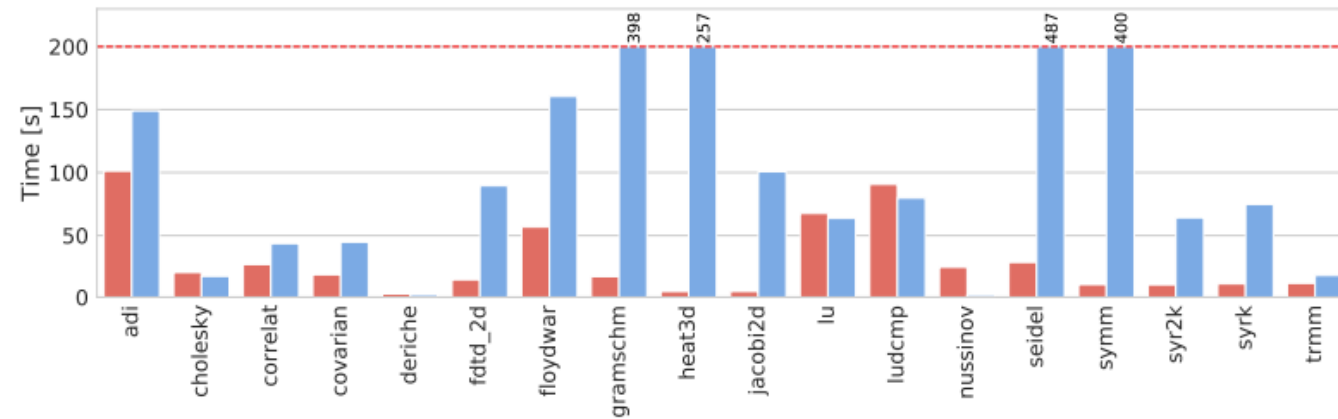
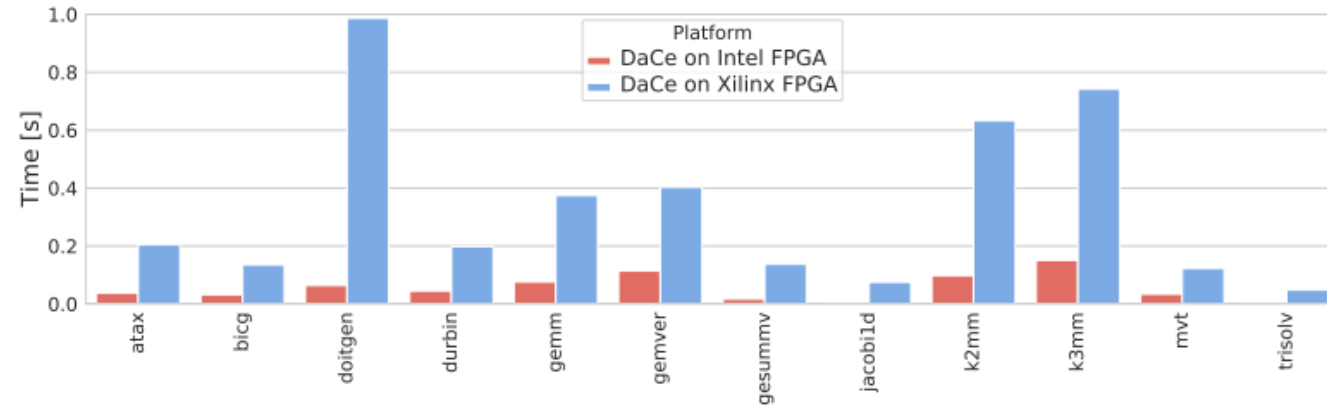
Apply other optimizations

DaCe and FPGA

Let's use DaCe to compute $y = A^T Ax$ (ATAx)

```
import dace

M, N = 24, 24
@dace.program
def atax(A: dace.float32[M, N],
        x: dace.float32[N],
        y: dace.float32[M]):
    return (A @ x) @ A
```



Polybench kernels code generated from Numpy for both Xilinx and Intel



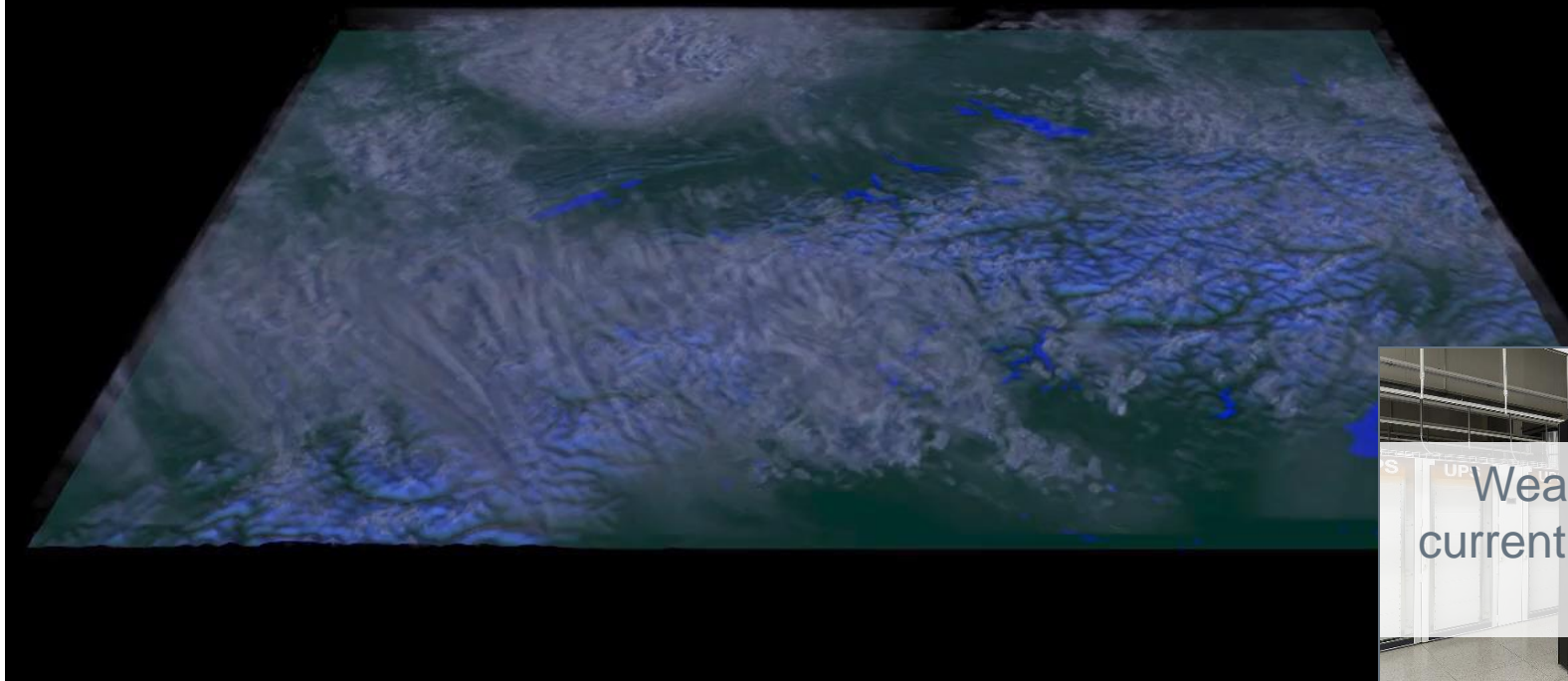
Weather simulation at MeteoSwiss

[1]

COSMO 1.1 km

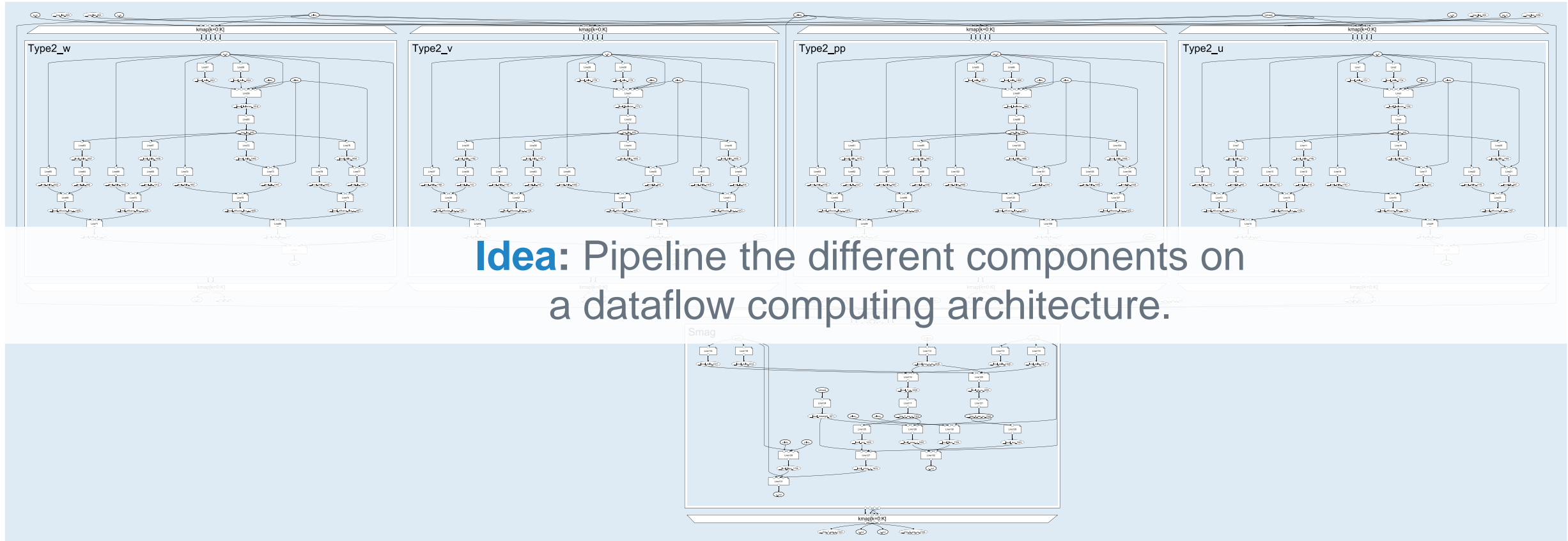
2018-05-29 00:00 UTC+2

ETH zürich

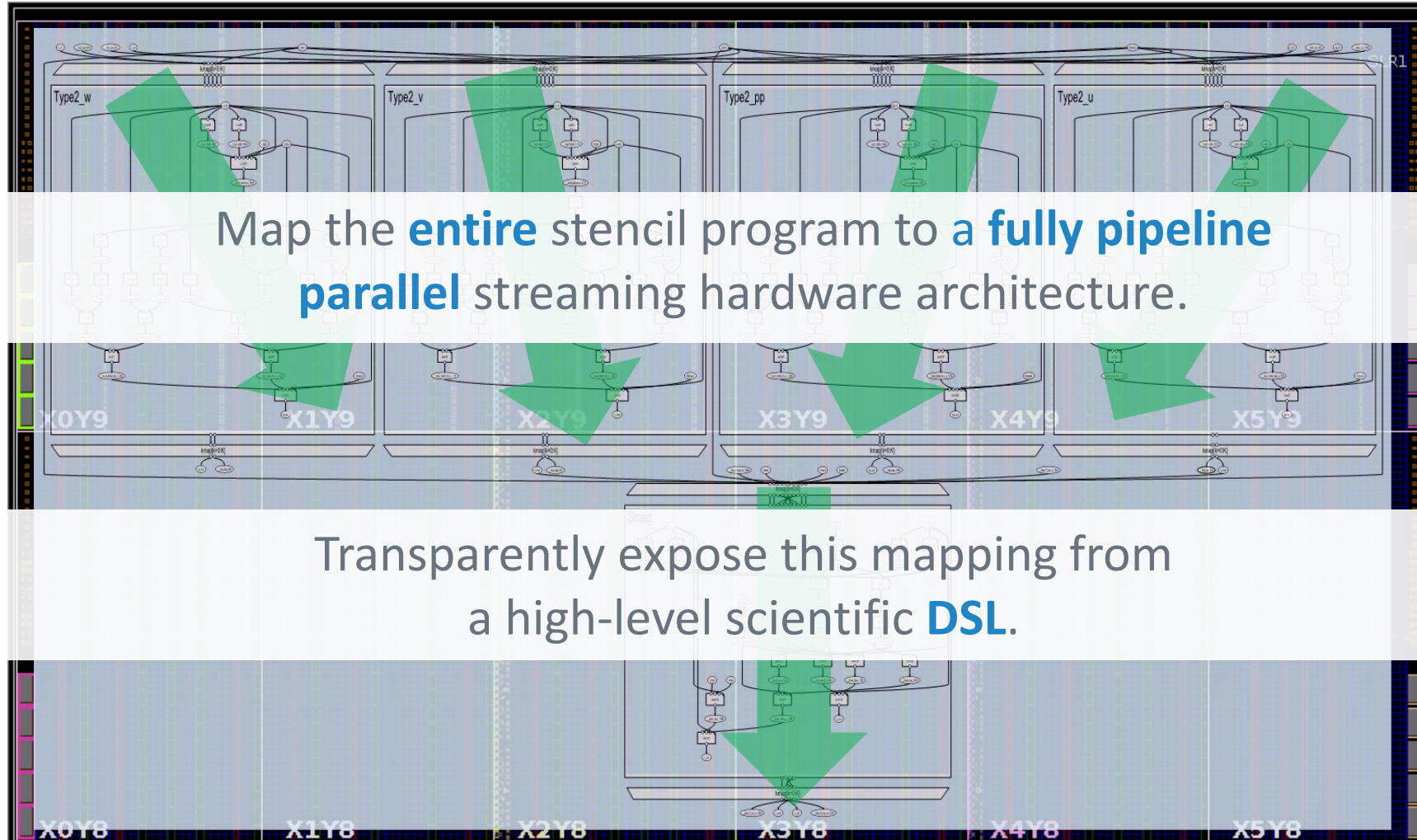


[1] Institute for Atmospheric and Climate Science and Computer Graphics Laboratory, ETH Zürich [<https://vimeo.com/389292423>]
[2] Swiss National Supercomputing Center (CSCS) [<https://www.cscs.ch/computers/arolla-tsa-meteoswiss>]

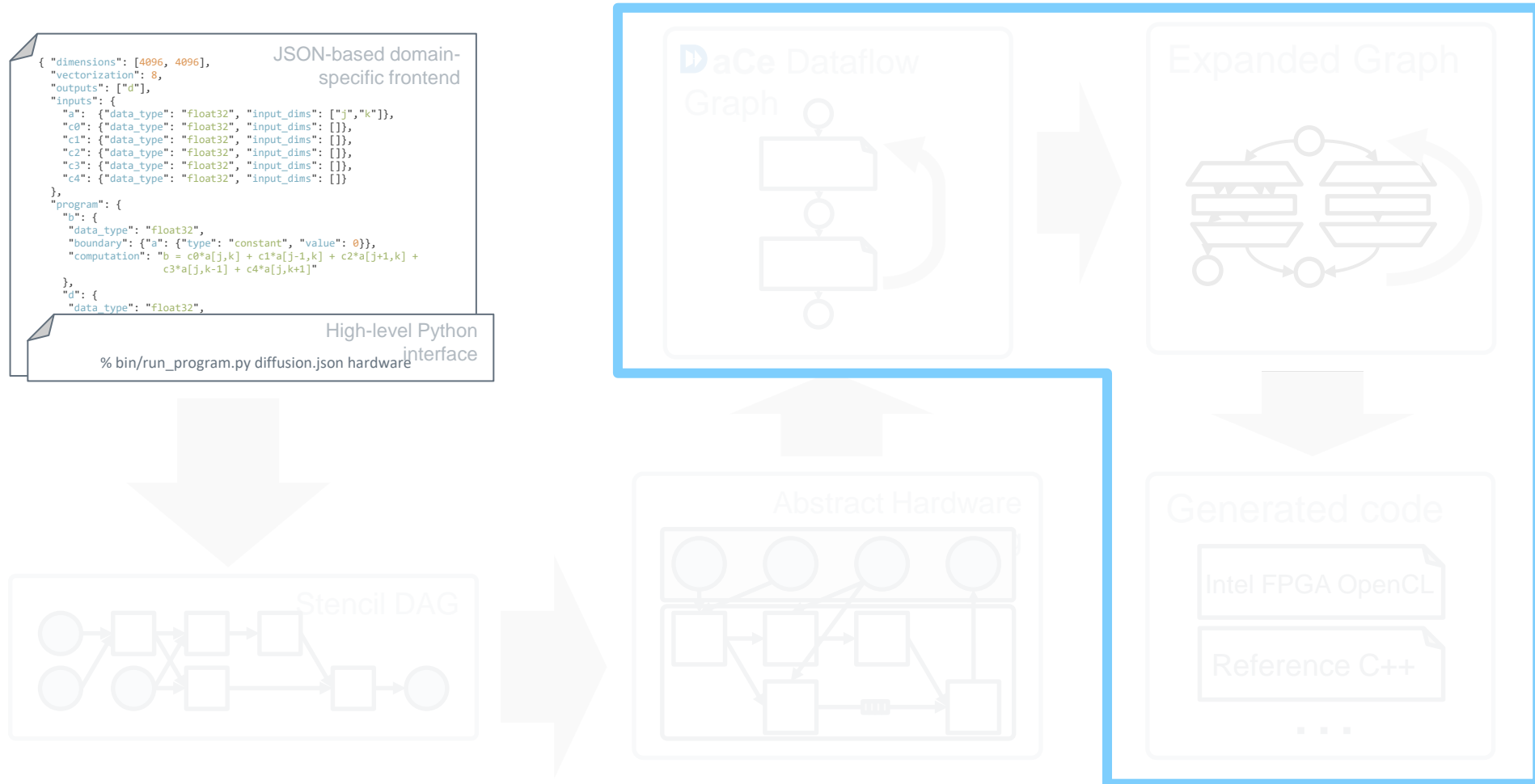
Weather programs



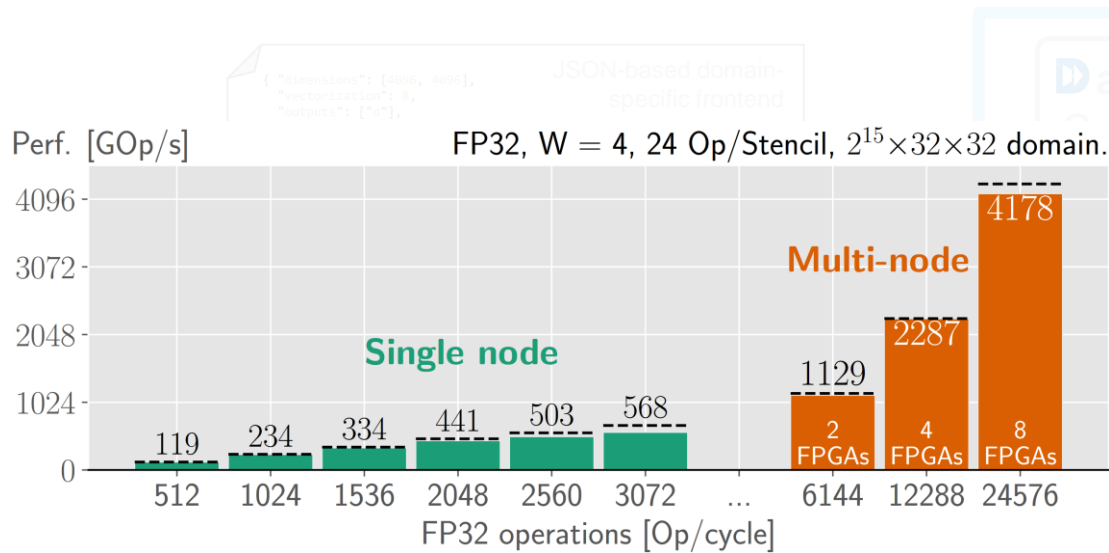
StencilFlow – FPGAs for Climate Modeling




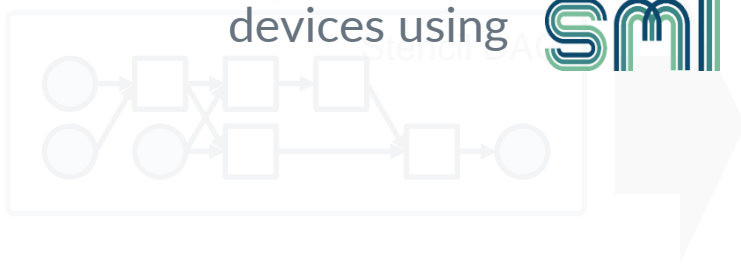
Stencilflow



Stencilflow



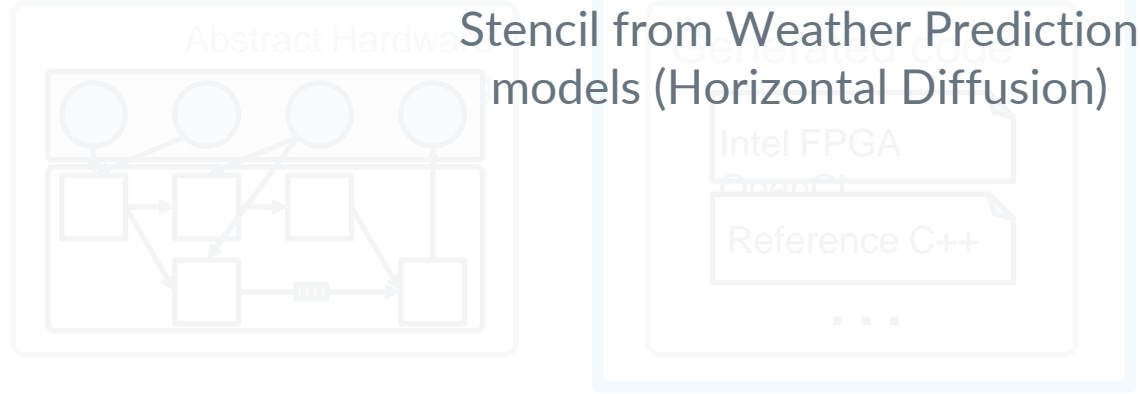
Iterative stencils (e.g., Jacobi, Diffusion, ...) on multiple FPGA devices using 



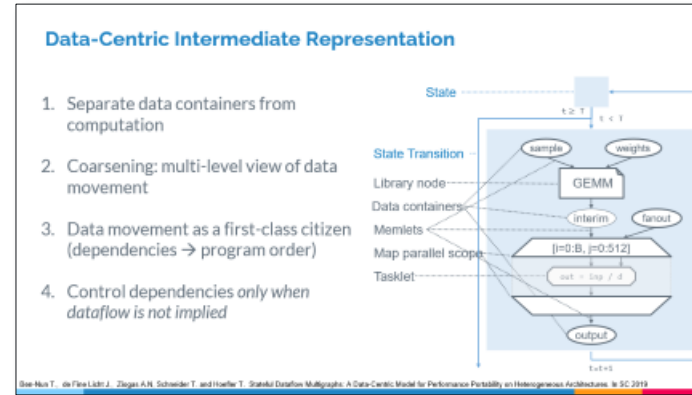
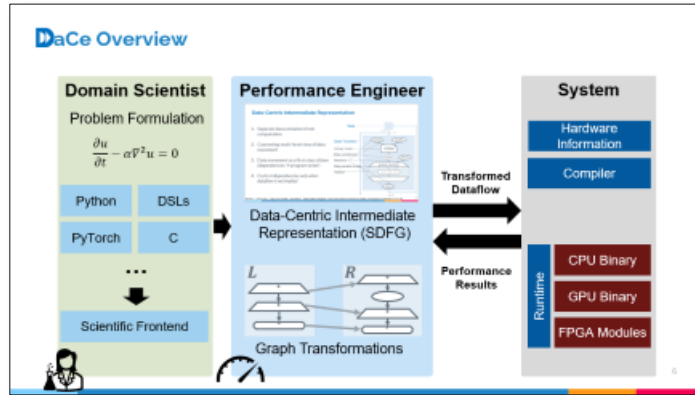
	Runtime	Performance	Peak BW.	%Roof.
Stratix 10	1,178 μ s	145 GOp/s	77 GB/s	52%
Stratix 10*	332 μ s	513 GOp/s	∞ GB/s	—
Xeon 12C	5,270 μ s	32 GOp/s	68 GB/s	13%
P100	810 μ s	210 GOp/s	732 GB/s	8%
V100	201 μ s	849 GOp/s	900 GB/s	26%

*Without memory bandwidth constraints.

Stencil from Weather Prediction models (Horizontal Diffusion)



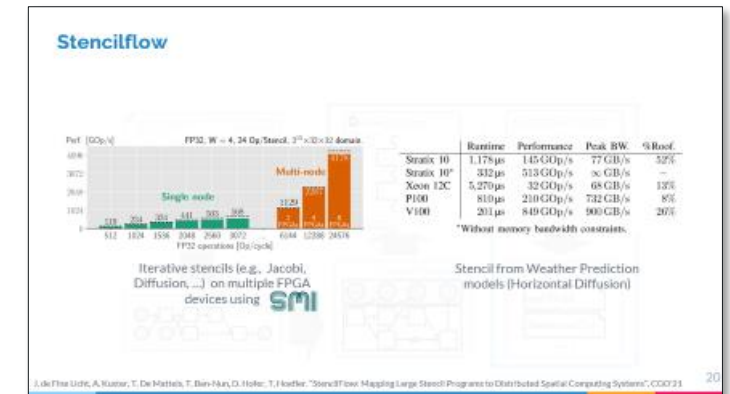
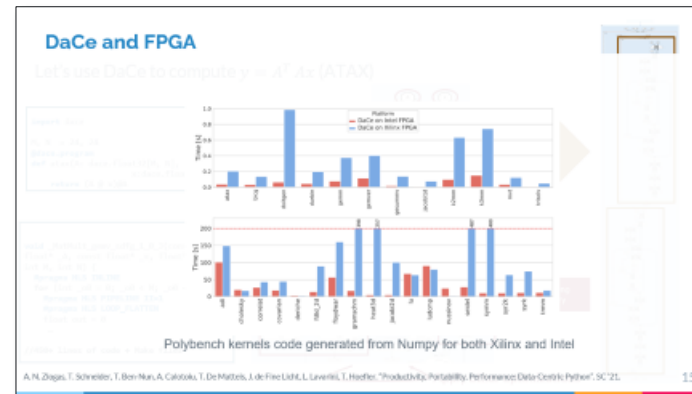
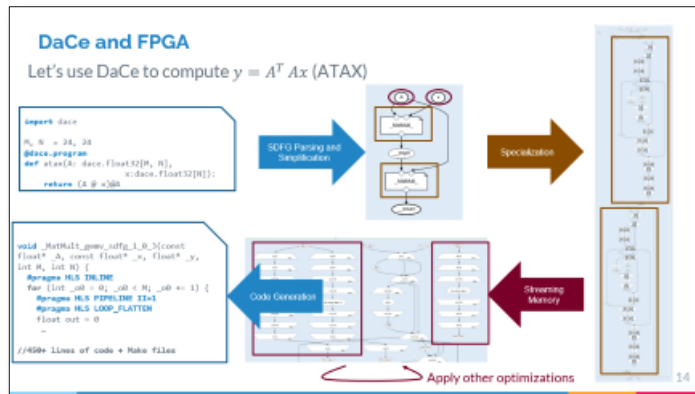
Conclusion



DaCe is in active development.
If you want to know more/contribute:



github.com/spcl/dace



t.de.matteis@vu.nl