

# StencilFlow: Mapping Large Stencil Programs to Distributed Spatial Computing Systems

Johannes de Fine Licht\*, Andreas Kuster\*, Tiziano De Matteis\*, Tal Ben-Nun\*, Dominic Hofer\*<sup>†</sup>, Torsten Hoefler\*

\*Department of Computer Science, ETH Zurich, Switzerland <sup>†</sup>MeteoSwiss, Switzerland

{definelj, kustera, tdematt, tbennun, dohofer, htor}@ethz.ch

Abstract-Spatial computing devices have been shown to significantly accelerate stencil computations, but have so far relied on unrolling the iterative dimension of a single stencil operation to increase temporal locality. This work considers the general case of mapping directed acyclic graphs of heterogeneous stencil computations to spatial computing systems, assuming large input programs without an iterative component. StencilFlow maximizes temporal locality and ensures deadlock freedom in this setting, providing end-to-end analysis and mapping from a high-level program description to distributed hardware. We evaluate our generated architectures on a Stratix 10 FPGA testbed, yielding 1.31 TOp/s and 4.18 TOp/s on single-device and multi-device, respectively, demonstrating the highest performance recorded for stencil programs on FPGAs to date. We then leverage the framework to study a complex stencil program from a production weather simulation application. Our work enables productively targeting distributed spatial computing systems with large stencil programs, and offers insight into architecture characteristics required for their efficient execution in practice.

### I. INTRODUCTION

Spatial architectures such as FPGAs, Intel's *Configurable Spatial Accelerator* [1], Xilinx' AI Engines [2], and the Cerebras deep neural network accelerator [3], are characterized by a large number of small processing units connected by a configurable network. These systems sacrifice generality and traditional coherence across hierarchical memory subsystems to achieve higher transistor efficiency than load/store architectures (a.k.a. von Neumann architectures), which is essential to continue scaling in the post-Dennard/Moore era.

Common to spatial architectures is their amenability to be programmed with dataflow abstractions, as this throws away notions of implicit accesses to off-chip resources and communication between parallel processing units, in favor of explicitly programmable off-chip and on-chip data movement. In this paradigm, computations are laid out spatially on the device, rather than existing as a temporal instruction stream, directly exposing notions of data locality. The simplest example of this is a pipeline, where each stage synchronously feeds the next. Systolic arrays add more complexity by extending this to a sequence of synchronous pipelines that communicate in a messaging fashion. Generally, directed acyclic graphs (DAGs) of pipelines allow arbitrary dataflow, where each node can be attached to multiple producers and consumers.

The temporal locality in iterative or dependent stencil computations is challenging to exploit on load/store architectures, as they require complex tiling schemes [4] and selective fusion



Fig. 1: Overview of the StencilFlow end-to-end system.

of code segments [5], [6]. In contrast, exploiting this reuse via dataflow is intuitive, as consecutive stages can be pipelined and synchronized via their fine-grained dependencies [7]. Implementations of stencils achieving high performance on reconfigurable hardware often assume idealized iterative stencils, as this enables temporal blocking of consecutive timesteps [8], [9], which maps naturally to pipelined architectures.

In this work, we consider the challenging case of arbitrary stencil DAGs, motivated by their existence in numerical climate and weather prediction, where each node is a (potentially complex) stencil operation reading from one or more input memories, and writing its output to one or more consumers. As a motivating case study, we target an application from the Consortium for Small-scale Modeling (COSMO). The consortium consists of eight national weather services which aim to develop, improve and maintain a non-hydrostatic local area atmospheric model. The COSMO model is used for both operational [10], [11] and research [12], [13] applications by the members of the consortium and many universities worldwide. The stencils used in these simulations are dominated by series of heterogeneous stencil computations. Unlike the uniform codes often evaluated in high-performance computing research, these programs run many different stencil operations on many different inputs of varying dimensionality, and exhibit complex dependency patterns between them.

We present a full-stack solution, from a high-level stencil DSL to low-level spatial program definitions, that are code generated for hardware execution, summarized in Fig. 1. We introduce a method that maps stencil programs to spatial architectures by using dataflow principles to form compositions that are deadlock free and maximize the number of active pipelines, based on an analysis of iteration patterns and the computational source code. *Fully code-generated* architectures emitted by StencilFlow evaluated on an FPGA testbed reach 1.32 TOp/s and 4.18 TOp/s in single-device and multi-device

experiments, respectively, which to the best of our knowledge is the *highest performance recorded for stencil programs executed on FPGA hardware to date*. The full source code is available on GitHub<sup>1</sup>, exposing productive high-level Python interfaces, while compiling to highly efficient hardware through the codegenerating backend.

Clear separation of concerns at multiple levels of the stack is a key concept in our approach. An input program is formulated as a high-level DSL, constraining the program to an analyzable and optimizable form. Input programs are first optimized on a domain-specific level, where we can perform specialized transformations, such as fusing consecutive stencil nodes. Then, programs are lowered to a dataflow representation represented in DaCe [14], where we can control and optimize for data movement. The dataflow representation is then specialized for the targeted architecture, and finally code-generated to be compiled and synthesized for hardware.

# II. DEFINITION OF A STENCIL PROGRAM

As the input format of StencilFlow, we define a "stencil program" as a *directed acyclic graph of stencil operations on a structured grid* (an example is shown in Fig. 2), where each node is either a stencil operation performed on the full output domain or a memory container, and edges are dependencies between stencils and memories: i.e., outputs produced by one stencil that are consumed by one or more other stencils, and/or are read from/written to memory. Each stencil takes one or more inputs, that are sourced either from off-chip memory, or fed by a previous stencil evaluation, and produces exactly one output. To support a broader class of computations present in weather models considered, we furthermore allow stencils to read from lower-dimensional inputs: e.g., a 3D stencil can read from a 2D, 1D, or even "0D" (scalar) arrays using subsets of its indices. A stencil node is defined by:

- A definition of each logical input that is read, which we refer to as "fields", with a corresponding data type, and a sequence of offsets relative to the center ("field accesses").
- A code segment describing the computation at each point in the iteration space, where only the specified input accesses (including 0D constants) can be used in computations. Since it is important to know the latency of computations, the code is restricted to be *analyzable* (i.e., no external data structures or external functions, with the exception of standard math functions). However, ternary functions/conditionals are allowed, **including data-dependent branches**.
- A series of boundary conditions, defining how out-of-bounds accesses should be handled.

Currently supported boundary conditions include: **constant**, where out of bounds accesses are replaced with a given constant value; **copy**, where out of bounds accesses are placed by the value at offset 0 in all dimensions (the "center" value); and **shrink**, where all computed values that read out of bounds values are simply ignored in the output. The former two are specified per input, whereas shrink is specified on the output.

To facilitate productive definition of stencil programs, we define a simple JSON-based input format, which only requires the minimum amount of information necessary to instantiate the stencil DAG to be specified explicitly. An example is shown in Lst. 1. In practice, the definition must additionally provide data sources for each input field. Stencil programs can have 1, 2, or 3 dimensions, but assume all stencils iterate over the same iteration space (although they can have variable constant offsets into the output field).

## III. MAPPING TO DISTRIBUTED HARDWARE

There is a substantial body of previous work on mapping single stencil operations to reconfigurable hardware [8], [7], [15], [9], where high performance is achieved by chaining many consecutive timesteps together as a rich source of temporal locality. Some of this methodology carries over to the more general scenario we consider here, but we must additionally consider forks and joins in the stencil program, inputs and outputs shared by multiple producers and consumers, heterogeneity and complexity in stencil computations, and mapping the graph to multiple devices.

# A. Mapping to Hardware

For our hardware mapping, we work from the base assumption that every stencil operation in the dependency graph is mapped to simultaneous dedicated logic (stencil units/operators), even if this requires the design to span multiple devices. All stencil operations are scheduled simultaneously, operating in a **fully pipeline parallel manner**. In this scenario, production and consumption rates are identical across the dataflow graph, allowing the runtime to be modeled as a single, deep pipeline (described in Sec. VIII-A).

Each stencil unit executes a pipeline, which processing a number of cells equal to the product of the input dimensions, where logic required to handle out-of-bound accesses is predicated into the pipeline. The next cell is evaluated as soon as all inputs required *for that cell* are ready. This way, all dependencies between stencils become fine-grained on a per-cell level. This spatial computing view is distinct from the load/store view, as we default to *perfect data reuse* (i.e., we exploit all available temporal locality). In contrast, the efficiency of computations on load/store architectures relies on



Lst. 1: Program description.

Fig. 3: Hardware mapping.

<sup>&</sup>lt;sup>1</sup>https://github.com/spcl/stencilflow

maintaining a task granularity suitable for the architecture (large number of identical threads on GPU, small number of large tasks on CPU). Kernel fusion is thus a critical optimization to achieve the right task granularity [6] for performance through spatial locality, whereas *StencilFlow programs are executed in a fully "fused" schedule*, but are instead concerned with satisfying the off-chip and on-chip memory constraints (fusing stencil operators takes a different meaning, described in Sec. V-B).

Inputs are provided to each stencil unit through on-chip channels with a *compile-time fixed size*, where the producer can either be another stencil unit (i.e., a dependency), or a memory unit reading directly from off-chip memory. If one or more inputs are not ready, the pipeline must stall while waiting for the remaining inputs to arrive. For any DAG that is not a multi-tree, this can result in deadlocks if channel capacities are insufficient to buffer inputs ready early until inputs ready later arrive, due to the circular dependency implied by each data exchange requiring the receiver and sender to not be *full* and not be *empty*, respectively. We must thus take all paths through the DAG into account when deciding the size of buffers between dependencies.

In the example shown in Fig. 4, the stencil unit computing C requires data from both stencil units A and B to begin streaming. The results streamed out of A are also required by B. On the left hand side, C is waiting for data from B (i.e., for the data stream to not be empty), B is waiting for additional data from A, and A is waiting for C to accept the data (i.e., for the data stream to not be full), thus forming a circular dependency. Without additional buffering, this results in a deadlock. By adding an appropriate buffer between A and C (right hand side), we can inject sufficient credits to tolerate the delay induced by the path through B. We describe how StencilFlow computes the buffer depths required to prevent deadlocks and ensure continuous streaming operation in Sec. IV-B.



Fig. 4: Preventing deadlocks by injecting buffers.

## B. Mapping to the Distributed Setting

To scale beyond the off-chip memory bandwidth, on-chip memory capacity, and logic resources available on a single chip, we let designs scale to multiple devices. For modeling and code generation, this means that certain inter-stencil connections



Fig. 5: Stencil program spanning two devices.

will cross devices, and thus imply communication across the network. Furthermore, data located in off-chip memory must be present on any device that accesses it, implying potential replication to multiple devices that require it. In the example shown in Fig. 5,  $a_2$  is accessed by stencils on either device, requiring it to exist in both DRAM memories.

To implement inter-node communication in practice, we leverage the Streaming Message Interface [16] (SMI), which exposes communication as channels with FIFO semantics, resulting in inter-node communication being nearly identical to intra-device communication between stencils in the generated code. With the target in mind, the following will describe the program analysis, and the central components of the StencilFlow stack, required to build these spatial architectures.

# IV. FROM DAG TO DATAFLOW

The StencilFlow framework analyses the stencil DAG, and uses this to construct a dataflow graph that maps to efficient hardware. Data reuse happens both internally in each stencil, facilitated by "internal buffers", and on the edges between stencil nodes, referred to as "delay buffers".

# A. Internal Buffers for Intra-Stencil Reuse

The most straightforward source of temporal locality comes from within each stencil operation, where the same input field is often accessed at multiple offsets relative to the center, illustrated in Fig. 6 for accesses  $\{[-1,0], [0,-1], [0,1], [1,0]\}$ in a 2D iteration space. Furthermore, in the global dataflow setting, the core assumption of StencilFlow is that data should only be loaded once, streaming directly between kernels without going through off-chip memory.

A stencil node has 0 or 1 internal buffers per field accessed, depending on whether there are multiple accesses to the given field within the stencil. The size of each buffer is determined by the *largest distance between any two offsets in memory order*, *plus one* (or plus the vector width, in the case of vectorized



Fig. 6: Internal data reuse buffer.

Fig. 7: Varying buffer shapes.

kernels) in the stencil iteration space: e.g., in a 3D iteration space of shape  $\{K, J, I\}$ , two accesses a[0, 1, 0] and a[0, -1, 0] require buffering two 1D rows (2I+W elements, where W is the vector width), while two accesses b[0, 0, 0] and b[1, 0, 0] require buffering a 2D slice (2IJ + W), shown in Fig. 7 top and bottom, respectively. In general, buffers sizes can be up to a constant number of (D-1)-dimensional slices for a D-dimensional stencil.

StencilFlow computes the internal buffer size for each field, for each stencil, independently. However, the schedule for when the pipeline starts writing each buffer is dependent on the other fields accessed. For example, if a stencil reads multiple fields with internal buffer sizes  $\{B_1, \ldots, B_F\}$ , each internal buffer can be only start to be filled after the first  $B_i - \max\{B_1, \ldots, B_F\}$  iterations (the largest buffer(s) will always start reading immediately), so it is synchronized with the other fields. Additional accesses *in between* the "highest" and "lowest" offset in memory order do not affect the total buffer size, although they can affect the buffer implementation in practice by adding more parallel accesses into the buffer.

Filling the internal buffers also affects the latency, and transitively the runtime, of the stencil program. A stencil node cannot begin computations before all operands are available, which only happens once all internal buffers have been filled. As the size of buffers is exactly the distance between the lowest and highest accessed index in order of the stencil iteration space, the *initialization phase* of a stencil is given by  $\max\{B_1, \ldots, B_F\}$ , which is crucial to the delay buffer calculation described in the following.

# B. Delay Buffers for Inter-Stencil Reuse

Edges between stencils in the DAG enable data reuse by replacing expensive round-trips to off-chip memory with direct dataflow. Furthermore, if multiple stencils require data from the same input field, it is sufficient to read it from memory once, and stream the data to all stencils requiring it. StencilFlow exploits all such opportunities, while preventing the deadlock scenario illustrated in Fig. 4. This requires synchronizing inputs to consumers by adding buffers that delay the data (i.e., inject sufficient credits) until all inputs are ready without blocking the producer(s). We annotate these delay buffers on edges in the dataflow graph, corresponding to FIFO channel depths.

There are two factors that determine delays in the DAG. First, the AST formed by computation of a stencil operation



Fig. 8: Delay buffers on edges enable reuse and deadlock freedom.

forms another DAG, whose critical path adds a delay between a sequence of inputs entering and exiting the pipeline. Computing the critical path requires latency information for each operation performed, which is both type and architecture dependent. As a result, these latencies can be provided as configuration to the framework, and default to conservative values to account for the worst case scenario. We note that these delays are typically small (<100 cycles), and do not contribute significantly to the overall fast memory usage, even when conservatively overestimated. More importantly, delays occur in the initialization phases within each stencil, where internal buffers are being filled before enough data is available to start computations. Each stencil node in the stencil program will contribute  $\max\{B_1, \ldots, B_F\}$  elements to this delay, where  $\{B_1,\ldots,B_F\}$  is the set of F internal buffer sizes for the given stencil.

To determine the size of delay buffers on the edges arriving at a given node, we traverse the DAG backwards from the node, computing the latency contributions along all possible paths, from all possible source nodes, and for each edge, *including the contribution of the initialization phase of the node itself*, recording the highest delay encountered per edge. The buffer size on each edge is then the highest delay found for that edge, subtracted from the highest delay found across *all* edges (it follows that each node will have at least one incoming edge with delay size zero). Similar to internal buffers, the maximum size of delay buffers is proportional to the size of a (D-1)-dimensional slice of the iteration space. An example of annotated delay buffers in shown in Fig. 8.

# C. Vectorization

When insufficient reuse is present in a target program, we can employ vectorization to increase parallelism and memory bandwidth utilization, in order to approach a compute logic or memory bandwidth bound. To this end, we allow StencilFlow input programs to specify a vectorization factor, which will not only affect the generated hardware, but also the dataflow analysis. Vectorizing by a factor of W reduces the number of iterations in the inner loop of all stencils in the program by a factor of W, which affects the size of initialization phases, and transitively the delay buffers in the system. In addition to directly increasing the bandwidth requirement and parallelism in the program, vectorization can also have the subtler effect of coarsening stencil nodes, increasing the ratio of "useful" compute logic to overhead logic. We can thus also use vectorization in time tiling-like scenarios to coarsen simple stencils and increase the achievable performance.

Once the stencil program has been enriched with the appropriate internal buffer and delay buffer sizes, the resulting graph is emitted to the data-centric backend for domain-specific and low-level optimizations.

## V. DATA-CENTRIC ABSTRACT REPRESENTATION

We use the Data-Centric (DaCe) [14] framework as a dataflow representation and backend for the hardware mapping.



Fig. 9: Data-centric representation (SDFG) of a Laplace stencil.

DaCe defines a graph-based development workflow that maintains a separation of concerns between domain scientists and performance engineers, based on the observation that the vast majority of hardware optimizations are centered around data movement reduction. DaCe generates high-performance code for both load/store architectures and reconfigurable hardware, and supports high-level synthesis (HLS) backends for both Xilinx and Intel FPGA architectures.

DaCe separates program definition from its optimization by using the the Stateful DataFlow multiGraph (SDFG) representation. In each multigraph, data movement (edges) is explicitly separated from data/FIFO containers and computations (nodes). These acyclic multigraphs are, in turn, nested within state machines (directed graphs) that represent the control flow of the application. An example of a two-dimensional Laplace operator is shown in Fig. 9.

In the DaCe workflow, a program is developed in a frontend language (e.g., Python, StencilFlow) by a domain scientist. All subsequent hardware mapping and optimizations are performed on the SDFG separately, by a performance engineer. Optimization and hardware mapping is achieved via graph rewriting rules, called *transformations*, for data movement reshaping, scheduling parallel subgraph scopes to processors, and modifying data storage/layout. Transformations are user-extensible and written in Python interfaces, allowing both domain-specific and general purpose optimizations, and enabling knowledge transfer between applications. Adaptations to the graph are saved separately from the source code, allowing the original source code to be modified without changing the optimization scheme.

# A. Extensions to DaCe

To support this work, we extend the DaCe framework, introducing a new type of dataflow node, pipelined scopes, and three new transformations. In particular, we extend the SDFG with the concept of *library nodes*. Library nodes function similarly to computational nodes, but encode domainspecific information and contain multiple implementation targets, which translate into different subgraphs upon expansion. The StencilFlow-specific library node Stencil was developed for this work, and will be used extensively throughout the following. Since the high-level semantics of library node types are known, they allow performance engineers to develop domain-specific transformations, such as algebraic contractions (e.g., double transposition) and others. With library node expansions potentially containing other library nodes, multilevel coarsening and transformations are thus enabled in SDFGs, inspired by the MLIR [17] stack.

As a useful shorthand for pipelined iteration spaces, we introduce the *pipeline scope*, augmented with information on initialization and draining phases, to easily allow the programmer to inject specialized behavior during initialization, streaming, and draining phases. For StencilFlow, this allows encoding the internal buffer initialization phase, and draining phases where results are still being computed only using data present in local buffers, thus omitting reads from inputs.

With the domain-specific concepts enabled by library nodes, we are now able to develop transformations for stencil programs on reconfigurable hardware. We develop both domain-specific and a general-purpose transformation, summarized in Fig. 10. NestDim reschedules stencil computations by taking multiple, parametrically-parallel stencils and creating one stencil, which can be mapped into different schedules on hardware. StencilFusion schedules multiple dependent stencils as one stencil with multiple statements, differing from standard map fusion by taking boundary conditions and redundancy into account. For general-purpose transformations, we add the MapFission transformation, which splits a parallel subgraph into multiple parallel subgraphs (which can in turn be rescheduled), introducing temporary storage between the subgraph components. The NestDim and MapFission transformations are used as a tool to extract stencil programs from existing SDFGs to analyze them in StencilFlow, while StencilFusion is an optimization for both load/store and spatial architectures, described in the context of StencilFlow below.

#### B. Spatial Stencil Fusion

On load/store architectures, fusing consecutive stencils is used to increase performance by improving data locality, reducing write/read roundtrips from off-chip memory, and reducing context scheduling overhead [18]. When applying the transformation on StencilFlow dataflow graphs, the effect is somewhat different, as the schedule of the spatial architecture is already fully "fused" into a global pipeline. Instead, fusing stencils has the following effects:



Fig. 10: Transformations used. (DS: Domain-Specific, GP: General-Purpose).



(a) Load/store stencil fusion. (b) Stencil fusion in StencilFlow.

Fig. 11: Unlike load/store fusion, spatial fusion only reduces latency.

- The critical path through the program can be reduced by combining the initialization phases (see Sec. IV-A) of two consecutive stencils.
- Internal buffers for the same input field are combined into a single internal buffer.
- Multiple smaller delay buffers can be combined into fewer, larger buffers, which affects hardware utilization, depending on the granularity of on-chip memory on the target platform.
- Combined code sections increase the opportunity for common subexpression elimination by the optimizing compiler.
- Coarser stencil nodes increase the ratio of "useful" logic to the number of pipelines instantiated, which can affect spatial resource overhead.

The difference between fusing tasks on load/store architectures and the spatial fusion performed here is illustrated in Fig. 11. On load/store architectures (Fig. 11a), the total number of scheduled kernels is reduced when fusing task 0 and task 2 into a single kernel. In Fig. 11b, all operators are already scheduled in parallel, but the initialization latency can be reduced *if the fused nodes*  $s_A$  and  $s_B$  are on the critical path.

In our dataflow canonicalization pass, we define a collection of heuristics for fusing two stencils so that these effects are observed. Firstly, the necessary conditions for fusion are checked, namely that the two stencils operate on the same data shape (correlating to iteration space) and that they have the same StencilFlow boundary condition definitions. Then, we only consider stencils that are connected by one data container node u with deg(u)=2, in order to ensure that all stencils (fused or otherwise) have a single output. Finally, we ensure no other instances of u exist in other states, so that it can be completely removed from the graph without adding an extra write to off-chip memory.

For the experiments in this work, we perform aggressive stencil fusion of input programs, as this is observed to reduce overall logic through the coarsening of stencil nodes, and slightly reduces runtime by pruning initialization latencies.

# VI. CODE GENERATION

StencilFlow relies on DaCe backends to generate the final kernel code, which is passed to an optimizing compiler. For the experiments performed in this work, we target the Intel FPGA SDK for OpenCL backend [19], which is a high-level synthesis (HLS) compiler, emitting RTL code from annotated OpenCL. Being an HPC-oriented framework, DaCe automatically performs necessary annotations for pipelining, unrolling, and coalescing loops emitted from parametric maps in the dataflow graph, splits parallel sections into processing elements (i.e., OpenCL kernels), annotates buffer depth properties for channels, declares kernels as autorun when possible, inlines constants, and performs conversions between vectorized and non-vectorized data types. Host code necessary to interface with the kernel and the necessary memory copies are generated, and the final program can be called by using the high-level Python interface. By using DaCe for code generation and by using the library node abstraction for stencil computations, supporting Xilinx FPGAs, emitting RTL code directly, or targeting other spatial systems entirely will only require adapting the stencil library node expansion, provided that support for the desired architecture is present in/added to the DaCe framework.

#### A. Intel FPGA Optimizations in DaCe

When targeting the Intel HLS compiler, delay buffers are represented as DaCe streams with a given buffer size, which are mapped to the Intel OpenCL channel abstraction, that in turn are mapped to FIFOs in hardware. We target the shift register pattern in Intel's OpenCL compiler to efficiently implement internal buffers within each stencil node. To achieve this in DaCe, a data container spanning the full width of each internal buffer is created, injecting and shifting elements every cycle. "Tap" points (constant offset accesses) into the array are then connected to the stencil, where offsets are generated from the distance between accesses flattened into a 1D iteration space.

The processing done per cell in an expanded stencil library node is shown in Fig. 12. The graph contains three consecutive components: a shift phase, containing a fully unrolled map scope (trapezoids) where a "tasklet" (octagon) shifts each entry of the shift register memory (ovals) by the vectorization width to i+W; an **update** phase, where new values are read from the input channel (dashed border) into the front of each shift register by a tasklet; and a **compute** phase, where the buffers are accessed at all tap points and fed to the main computation tasklet, which is parametrically unrolled to treat each element in the vector with potentially different boundary conditions, and passes through another tasklet that conditionally writes the output stream if the stencil is not in the initialization phase. This full graph will be wrapped in a parametric scope that defines the iteration space of the stencil program, which is fully pipelined, such that all three phases are executed in a pipeline parallel manner. The input and output streams (dashed borders) are connected to the appropriate producers and consumers in



Fig. 12: Data movement occuring per iteration of a stencil.

the global dataflow graph. Source nodes are instantiated as dedicated prefetchers that can read ahead of computations, and dedicated writers are instantiated at sink nodes that can buffer data while waiting for DRAM writes.

## B. Generating Distributed Programs

To code-generate distributed implementations, we integrate an OpenCL implementation of the *Streaming Message Interface* [16] into the DaCe backend. SMI extends HLS with a distributed memory programming model for reconfigurable hardware that unifies message passing with pipelined streambased communication of data, such that cross-chip communication is expressed the same way as on-chip communication.

When a stencil program spans multiple devices, the computation running on each device is represented by a separate DaCe program, as it will compile to separate bitstreams that must be configured to each device in the sequence. Devices communicate via *remote streams*, which are DaCe streams annotated with having a source/destination located on a different device, which will trigger the SMI backend to code generate the relevant networking code and emit streaming message communication.

If multiple network connections are present between two endpoints, SMI can split a communication stream into two or more substreams following different channels across the network, and recombine them at the other end, allowing for a multiplicative increase in achievable bandwidth. In StencilFlow, we exploit this to increase the vectorization width and number of channels spanning across devices.

# C. Reference Code

Using domain-specific stencil nodes within the DaCe framework, we are able to maintain a high-level view of the program, which enables optimization for different architectures. While exploring CPU and GPU performance is out of the scope of this work, we exploit this capability to generate reference CPU-executed graphs where stencil evaluations are executed sequentially in topological order (i.e, no fusion or parallelism between stencil evaluations), which we can verify against the generated hardware kernels.

# VII. WORKFLOW AND ARTIFACTS

To summarize the stack described throughout the above, an overview of the StencilFlow workflow is shown in Fig. 13. The StencilFlow framework is a pure Python code ( $\approx$ 5,300 SLOC at the time of writing) developed for the purpose of this work. The DaCe framework was extended with Python and C++ features to support domain-specific nodes.

The input program to StencilFlow can either be given as the JSON-based program description described in Sec. II, or as a DaCe dataflow graph containing domain-specific stencil nodes. In the latter case, we developed software that performs canonicalization passes to the DaCe graph, before extracting the stencil pattern to the standard program description format. This allows us to read in external programs, which will be required for the case study in Sec. IX.



Fig. 13: Workflow overview, with code artifacts annotated on arrows. Dashed outline indicates an existing feature that was extended.

StencilFlow can directly run the stencil program from the input description, transparently executing parsing, dependency analysis, buffering analysis, SDFG generation, domain-specific optimization, library node expansion, general purpose optimization, code generation, compilation of the host code, compilation of the kernel (requiring the full synthesis, placement and routing flow if FPGAs are targeted), execution of the program, and validation of results.

#### VIII. BENCHMARKS

We benchmark the architectures emitted by StencilFlow to establish the highest achievable performance and bandwidth on a testbed platform, which we can use to analyze the characteristics required to push performance of stencil applications.

#### A. Computing Expected Runtime

We annotate benchmarks with the "expected" runtime, given by the lower bound on number of cycles required to evaluate the program, assuming all data is available at the earliest possible cycle. Because the full stencil DAG is executed in a pipeline parallel manner, we can model the runtime as a single, global pipeline. It is generally true for a pipelined circuit that the number of cycles required to process N inputs is

$$C = L + I \cdot N,\tag{1}$$

where L is the latency of the pipeline, and I is the initiation interval (i.e., the number of cycles between allowing a new set of inputs to the circuit) [20]. All architectures emitted by StencilFlow are fully pipelined, so we fix  $I=1 \frac{\text{cycles}}{\text{operand}}$ . N is the product of the domain dimensions (number of iterations in the iteration space), divided by the vectorization width Wwhen applicable. L is computed from the circuit latency and initialization delay described in Sec. IV-B. N and L compose differently: N covers the streaming section where stencils can operate in a pipeline parallel fashion, whereas L covers the initialization phase where stencil units are not feeding downstream consumers. The depth of the DAG thus adversely affects the performance upper bound, while the size of the domain affects it favorably, increasing the relatively number of "useful" cycles to cycles spent in initialization. Since L is only proportional to D-1 or fewer dimensions (see Sec. IV-B), it becomes negligible when the domain is large relative to

the depth of the stencil DAG. However, we include it when computing expected runtime for completeness.

#### **B.** Experimental Platform

To evaluate the efficiency of dataflow architectures laid out by StencilFlow, we map them to a state-of-the-art FPGA platform. We target the BittWare 520N PCI-e attached board, with an Intel GX 2800 Stratix 10 processor, 4 DDR4 memory banks with a combined peak bandwidth of 76.8 GB/s, and four networkattached QSFP ports rated at 40 Gbit/s. The annotated OpenCL code generated from DaCe is compiled with version 19.1.0 of the Intel FPGA OpenCL SDK and Quartus compiler, targeting the p520 max sg2801 shell offered by BittWare. This shell supports networking via OpenCL channels, which we target using the SMI library (Sec. VI-B). The FPGAs are installed in the Noctua cluster at the Paderborn Center for Parallel Computing, which exposes a programmable, fully connected optical switch, allowing us to chain FPGAs together in a sequence with two 40 Gbit/s links between each consecutive device to explore multi-device scaling. Our benchmarks focus on 32-bit precision, as this is used in production by our motivating weather simulation example, and because this precision is supported natively on the Stratix 10. However, all parts of the StencilFlow stack support any data type recognized by the underlying compiler, including double precision floating point and integer types.

# C. Iterative Stencil Performance

StencilFlow is built to handle complex stencil kernels, but is also capable of processing traditional, iterative-style stencil codes. We produce benchmarks using such kernels to establish the highest floating point performance reachable by StencilFlow, which can be compared to previous work. This is achieved by chaining together long linear sequences of stencils executed on a large input domain, analogous to time-tiled iterative stencils.



Fig. 15: Performance scaling with 4-way vectorization.

To evaluate scaling behavior of an iterative stencil, we gradually increase the number of chained stencil computations until a single FPGA device is fully utilized, then we continue the chain across multiple devices by replacing accesses to on-chip FIFOs with network channels. We repeat the experiment with and without vectorization, to see the effect of coarsening stencil stages. The resulting benchmarks are shown in Fig. 14 and Fig. 15 without and with vectorization, respectively.

Without vectorization, the highest performing bitstream yields 264 GOp/s on a single device, and scales up to 1.5 TOp/s across 8 FPGAs. A 4-way vectorized code reaches 568.2 GOp/s and 4.2 TOp/s on single and multi-device, respectively. Vectorization thus proves to be crucial to achieve high utilization of compute resources on the Stratix 10, as it reduces the ratio of overhead logic to computational logic. This further motivates the necessity of the stencil fusion transformation (Sec. V-B) on input programs to coarsen the granularity of stencil nodes. Frequencies across all benchmarks are consistently in the range 292-317 MHz, which is factored into the upper bound calculation shown as dashed black lines, computed from Eq. 1 as C/f, where f is the design frequency.

We additionally measure the highest performance achievable *without* networking on a single device, as we are unable to vectorize the stencils in the distributed experiment further due to the network bandwidth bottlenecking the computation, included in Tab. I. As a non-relative measure of device utilization, the table includes resource usage for the maximum performing stencil for each data type. The highest measured stencil performance of 1.3 TOp/s and 4.2 TOp/s marks a  $9.4 \times$  and  $30 \times$  speedup over the stencil performance reported for a single VCU1525 device in the original work on DaCe for single-device and multi-device, respectively (which in turn outperformed a state-of-the-art HLS compiler by five orders of magnitude, showing in inability of HLS compilers to yield satisfactory out-of-the-box performance).

 TABLE I

 HIGHEST PERFORMING KERNELS AND THEIR RESOURCE USAGE.

	Performance	ALM	FF	M20K	DSP	
Total		103 M	3.7 M	11.7 K	5760	
Avail.		$692~\mathrm{K}$	$2.8 \mathrm{M}$	$8.9~\mathrm{K}$	4468	
Jacobi 3D	265 COm /a	233 K	534 K	1495	784	
(Ours)	205 GOp/s	33.6%	19.3%	16.7%	17.6%	
Jacobi 3D	001 CO- /-	$437~\mathrm{K}$	$1207~\mathrm{K}$	2285	3072	
W=8 (Ours)	921 GOp/s	63.1%	43.6%	25.5%	68.8%	
Diffusion 2D	1 212 COm /a	449 K	$1329~\mathrm{K}$	2565	2304	
W=8 (Ours)	1,515 GOp/s	64.8%	48.0%	28.6%	51.6%	
Diffusion 3D	$1,\!152{\rm GOp/s}$	567 K	1606 K	5357	3072	
W=8 (Ours)		81.9%	57.9%	59.8%	68.8%	
Diffusion 2D (Zohouri et. al. [8]) Diffusion 3D (Zohouri et. al. [8])	$913\mathrm{GOp/s}$	471.4 K	1173.6 K	2204	3844	
		68.0%	42.3%	24.6%	86.0%	
	$934{\rm GOp/s}$	450.5  K	1078.2 K	8684	3592	
		65.0%	38.9%	97.0%	80.4%	
Waidyasooriya and Hariyama [21]	$630{ m GOp/s}$	Arria 10 GX 1150				
SODA [9]	$135  \mathrm{GOp/s}$	ADM-PCIE-KU3				
Niu et al. [22]	$119  \mathrm{GOp/s}$	Virtex-6 SX475T				
Ben-Nun et al. [14]	$139\mathrm{GOp/s}$	Virtex UltraScale+ VCU1525				

For a more direct comparison on the Stratix 10 platform, we compare StencilFlow to a handwritten stencil implementation. Zohouri et al. [8] combine spatial and temporal blocking in an HLS-based design to achieve high performance on stencil codes on an Arria 10 FPGA. We extend the authors' work by building their code<sup>2</sup> for Stratix 10, using the Diffusion 2D and 3D stencil codes. On advice from the authors, we configure parameters to a vectorization width of 16, run enough repetitions that the kernel runs for multiple seconds to hide initialization overhead, and disable burst interleaving. We include the resulting performance in Tab. I, along with other previous results by Niu et al. [22] and Waidyasooriya and Hariyama [21], showing that StencilFlow is competitive even with hand-tuned code. We also consider frameworks emitting stencil FPGA code, including the Jacobi 3D result of SODA [9], which is the stencil backend of HeteroHalide [23] and HeteroCL [24]. For previous work we note the FPGA used for evaluation by the respective authors. We do not compare quantitatively to HeteroCL and Wang and Liang [25], as the authors do not report absolute performance numbers.

# D. Off-Chip Memory Bandwidth

To measure achievable off-chip memory bandwidth by StencilFlow programs, we run two series of benchmarks: first, we measure the effective bandwidth utilization when scaling up number of accesses, but accessing only 32-bits per cycle at each access point. This stresses the routing on the device to deliver data to all end-points every cycle. Second, we request the same total number of 32-bit operands, but at fewer, vectorized endpoints, requiring more operands per cycle per endpoint. We found the -global-ring and -duplicate-ring options to the Intel FPGA OpenCL compiler to significantly increase the number of parallel access points supported in the architecture before designs dropped in frequency. The resulting benchmarks, along with the analytically computed performance upper bound, are shown in Fig. 16. For the non-vectorized green bars, the x-axis corresponds to the number of access points, while the number of access points for the vectorized orange bars is the number of operands divided by the vector size of 4 (i.e., up to 12 access points are depicted).

After 24 parallel access points, we see a decrease in effective memory performance relative to peak, flattening out



Fig. 16: Effective bandwidth with number of operands requested per cycle (i.e., number of operands served if infinite bandwidth).

<sup>2</sup>https://github.com/zohourih/Diffusion\_FPGA, commit 96588e2.

at 36.4 GB/s, which is 36.4/76.8 = 47% of peak bandwidth. This marks the limit of the memory controller crossbar, and of routing a large number of memory accesses across the device. The 4-way vectorized scenario allows for higher achievable bandwidth, but experiences a drop in efficiency at a lower number of access points ( $0.94 \times$  at 12 access points), and flattens out at 58.3 GB/s, which is 76% of peak bandwidth. No further increase was seen with more access points, and 8-way vectorized programs achieve similar bandwidth.

#### IX. WEATHER SIMULATION APPLICATION STUDY

To stress the full capability of the StencilFlow stack we evaluate the horizontal diffusion stencil program, a large reallife weather simulation kernel from the COSMO weather model. Horizontal diffusion is a 4th order explicit method performed on a staggered latitude-longitude grid with Smagorinsky diffusion to smoothen wind velocity components [26]. We obtain the program from an input SDFGs using stencil library nodes, shown for horizontal diffusion in Fig. 17a, applying the NestDim and MapFission transformations described in Sec. V-A, resulting in an SDFG as the one shown in Fig. 17b, from which the stencil program is extracted. The DAG in Fig. 17c is created after aggressively fusing consecutive stencils (see Sec. V-B). In the fully fused program, initialization latency (L in Eq. 1) accounts for  $\sim 0.7\%$  of the total number of iterations required to evaluate the program, and is thus a negligible overhead. This program is run in production by the Swiss Federal Office of Meteorology and Climatology (MeteoSwiss), where simulations are performed with 32-bit floating point on an NVIDIA Pascal Tesla K80 cluster. We compare StencilFlow to the stronger TSMC 16 nm Tesla P100 GPU on the same architecture (comparable release window to the Stratix 10), a TSMC 12 nm Tesla V100 Volta GPU, and a 12-core Xeon CPU.

#### A. Horizontal Diffusion Analysis

The horizontal diffusion DAG characterized by a high number of stencils reading the same input locations (28 accesses of 10 unique fields), allowing for the communication volume between them to be consolidated via delay buffers, as well as complex dependencies between stencil nodes (each nonsource stencil receives data from 2-6 other stencil nodes). This requires the full complexity of an arbitrary DAG, and allows us to stress the full stack of StencilFlow.

Floating point operations in the DAG include 87 additions, 41 multiplications, and 2 square roots, in addition to 2 minimum and 2 maximum operations, and ternary operations resulting in 20 data-dependent branches. With maximum reuse of all input fields and all computed fields (i.e., perfect locality), the program reads 5IJK+5I operands and writes 4IJK operands, for a total of 9IJK+5I operands. Considering floating point arithmetic only, this implies a upper bound arithmetic intensity of (square root is counted as one operation):

$$\frac{(87+41+2)IJK \text{ [Ops]}}{9IJK+5I \text{ [operands]}} \approx \frac{130}{9} \left[\frac{\text{Ops}}{\text{operand}}\right]$$



(b) Canonicalized SDFG.



(c) DAG of the transformed program inferred by StencilFlow.

Fig. 17: Horizontal diffusion stencil program from the COSMO weather and climate model.

which for 32-bit floating point corresponds to

$$\frac{130/9 \text{ [Ops/operand]}}{4 \text{ [B/operand]}} = \frac{65}{18} \left[\frac{\text{Ops}}{\text{B}}\right].$$
 (2)

Using the benchmark of practically achievable bandwidth presented in Sec. VIII-D for the Stratix 10 FPGA, the highest achievable performance in roofline model [27] terms is:

$$\frac{65}{18} [\text{Ops/B}] \cdot 58.3 [\text{GB/s}] = 210.5 [\text{GOp/s}], \quad (3)$$

or 277.3 GOp/s at the peak data sheet bandwidth of 76.8 GB/s. This is well below what is achievable by a stencil program with higher arithmetic intensity (see Sec. VIII-C), indicating that high bandwidth is required to shine in realistic stencil applications. We compute the bandwidth required to saturate the compute performance measured in Sec. VIII-C for the arithmetic intensity of the studied program to be:

$$\frac{917.1 \ [\text{GOp/s}]}{65/18 \ [\text{Op/B}]} = 254.0 \ [\text{GB/s}].$$
(4)

The ideal logic to bandwidth ratio is thus off from the ideal ratio by a factor of  $\sim 3-4$  on the target Stratix 10 platform. To explore the performance potential of the Stratix 10 without this memory bottleneck, we will include experiments with simulated "infinite" memory bandwidth, by replacing memory accesses with compile-time constants fed to the computational circuit (and omitting validation of functional correctness).

## B. Horizontal Diffusion Benchmark

We compile the DAG in Fig. 17c for the Stratix 10 from the constructed dataflow graph by StencilFlow. As shown in the analysis above, the program is bandwidth-bound on this platform, which requires us to saturate the bandwidth to maximize performance. Without vectorization, the pipelined circuit requires approximately 9 operands/cycle, corresponding to 10.8 GB/s at 300 MHz for single precision floating point. We thus vectorize the program by a factor of 8 for a maximum bandwidth of 86.4 GOp/s, in addition to building a 16-way vectorized kernel with simulated input memory to evaluate performance without the memory bottleneck. We target a  $128 \times 128 \times 80$  domain size, which is used for performance benchmarking by MeteoSwiss. Specifically, a  $128 \times 128$  horizontal domain is stacked in 80 vertical layers. In addition to runtime and the effective performance, we consider peak memory bandwidth and the associated fraction of highest achievable performance *for the given arithmetic intensity* computed according to Eq. 2 (%Roof.). The results are listed in Tab. II.

TABLE II HORIZONTAL DIFFUSION BENCHMARKS.

	Runtime	Performance	Peak BW.	%Roof.
Stratix 10	$1,\!178\mu\mathrm{s}$	$145\mathrm{GOp/s}$	$77\mathrm{GB/s}$	52%
Stratix 10*	$332\mu s$	$513{ m GOp/s}$	$\infty{ m GB/s}$	—
Xeon 12C	$5{,}270\mu{ m s}$	$32\mathrm{GOp/s}$	$68\mathrm{GB/s}$	13%
P100	$810\mu s$	$210  \mathrm{GOp/s}$	$732\mathrm{GB/s}$	8%
V100	$201\mu s$	$849\mathrm{GOp/s}$	$900\mathrm{GB/s}$	26%

<sup>\*</sup>Without memory bandwidth constraints.

We include CPU and GPU performance as a point of comparison, using a 12-core Intel Xeon 2.60/3.50 GHz E5-2690V3 CPU, and NVIDIA Tesla P100 and V100 GPUs, compiled with CUDA v10.1 and gcc 8.3.0. The application is synthesized using the MeteoSwiss Dawn [28] stencil-optimizing compiler toolchain<sup>3</sup>, which was also used to generate the StencilFlow input program. Dawn is specifically designed to optimize weather and climate stencil programs for GPU and CPU, employing data movement optimizations, GPU kernel fusion, CPU multi-threading, vectorization, and efficient GPU boundary scheduling. The domain size of  $128 \times 128 \times 80$  is sufficient for saturating the GPU thread scheduler (i.e., larger domains do not significantly increase GPU performance). The horizontal diffusion program emitted by Dawn for CPU and

<sup>&</sup>lt;sup>3</sup>https://github.com/MeteoSwiss-APN/dawn, commit 4ae6dc0.

GPU executes five components of horizontal diffusion as distinct kernels. We omit kernel launch overhead and report the raw kernel execution time only, included in Tab. II.

The FPGA platform outperforms the CPU by  $4.5 \times$  and is outperformed by either GPU, but comes closest to the upper bound (Eq. 3) imposed by its roofline characteristics at the given arithmetic intensity: 52% of the bandwidth upper bound (69% of the highest measured bandwidth in Sec. VIII-D), at 26% ALMs, 27% DSPs, and 20% M20K utilization, respectively. The benchmark simulating infinite memory bandwidth shows significant headroom for pushing the performance at this arithmetic intensity with higher bandwidth off-chip memory: without the bandwidth bottleneck, the Stratix 10 would outperform the P100, but falls at 60% of the performance of the V100, at 46% ALMs, 48% DSPs, and 20% M20Ks.

# C. Silicon Efficiency

The Stratix 10 is estimated to be a  $700 \text{ mm}^2$  die [29] (half the Stratix 10M, which fuses two Stratix 10 chiplets) on Intel's 14 nm process, compared to  $610 \text{ mm}^2$  on TSMC 16 nm and  $815 \text{ mm}^2$  on TSMC 12 nm for the P100 and V100, respectively. Using the benchmarks from Tab. II, this amounts to a silicon efficiency of 0.21 and 0.71  $\frac{\text{GOp}/\text{s}}{\text{mm}^2}$  with and without the memory bottleneck for the Stratix 10, respectively;  $0.34 \frac{\text{GOp}/\text{s}}{\text{mm}^2}$  for P100; and  $1.04 \frac{\text{GOp}/\text{s}}{\text{mm}^2}$  for the V100, when performing the horizontal diffusion experiment.

## D. Spatial Tiling

We have not considered spatial tiling, as on-chip memory requirements were not a restriction for building the large weather stencil program evaluated. Both memory bandwidth and logic were bottlenecks before on-chip memory capacity, despite minimizing off-chip memory bandwidth in the program. Eventually, increasing the domain size will scale the internal buffer and delay buffer sizes beyond what is feasible to buffer in on-chip memory. Spatial tiling can be employed in this scenario, introducing redundant computation at the domain boundaries proportional to the DAG depth and the tile surfaceto-volume ratio. This is primarily a scheduling challenge, which can be efficiently solved in practice [8].

#### X. RELATED WORK

There are numerous works on stencil accelerators on FP-GAs [15], [8], [21], including for multi-device settings on up to 9 interconnected FPGAs [7], all of which we have considered throughout this work. Other frameworks generating stencil architectures have also been proposed [9], [24], [23], [25], which we consider in Sec. VIII-C. Common to these works is that they treat a single stencil operation applied iteratively, allowing them to unroll the time dimension as a source of temporal locality. StencilFlow is on a par or outperforms all the above on simple iterative stencils, and treats a much wider range of input programs. Niu et al. [22] explore runtime reconfiguration of an FPGA to eliminate idle operators during program execution. Runtime reconfiguration is not beneficial for stencil programs considered by StencilFlow, as all operators

are assumed to operate in the same iteration space and fully in parallel after the initialization phase.

Darkroom [30] is a framework producing spatial accelerators of image processing pipelines from a high-level input DSL. StencilFlow takes a similar approach, but accepts a wider scope of input programs: in particular arbitrary DAGs of stencils, and 3D input/output domains. Other DSLs [31], [32] do not consider spatial computing architectures.

For the application study, Singha et al. [33], [34] present a hand-tuned implementation of the horizontal diffusion application targeting an FPGA+CPU coherent system. The authors report 129.9 GOp/s on an ADM-PCIE-9V3 board with the NARMADA accelerator, and 485.4 GOp/s on an ADM-PCIE-9H7 board with the NERO accelerator, the latter owing its large increase in performance to the introduction of HBM memory, effectively eliminating the memory bottleneck described by Eq. 4. The fully code generated kernels emitted by StencilFlow outperform the DDR4-based accelerator when memory bound, and the HBM-based when compute bound (i.e., when high memory bandwidth is simulated).

## XI. CONCLUSION

We introduced StencilFlow, an end-to-end analysis, optimization and code-generation stack built on the DaCe framework, enabling the generation of complex high-performance stencil programs on spatial architectures from a high-level input DSL. Based on a DAG representation, StencilFlow automatically insert buffers within and between stencil operations to achieve perfect reuse of all data in the program. Architectures emitted by StencilFlow achieve the highest recorded single-device performance of 1.31 TOp/s, and the highest recorded multidevice performance of 4.18 TOp/s on 8 FPGAs. We demonstrated the domain complexity supported by the framework by treating a large stencil program used in production for weather prediction, comparing the generated architecture to state-ofthe-art GPU and CPU performance. We release StencilFlow as open source software, enabling reproducibility and allowing scientists to easily target spatial computing accelerators with complex stencil programs.

# Acknowledgements

We authors wish to thank Tobias Kenter, Christian Plessl, and the Paderborn Center for Parallel Computing ( $PC^2$ ) for generously providing support and compute hours on the Noctua FPGA cluster; the Swiss National Supercomputing Center (CSCS) for providing computing infrastructure; and Jakub Beránek for last-minute engineering contributions. This work was supported by the European Research Council under the European Union's Horizon 2020 programme (grant agreement DAPP, No. 678880). Tal Ben-Nun is supported by the Swiss National Science Foundation (Ambizione Project No. 185778).

#### REFERENCES

 Intel's exascale dataflow engine drops x86 and von Neumann. Published on The Next Platform, August 2018. Accessed June 5, 2020. [Online]. Available: https://www.nextplatform.com/2018/08/30/ intels-exascale-dataflow-engine-drops-x86-and-von-neuman/

- [2] B. Gaide, D. Gaitonde, C. Ravishankar, and T. Bauer, "Xilinx adaptive compute acceleration platform: Versal architecture," in *Proceedings of* the 2019 ACM/SIGDA International Symposium on Field-Programmable Gate Arrays (FPGA'19), 2019.
- [3] K. Rocki, D. V. Essendelft, I. Sharapov, R. Schreiber, M. Morrison, V. Kibardin, A. Portnoy, M. S. Jean Francois Dietiker, and M. James, "Fast stencil-code computation on a wafer-scale processor," in *Proceedings* of the International Conference for High Performance Computing, Networking, Storage and Analysis (SC'20), 2020.
- [4] K. Matsumura, H. R. Zohouri, M. Wahib, T. Endo, and S. Matsuoka, "AN5D: automated stencil framework for high-degree temporal blocking on GPUs," in *Proceedings of the 18th ACM/IEEE International Symposium on Code Generation and Optimization (CGO'20)*, 2020.
- [5] T. Gysi, T. Grosser, and T. Hoefler, "MODESTO: Data-centric analytic optimization of complex stencil programs on heterogeneous architectures," in *Proceedings of the 29th ACM on International Conference on Supercomputing (ICS'15)*, 2015, pp. 177–186.
- [6] —, "Absinthe: Learning an analytical performance model to fuse and tile stencil codes in one shot," in 2019 28th International Conference on Parallel Architectures and Compilation Techniques (PACT'19). IEEE, 2019, pp. 370–382.
- [7] K. Sano, Y. Hatsuda, and S. Yamamoto, "Multi-FPGA accelerator for scalable stencil computation with constant memory bandwidth," *IEEE Transactions on Parallel and Distributed Systems (TPDS)*, vol. 25, no. 3, 2013.
- [8] H. R. Zohouri, A. Podobas, and S. Matsuoka, "Combined spatial and temporal blocking for high-performance stencil computation on FPGAs using OpenCL," in *Proceedings of the 2018 ACM/SIGDA International Symposium on Field-Programmable Gate Arrays (FPGA'18)*, 2018.
- [9] Y. Chi, J. Cong, P. Wei, and P. Zhou, "SODA: Stencil with optimized dataflow architecture," in 2018 IEEE/ACM International Conference on Computer-Aided Design (ICCAD'18), 2018.
- [10] M. Baldauf, A. Seifert, J. Förstner, D. Majewski, and M. Raschendorfer, "Operational convective-scale numerical weather prediction with the COSMO model: Description and sensitivities." *Monthly Weather Review*, 139:3387–3905, 2011.
- [11] T. Weusthoff, F. Ament, M. Arpagaus, and M. W. Rotach, "Assessing the benefits of convection-permitting models by neighborhood verification: Examples from map d-phase." *Monthly Weather Review*, 138:3418–3433, 2010.
- [12] W. Langhans, J. Schmidli, O. Fuhrer, S. Bieri, and C. Schär, "Longterm simulations of thermally driven flows and orographic convection at convection-parameterizing and cloudresolving resolutions." *Journal of Applied Meteorology and Climatology*, 52:1490–1510, 2013.
- [13] A. Possner, E. Zubler, O. Fuhrer, U. Lohmann, and C. Schär, "A case study in modeling lowlying inversions and stratocumulus cloud cover in the bay of Biscay." Weather and Forecasting, 29(2):289–304, 2014.
- [14] T. Ben-Nun, J. de Fine Licht, A. N. Ziogas, T. Schneider, and T. Hoefler, "Stateful dataflow multigraphs: A data-centric model for performance portability on heterogeneous architectures," in *Proceedings of the International Conference for High Performance Computing, Networking, Storage and Analysis (SC'19),* 2019.
- [15] H. Fu and R. Clapp, "Eliminating the memory bottleneck: an FPGAbased solution for 3D reverse time migration," in *Proceedings of the* 2018 ACM/SIGDA International Symposium on Field-Programmable Gate Arrays (FPGA'11), 2011.
- [16] T. De Matteis, J. de Fine Licht, J. Beránek, and T. Hoefler, "Streaming Message Interface: High-performance distributed memory programming on reconfigurable hardware," in *Proceedings of the International Conference for High Performance Computing, Networking, Storage and Analysis* (SC'19), 2019.
- [17] C. Lattner, J. Pienaar, M. Amini, U. Bondhugula, R. Riddle, A. Cohen, T. Shpeisman, A. Davis, N. Vasilache, and O. Zinenko, "MLIR: A compiler infrastructure for the end of Moore's law," *arXiv preprint abs/2002.11054*, 2020.

- [18] P. Hudak and B. Goldberg, "Serial combinators: "optimal" grains of parallelism," in *Conference on Functional Programming Languages and Computer Architecture (FPCA'85)*, 1985.
- [19] T. S. Czajkowski, U. Aydonat, D. Denisenko, J. Freeman, M. Kinsner, D. Neto, J. Wong, P. Yiannacouras, and D. P. Singh, "From OpenCL to high-performance hardware on FPGAs," in 22nd international conference on field programmable logic and applications (FPL'12), 2012.
- [20] J. de Fine Licht, M. Besta, S. Meierhans, and T. Hoefler, "Transformations of high-level synthesis codes for high-performance computing," *IEEE Transactions on Parallel and Distributed Systems (TPDS)*, vol. 32, pp. 1014–1029, May 2021.
- [21] H. M. Waidyasooriya and M. Hariyama, "Multi-FPGA accelerator architecture for stencil computation exploiting spacial and temporal scalability," *IEEE Access*, vol. 7, 2019.
- [22] X. Niu, T. C. P. Chau, Q. Jin, W. Luk, Q. Liu, and O. Pell, "Automating elimination of idle functions by runtime reconfiguration," ACM Transactions on Reconfigurable Technology and Systems (TRETS), vol. 8, no. 3, 2015.
- [23] J. Li, Y. Chi, and J. Cong, "HeteroHalide: From image processing DSL to efficient FPGA acceleration," in *The 2020 ACM/SIGDA International Symposium on Field-Programmable Gate Arrays (FPGA'20)*, 2020.
- [24] Y.-H. Lai, Y. Chi, Y. Hu, J. Wang, C. H. Yu, Y. Zhou, J. Cong, and Z. Zhang, "HeteroCL: A multi-paradigm programming infrastructure for software-defined reconfigurable computing," in *Proceedings of the* 2019 ACM/SIGDA International Symposium on Field-Programmable Gate Arrays (FPGA'19), 2019.
- [25] S. Wang and Y. Liang, "A comprehensive framework for synthesizing stencil algorithms on FPGAs using OpenCL model," in *Proceedings of* the 54th Annual Design Automation Conference (DAC'17), 2017.
- [26] J. Smagorinsky, "General circulation experiments with the primitive equations," *Monthly Weather Review, Vol. 91, No. 3*, 1963.
- [27] S. Williams, A. Waterman, and D. Patterson, "Roofline: an insightful visual performance model for multicore architectures," *Communications* of the ACM, vol. 52, no. 4, pp. 65–76, 2009.
- [28] C. Osuna, T. Wicky, F. Thuering, T. Hoefler, and O. Fuhrer, "Dawn: a high-level domain-specific language compiler toolchain for weather and climate applications," *Supercomputing Frontiers and Innovations*, vol. 7, no. 2, 2020.
- [29] Intel introduces world's largest FPGA with 43.3 billion transistors. Published on Tom's Hardware, November 2019. Accessed October 27, 2020. [Online]. Available: https://www.tomshardware.com/news/ intel-introduces-worlds-largest-fpga-with-433-billion-transistors
- [30] J. Hegarty, J. Brunhaver, Z. DeVito, J. Ragan-Kelley, N. Cohen, S. Bell, A. Vasilyev, M. Horowitz, and P. Hanrahan, "Darkroom: compiling highlevel image processing code into hardware pipelines," ACM Transactions on Graphics (TOG), vol. 33, no. 4, 2014.
- [31] T. Gysi, C. Osuna, O. Fuhrer, M. Bianco, and T. C. Schulthess, "STELLA: A domain-specific tool for structured grid methods in weather and climate models," in *Proceedings of the International Conference for High Performance Computing, Networking, Storage and Analysis (SC'15)*, 2015.
- [32] B. Hagedorn, L. Stoltzfus, M. Steuwer, S. Gorlatch, and C. Dubach, "High performance stencil code generation with Lift," in *Proceedings of* the 2018 International Symposium on Code Generation and Optimization (CGO'18), 2018.
- [33] G. Singh, D. Diamantopoulos, C. Hagleitner, S. Stuijk, and H. Corporaal, "NARMADA: Near-memory horizontal diffusion accelerator for scalable stencil computations," in 2019 29th International Conference on Field Programmable Logic and Applications (FPL'19), 2019.
- [34] G. Singh, D. Diamantopoulos, C. Hagleitner, J. Gómez-Luna, S. Stuijk, O. Mutlu, and H. Corporaal, "NERO: A near high-bandwidth memory stencil acceleratorfor weather prediction modeling," in 2020 30th International Conference on Field Programmable Logic and Applications (FPL'20), 2020.