

Domain-Specific Optimization and Generation of High-Performance GPU Code for Stencil Computations

Prashant Singh Rawat[†], Miheer Vaidya[†], Aravind Sukumaran-Rajam[†], Mahesh Ravishankar[§], Vinod Grover[§],
Atanas Rountev[†], Louis-Noël Pouchet[‡], P. Sadayappan[†]

[†] The Ohio State University

{rawat.15, vaidya.56, sukumaranrajam.1, rountev.1, sadayappan.1}@osu.edu

[§] NVIDIA Corporation

{mravishankar, vgrover}@nvidia.com

[‡] Colorado State University

pouchet@colostate.edu

Abstract—Stencil computations arise in a number of computational domains. They exhibit significant data parallelism and are thus well suited for execution on Graphical Processing Units (GPUs), but can be memory-bandwidth limited unless temporal locality is utilized via tiling. This paper describes how effective tiled code can be generated for GPUs from a DSL for stencils. Experimental results demonstrate the benefits of such a domain-specific optimization approach over state-of-the-art general-purpose compiler optimizations.

Index Terms—Stencil computations, tiling, Resource optimization, streaming, GPGPU

I. INTRODUCTION

Stencil computations arise in scientific applications in many domains. Due to the large number of variants of the stencil operators that arise in applications, it is not feasible to construct efficient libraries for stencil computations. Hence, stencils have been the focus of a number of domain-specific languages and frameworks [1]–[15]. GPUs (Graphical Processing Units) are an attractive architectural target for stencil computations because of the high degree of data parallelism. However, achievable performance for many stencil computations is constrained by bandwidth to global memory. If the number of floating-point operations per point is not sufficiently high, reuse across multiple time steps for an iterated stencil or across a sequence of stencils is essential to achieve high performance.

Tiling is a fundamental technique for data locality enhancement, and can enable significant reduction in the amount of data movement from/to global memory. However, the nature of data dependences that arise with stencils does not permit simple rectangular tiling across multiple time steps of an iterative stencil computation. A number of research efforts have therefore addressed the topic of effective tiling of stencil computations [2]–[8], [12]–[14], [16]–[29].

In this paper, we present a detailed analysis of constraints in achieving high performance with stencils on GPUs and describe domain-specific and GPU-target-specific optimization strategies to generate high-performance GPU code for stencils.

While the GPU algorithms discussed in this paper use CUDA terminology and the implementations use NVIDIA GPUs as reference hardware, the algorithmic strategies are not limited to CUDA or NVIDIA GPUs. These optimizations can be implemented using any language that allows mapping computations to GPU threads, e.g., OpenCL, and are applicable to GPUs from other vendors, e.g., AMD.

Several considerations are important in achieving high performance for any computation on GPUs: movement of contiguous chunks of data (coalesced accesses) from/to global memory; reuse of data in registers and on-chip scratch-pad memory (shared memory in CUDA terminology; called local memory in OpenCL); sufficient concurrency, orders of magnitude larger than the number of physical cores; and minimization of control-flow divergence among threads. The capacity of on-chip scratch-pad memory in an SM (Streaming Multiprocessor) is quite low, typically under 100 Kbytes. As elaborated later, the low shared memory capacity limits the maximum tile size. This further limits both the amount of reuse achievable for data in shared memory, and the maximum number of simultaneously active thread blocks in an SM. Low occupancy can have a deleterious effect on performance, because the number of instructions across the active warps that can be issued without stalling on dependent data might be insufficient to effectively overlap and mask the high global memory access latency. The total capacity in the register file per SM on GPUs is almost an order of magnitude larger than shared memory capacity. Further, access latency from registers is lower than shared memory. In this paper, we present details on a code generation strategy for efficient execution of stencils on GPUs, which effectively utilizes available register and shared memory resources to achieve high data reuse and thereby realize high performance.

The paper is organized as follows. Sec. II provides an overview of the main issues in achieving high performance with stencil computations. Sec. III describes a Domain-Specific Language (DSL) for specifying stencil computations: the STENCILGEN language. Sec. IV details the impact of

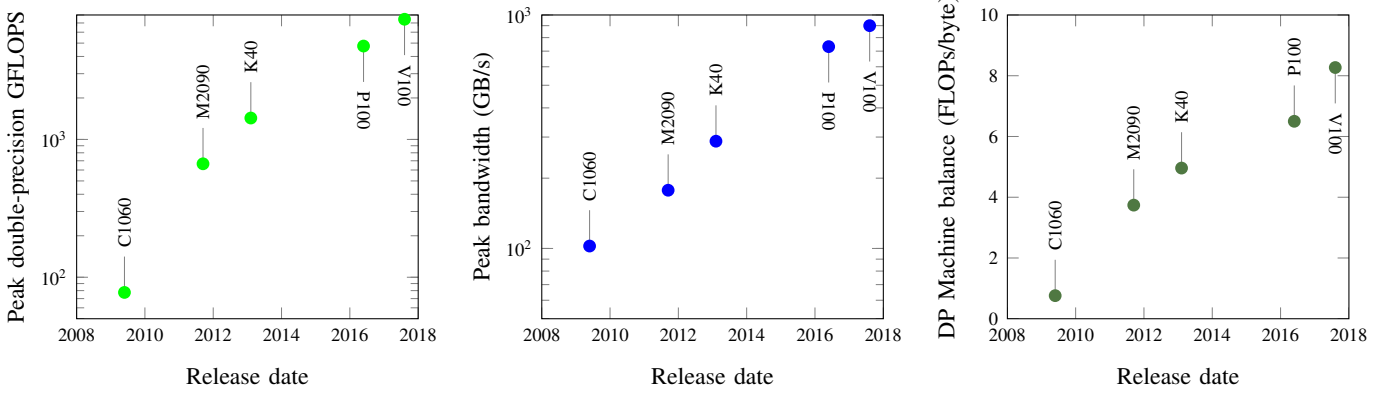


Fig. 1: Theoretical peak double-precision compute performance, machine balance and theoretical peak bandwidth for NVIDIA devices over the last decade

the GPU hardware constraints on the stencil optimization strategy. GPU code generation algorithms are described in Sec. V. Sec. VI discusses fusion across multiple stencils. Sec. VIII presents an experimental evaluation of the described stencil code generation approach, comparing it with several current general-purpose and special-purpose code generators and optimizers for multicore processors, manycore processors, and GPUs. Sec. VII discusses related work, and Sec. IX presents our conclusions.

II. TILING OF STENCIL COMPUTATIONS

In this section, we review the nature of the data dependences and potential for data reuse in executing stencil computations. Many stencil computations are fundamentally limited by memory bandwidth unless temporal locality is exploited across a sequence of stencils in a processing pipeline, or across repeated application of a stencil in an iterated computation.

A. Machine Balance and Operational Intensity

Fig. 1 provides data on the peak double-precision floating-point performance (Fig. 1a) and peak global-memory bandwidth (Fig. 1b) for GPUs from NVIDIA. Five generations of GPUs are shown: pre-Fermi (C1060), Fermi (M2090), Kepler (K40), Pascal (P100), and Volta (V100). The peak machine double-precision floating-point performance ($GFLOPS_{mc}$) has increased by over two orders of magnitude, from under 100 GFLOPS for the C1060 to around 10 TFLOPS for the Volta V100. The peak global memory bandwidth (BW_{mc}) has also increased across successive GPU generations, but not to the same extent: from 100 Gbytes/sec to nearly 1000 Gbytes/sec. Therefore, as seen in Fig. 1c, there has been a rise in the machine balance parameter, $MB = \frac{GFLOPS_{mc}}{BW_{mc}}$, from about 0.5 FLOPs/byte for the C1060 to around 8 FLOPs/byte for the Volta V100. As elaborated below, in order to achieve close to peak performance, the operational intensity (OI) of a computation (ratio of the number of arithmetic/logic operations to the number of bytes of data movement from/to memory) must be higher than the machine balance parameter.

```
for(i=1; i<N-1; i++)
  b[i] = c1*(a[i-1]+a[i]+a[i+1]);
```

(a) 1D stencil

```
for(i=1; i<N-1; i++)
  for(j=1; j<N-1; j++)
    b[i][j] = c2*(a[i-1][j]+a[i+1][j]+
      a[i][j-1]+a[i][j+1]+
      a[i][j]);
```

(b) 2D stencil

```
for(i=1; i<N-1; i++)
  for(j=1; j<N-1; j++)
    for(k=1; k<N-1; k++)
      b[i][j][k]=c3*(a[i-1][j][k]+a[i+1][j][k]+
        a[i][j-1][k]+a[i][j+1][k]+
        a[i][j][k-1]+a[i][j][k+1]+
        a[i][j][k]);
```

(c) 3D stencil

Fig. 2: Examples of 1D, 2D and 3D stencils

Consider the examples of 1D, 2D, and 3D stencils in Fig. 2. For the 1D stencil, $3(N-2)$ floating-point operations are executed and a total number of $2N-2$ distinct array elements are accessed (N elements of a and $N-2$ elements of b). If full reuse in registers/cache is achieved for each accessed array element, the achieved operational intensity would be $\frac{3(N-2)}{8(2N-2)}$, i.e., approximately $\frac{3}{16}$ FLOPs/byte for double-precision computation (8 bytes per data element). For the 1D stencil computation, the maximal possible operational intensity of $\frac{3}{16}$ FLOPs/byte is considerably lower than the machine balance parameter of modern GPUs. If an implementation of some algorithm has an upper bound of OI_{alg} for achievable operational intensity, and the peak memory bandwidth of the machine is BW_{mc} , the maximum achievable computational performance is $OI_{alg} \times BW_{mc}$. Equivalently, the maximum fraction of the machine's peak performance $GFLOPS_{mc}$ that an implementation of an algorithm can achieve is $\frac{OI_{alg}}{MB}$.

As the dimensionality of the grid increases, the maximum achievable operational intensity also increases. For the 2D stencil in Fig. 2(b), the number of floating-point operations is $5(N-2)^2$, and the number of accessed array elements

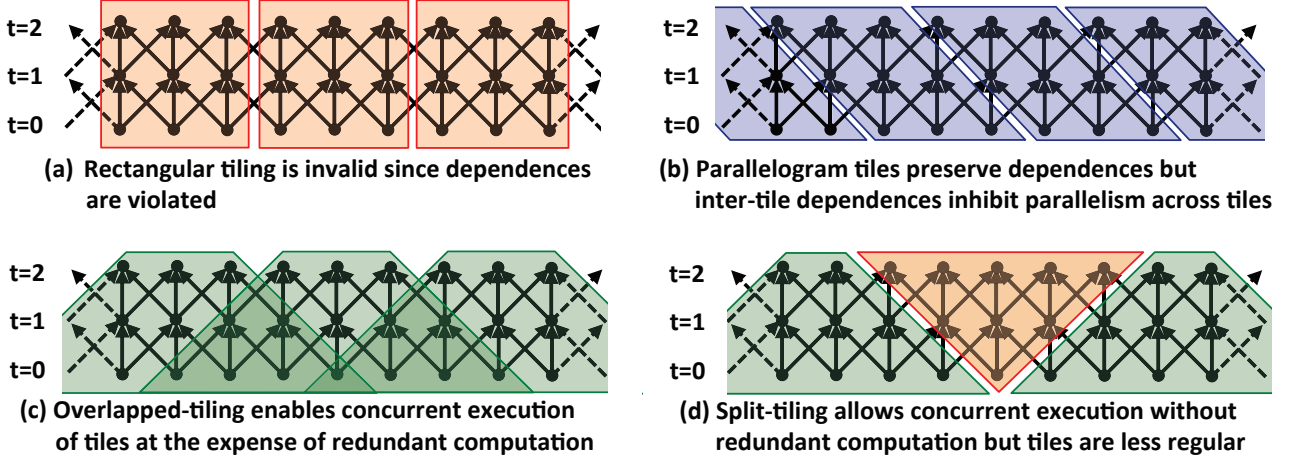


Fig. 3: Time-tiling to enhance data reuse with stencil computations

is $N^2 + (N - 2)^2$, with an OI upper bound of around $\frac{5}{16}$ FLOPs/byte for double-precision computation. For the 3D stencil Fig. 2(c), the OI upper bound rises to $\frac{7}{16}$ FLOPs/byte, but is still an order of magnitude lower than the machine balance parameter for current GPUs.

B. Time-tiling of Stencils

In order to increase performance with such stencils, it is essential to exploit reuse across time steps in an iterated stencil computation. With time-iterated stencils, the output array produced by the application of the stencil operation on the input array is in turn used as the input array for the next application of the stencil. Fig. 3 illustrates a 1D 3-point stencil iterated over three time steps. Since the arrays subject to the stencil computation are typically much larger than cache, reuse of data across the iterated application cannot be achieved without tiling, i.e. computing for a number of time steps for a subset of the array elements within a tile before accessing other array elements for the first time step. However, as seen in Fig. 3(a), a simple rectangular 2D tiling strategy is not feasible because dependencies would be violated if all iteration space points within one tile were contiguously executed.

Fig. 3(b) shows a valid tile shape for the 1D stencil example. All data dependencies to a tile are from iteration-space points in a tile to the left. As long as the tiles are executed from left to right, all dependencies will be preserved. However, the serializing inter-tile dependencies are undesirable. Fig. 3(c) shows overlapped trapezoidal tiles. The iteration-space points within a tile depend only on other points included within the same tile. Thus, concurrent execution of tiles is feasible. However, some redundant computation is necessary, as can be seen in Fig. 3(c) with the iteration-space points that belong in multiple tiles. Fig. 3(d) shows another strategy called split-tiling that also enables concurrent tile execution. With this strategy, a tile belongs in one of multiple groups (the number of groups grows with the dimensionality of the arrays and stencil operator). Tiles within the same group can execute concurrently, and the different groups are scheduled in order. In Fig. 3(d), all green-colored tiles execute first,

followed by the orange-colored tiles. Split-tiling has no redundant computational overhead. However, the tiles exhibit greater irregularity than with overlapped tiling, making it more challenging to achieve high performance on GPUs. In this paper, we present details on a stencil optimizer for GPUs that uses the overlapped tiling strategy.

C. Combining Overlapped Tiling with Streaming

The inter-tile concurrency enabled by overlapped-tiling comes at the price of redundant computation (and redundant data movement). While the fraction of redundant computation can be made very small for 1D overlapped tiling, the overhead can be very significant for 3D stencils. Consider a thread block of $8 \times 8 \times 8$ threads that computes two time steps of a 3D order-1 stencil. The thread block generates output after two time steps for $8 \times 8 \times 8$ elements, which means that the output after the first time step must be generated for a domain of size $10 \times 10 \times 10$, which in turn requires that input over a $12 \times 12 \times 12$ domain must be read. Thus, 12^3 data elements must be read in to produce results for 8^3 elements, almost half of them being also read by neighboring thread blocks. The amount of computation for the intermediate time step is proportional to $10^3 = 1000$, out of which $1000 - 512 = 488$ are redundant. This problem of redundant computation overhead for 3D stencils gets worse with an increase in stencil order and time tile size. Increasing the tile size can lower the overhead, but the maximum limits of thread-block size and shared-memory on an SM prohibit tiles much larger than $8 \times 8 \times 8$.

Since overlapped tiling for all the dimensions of a 3D stencil incurs high overhead on GPUs, an alternative is to use overlapped tiling along two of the three spatial dimensions, and sequentially *stream* along the third. We first describe overlapped tiling with streaming for a 2D stencil. Fig. 4a shows a 2D stencil where overlapped tiling is applied along both x and y dimensions with 2D blocks; the figure only shows overlap along the x dimension. Next, Fig. 4b shows streamed execution along the y dimension with 1D block. We can imagine the 2D domain sliced into 1D lines along the streaming dimension. We observe that due to the dependence

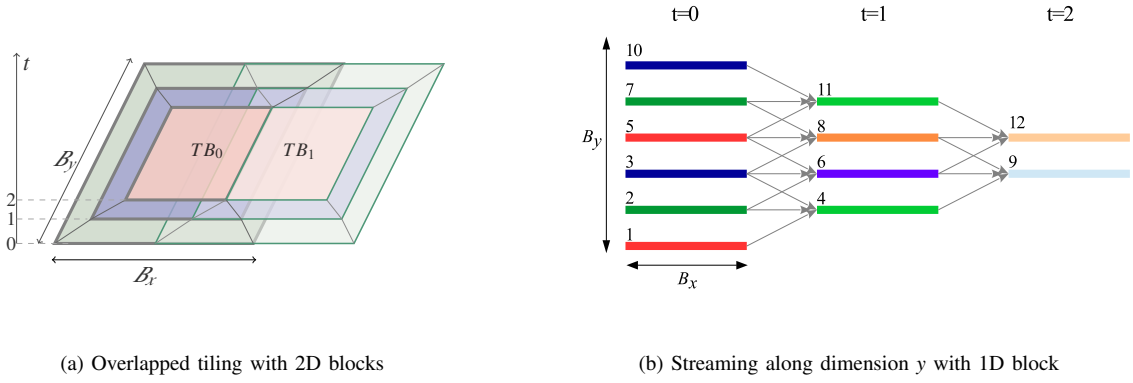


Fig. 4: Overlapped tiling and Streaming for 2D stencil

pattern, only three consecutive lines need to be read from the input domain to compute one line of the output domain. In general, an order- k 2D stencil will need to read $2k + 1$ consecutive input lines to compute one output line. The lines are cached in three distinct shared memory buffers per time steps, represented by different colors in Fig. 4b. The number over the buffers indicates the logical timestamp (LT) at which the particular buffer is populated.

In the prologue, three input lines are loaded into the buffers at $t = 0$, so that an intermediate result can be computed at $t = 1$ (LT = 4). After LT = 4, the data held in the buffer populated at LT = 1 is no longer needed, and this buffer can be subsequently used to cache a new input line. After this initial prologue, we can compute one intermediate result at $t = 1$ for every new input line read at $t = 0$. When we have three intermediate result lines available at $t = 1$, we can compute one line of the final output at $t = 2$ at LT = 9. In the steady state, we compute one output line at $t = 2$ by reading one input line at $t = 0$, and computing one intermediate line at $t = 1$.

This scheme requires only three buffers per time step, irrespective of the number of lines. In general, for an order- k stencil, $2k + 1$ buffers will be required per time step. The buffer that holds the oldest value at LT = r can be reused to store the newest value at LT > r . If the different lines along y are computed sequentially, one after the other, we refer to it as *serial-streaming*. However, the iteration space could also be partitioned along the y axis, with the different partitions being assigned across thread blocks. Each thread block serially streams through its portion of the iteration space, but different blocks can concurrently stream through their partition of iteration space. We refer to this as *concurrent-streaming*.

Streaming can be extended to 3D Jacobi stencils by interpreting the line in Fig. 4b as a plane. Since streaming eliminates redundant computations along the streaming dimension, and reduces the amount of shared memory, it is attractive to combine it with overlapped tiling. For a 3D stencil, overlapped tiling along two of the dimensions and streaming along the remaining dimension enables sufficient inter-tile concurrency, and significantly reduces the amount of redundant operations (computations and global-memory loads) since the required shared memory capacity per thread block corresponds to a 2D slice and not a 3D slice.

If the time tile size is 1, only spatial reuse is exploited. Such spatial blocking, also known as 2.5-D blocking, was used by Micikevicius [17] to optimize 3D finite difference computation. Nguyen et al. [16] used the corresponding time tiled blocking, also known as 3.5-D blocking, to accelerate LBM and 7-point 3D Jacobi stencil on multi-core CPU and GPU. Sec. V describes a systematic approach to determine the best tiling strategy for a given problem.

III. OVERVIEW OF A STENCIL DSL

To allow specification of stencil computations, we define STENCILGEN, a domain-specific language used as input to the code generation techniques described in subsequent sections. The use of a DSL enables the easy identification of the stencil patterns so that stencil-specific optimizations may be performed. A code generator can easily determine the high-level semantic properties of a DSL code region. In principle, such properties could also be inferred from an equivalent code region written in a general-purpose language. However, in such a scenario it is necessary that a variety of precise compiler analyses establish important properties such as lack of aliasing, values of induction variables, etc. The use of a domain-specific language helps avoid these problems. The STENCILGEN language describes:

- the point operation for each data element, as well as boundary conditions and time iterations
- the region of data on which a stencil is applied

Note that general-purpose programming is allowed outside the stencil regions, i.e., STENCILGEN can be used both as a stand-alone and as an embedded DSL.

The code example in Listing 1 highlights some of the main features of STENCILGEN. The read-only integer parameters M and N are used to describe the dimensions of the input and output arrays (line 3). Line 2 declares iterators, each of which will be mapped to a unique dimension of the computational loop nest. The iterators must be immutable within the body of the loop nest, and are assumed to be incremented in unit steps by the increment condition of the loop. All declared arrays and scalars (e.g., `in` and `out` in line 3) will be passed as arguments to the host function that will be generated from the STENCILGEN input. Line 4 specifies the arrays and scalars that need to be copied from host to device.

Listing 1: A representative stencil in STENCILGEN language

```

1 parameter M, N;
2 iterator i, j;
3 double in[M,N], out[M,N];
4 copy-in in;
5 stencil five_point_avg (A, B) {
6     double ONE_FIFTH = 0.2;
7     B[i][j] = ONE_FIFTH*(A[i-1][j]
8         + A[i][j-1]+A[i][j]+A[i][j+1]+A[i+1][j]);
9 }
10 stencil boundary (A, B) {
11     B[i][j] = A[i][j];
12 }
13 iterate (1:10) {
14     [0 : M-1][0 : N-1] : five_point_avg (in, out);
15     [1 : M-2][1 : N-2] : five_point_avg (out, in);
16 }
17 [0 : 1 ][0 : N-1] : boundary (in, out);
18 [M-2 : M-1][0 : N-1] : boundary (in, out);
19 [0 : M-1][0 : 1 ] : boundary (in, out);
20 [0 : M-1][N-2 : N-1] : boundary (in, out);
21 copy-out out;

```

Lines 5–9 define a stencil `five_point_avg`. The stencil definition has as arguments the input and output arrays/scalars used in the computation. This stencil averages the 5 neighboring data elements from array `A` and writes the result to an element of array `B`. All memory accesses in the stencil function must be scalars or array elements. Lines 10–12 define a stencil `boundary` that performs point-wise copy of elements from array `A` to array `B` at the boundary. Lines 13–20 invoke the stencil functions over subsets of the problem domain. The `iterate` construct at line 13 defines a time loop over the stencil calls. We explicitly unroll the call to stencil `five_point_avg` by a minimum number of iterations to obviate the need for exchanging the input and output after each time step. For such time-iterated stencils, the STENCILGEN-based code generator automatically determines the optimal degree of unrolling for the time loop, fuses the execution of two or more time steps, an optimization termed *time-tiling*, to exploit the temporal reuse exposed by the producer-consumer relationship between the time steps. Finally, line 21 defines the arrays and scalars that need to be copied back from device to host.

Embedding STENCILGEN in C/C++ code: Every embedded STENCILGEN region is delineated by `#pragma` annotations. The starting `#pragma` specifies additional code generation parameters, including the size of the GPU thread block, the dimension along which spatial streaming is to be performed, and the size of shared memory available on the target GPU.

When embedding STENCILGEN regions, the data arrays are assumed to be allocated by the outside code, and to be stored contiguously in the C/C++ memory space. The right-hand-side expressions of the assignment statements in the stencil body are assumed to be side-effect-free. Those expressions are allowed to contain standard side-effect-free math functions such as `sin`, `sqrt`, etc.

Additional details about STENCILGEN, including its grammar, are available in the first author’s Ph.D. dissertation [30].

IV. GPU CONSTRAINTS FOR STENCIL COMPUTATION

This section discusses architectural constraints to be considered when optimizing stencil computations on GPUs. The basic computational unit of a GPU is a thread. Going up the hierarchy, threads are grouped into a warp, warps are grouped into thread blocks, and thread blocks are grouped into a grid. The number of threads per warp is architecture-specific; 32 and 64 threads form a warp in NVIDIA and AMD devices, respectively. The GPU host code explicitly specifies the grid and thread block dimensions when the kernel is launched. All threads in a warp execute instructions in a synchronized, lock-step fashion. Warps within a thread block are mapped to the same *streaming multiprocessor* (SM), and can synchronize among themselves using synchronization primitives. Each thread has a fixed number of registers that can be varied at compile time. Threads in a thread-block can exchange data via shared memory. In recent GPU architectures, threads within the same warp can also exchange data held in registers via shuffle intrinsics.

Efficient tiling schemes for GPUs must

- perform coalesced access to global memory;
- have sufficient parallelism to tolerate memory access latencies; and
- judiciously use faster storage, like shared memory and registers, for caching data.

Shared memory has lower access latency than global memory, and is well suited to cache data that is accessed by multiple threads in a thread block. Registers are local to a thread, and therefore well suited to cache data that is accessed by a single thread. Registers are the fastest, and are more plentiful than shared memory in most GPU architectures. Therefore, it is often beneficial to offload some storage from shared memory to registers. However, excessive register usage may result in lower occupancy or expensive register spills.

In this section, we discuss GPU resource considerations in determining size/shape of thread blocks and grid for spatial tiling of stencil computations.

A. Constraints on thread block and grid size

GPUs have hardware limits on i) the maximum number of concurrently loaded threads per SM (T_{sm}), ii) the maximum number of threads in a thread-block (T_b), iii) total shared memory per SM (M_{sm}), iv) the maximum number of concurrently loaded thread blocks per SM (B_{sm}), and v) the register file size per SM (R_{sm}). For instance, $T_{sm} = 2^{11}$, $T_b = 2^{10}$, $M_{sm} = 48\text{KB}$, and $B_{sm} = 16$, and $R_{sm} = 2^{16}$ for an NVIDIA Tesla K20c GPU. The threads in an SM can be grouped in various ways, e.g., 2 blocks of 1024 threads, 16 blocks of 128 threads, etc.

The typical approach to data-parallel execution on GPUs is to assign one thread for computation of one element of the result array. In order to tolerate the very high latency to global memory (several hundred clock cycles), massive parallelism must be utilized. If a sufficient number of independent instructions are ready to be issued amongst the collection of active warps in an SM, the memory access latency can be fully overlapped. But even if memory access latency can be fully overlapped, the peak memory bandwidth can limit

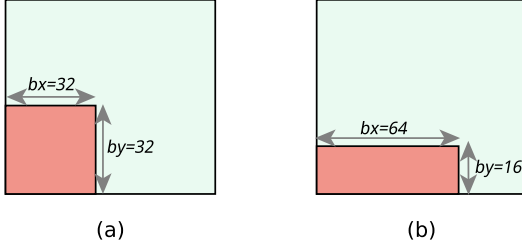


Fig. 5: Different thread block configurations for a fixed thread block size of 1024

performance, as discussed earlier. The only way to overcome the performance limitation from global memory bandwidth is to exploit temporal reuse on the accessed data, by buffering it in faster shared memory or registers. However, the shared memory and/or registers used for such buffering may result in a reduction in warp occupancy because fewer thread-blocks may now be concurrently schedulable on an SM due to the register/shared memory usage per thread-block.

Shared memory allows significantly faster access than global memory, and unlike registers, allows sharing of data across threads in a thread block. However, the more shared memory a thread block uses, the fewer thread blocks can be simultaneously scheduled on an SM. If m_b is the bytes of shared memory used by a thread block, the maximum number of concurrently active thread blocks cannot exceed $\frac{M_{sm}}{m_b}$.

Registers are the fastest storage resource available to a thread. The maximum number of registers per thread can be controlled via compiler flags for GPUs. If the maximum possible number of threads are to be active on an SM, the number of registers used per thread must be $\leq \frac{R_{sm}}{T_{sm}}$. If a thread uses t_{reg} registers, the maximum number of active threads per SM can be no more than $\min\left(T_{sm}, \frac{R_{sm}}{t_{reg}}\right)$.

If the size of a thread block (number of threads) is sz_b , the maximum number of concurrently schedulable thread blocks per SM, max_b , is upper-bounded by $max_b \leq \min\left(B_{sm}, \frac{T_{sm}}{sz_b}, \frac{M_{sm}}{m_b}, \frac{R_{sm}}{t_{reg} \cdot sz_b}\right)$. For bandwidth-bound computations, unless thread coarsening is utilized to increase instruction level parallelism (ILP), the block size sz_b is often chosen to maximize occupancy, and consequently the thread level parallelism (TLP), i.e., $sz_b \times max_b$ is set as close to T_{sm} as possible.

Coalesced access to global memory is important since it minimizes the number of memory transactions for accessing a given number of data elements. This requires that the fastest varying dimension of a multi-dimensional thread block be aligned with respect to data access with the fastest varying dimension of the accessed multi-dimensional array. The fastest varying dimension of the thread block is usually chosen to be a multiple (or factor) of the warp size to benefit from coalescence.

B. Partitioning threads across thread block dimensions

For a 2D thread block of size sz_b , many choices exist for the extents along the x and y dimensions. The performance of

a stencil computation can be quite sensitive to this choice. Let $sz_b = b_x \times b_y$. Two possible configurations for $sz_b = 1024$ are shown in Fig. 5. For an order- k stencil over an N^2 domain, the number of global memory loads is $N^2(32+2k)^2/1024$ for configuration (a) and $N^2(64+2k)(16+2k)/1024$ for configuration (b). For any value of k , configuration (a) requires fewer global load transactions for performing the same number of arithmetic operations. Consider, for example, a 3D single-precision 7-point Jacobi stencil executed on a GP100 device. The version with block size 32×32 incurs $5.91E+07$ global load transactions, achieving 1.2 TFLOPS. In contrast, the version with block size 64×16 incurs much higher $7.59E+07$ global load transactions, achieving only 989 GFLOPS.

V. RESOURCE OPTIMIZATION STRATEGIES

In this section, we discuss various factors that affect the best choice of tiling strategy for a given stencil computation. Sec. V-A describes considerations for the 2D case, for streamed overlapped tiling. Sec. V-B characterizes 3D stencils based on their access patterns, and describes optimization techniques for them. We consider a stencil of order $k = 1$, with time tile size $T = 4$, and all operations in single precision. Tiled code must be appropriately tuned for different GPUs; for the discussion below, we use hardware parameters of an NVIDIA Tesla K20c as the target GPU.

A. Overlapped Tiling for 2D Stencils

Many prior code generators that optimize 2D stencils perform overlapped tiling along both the dimensions [3], [4]. In contrast, we partition the N^2 input domain into overlapping strips of size $B_x \times N$: we stream along the y^1 dimension to eliminate redundant computations, performing overlapped tiling along x for inter-tile concurrency. Fig. 4a illustrates this strategy. The required amount of per-thread-block shared memory and registers depend on the thread block geometry.

1) *Streaming with 2D thread blocks*: In the steady state, at each step, a $B_x \times B_y$ thread-block reads B_y input lines along the y dimension to compute B_y points of the iteration space. Each thread traverses the iteration space at a stride of B_y along the y dimension. Each thread block stores $T(B_y + 2k)$ lines in shared memory. Consider two possible choices for B_y :

- $B_y = 1$ (1D block): For maximum occupancy, each thread block must have $\frac{2048}{16} = 128$ threads. In each iteration, a thread block loads one line from input at $t = 0$ and generates one line of output at $t = 4$. For this, it needs 6KB of shared memory. From Sec. IV-A, it follows that an SM can have at most $\min\left(16, \frac{2048}{128}, \frac{48KB}{6KB}\right) = 8$ blocks, a 50% loss of occupancy.
- $B_y = 32$ (2D block): Two thread blocks of size 32×32 can theoretically be active per SM. Since a thread block now operates on a chunk of 32 lines instead of 1, it needs 16KB of shared memory. Each SM can have $\min\left(16, \frac{2048}{1024}, \frac{48KB}{16KB}\right) = 2$ active blocks, which implies maximum occupancy.

¹Streaming along the x dimension is usually not beneficial because it entails non coalesced accesses while loading the input data.

Clearly, better occupancy is achieved with 2D blocks. For a $B_x \times B_y$ block, the computation proceeds as shown in Fig. 4b: in the steady state, $B_y - 2k$ output lines are generated at time step t , using B_y lines from the shared memory buffer corresponding to time step $t - 1$. With a *sliding-window* approach, the B_y oldest lines in the buffer can be reused to cache the new lines. The buffer can be implemented as a circular array, with modulo operations to find the top and bottom k rows. Since modulo operations are costly on GPUs as they compile to multiple instructions, it is beneficial to make B_y a power of 2, so that they can be implemented by bitwise operators² which has a very high throughput.

2) *Concurrent streaming with 1D thread block and using registers for storage*: As demonstrated above, serial streaming with 1D thread blocks suffered lowered occupancy due to high usage of shared memory. If z is the streaming dimension, then a stencil is termed *diagonal-access free* along z dimension if all stencil offsets (x_0, y_0, z_0) for access to points on different planes along the z dimension are strictly of the form $(0, 0, z_0)$. For such diagonal-access-free stencils, the shared memory requirement can be decreased by using registers instead of shared-memory to hold the $2k$ accesses to lines [17]. For a thread block of 128 threads, $k = 1$ and $T = 4$, a thread block now needs 8 registers and 2KB bytes of shared memory. Thus, the number of concurrently loadable thread blocks per SM is $\min\left(16, \frac{2048}{128}, \frac{48KB}{2KB}\right) = 16$. While occupancy is maximized, this strategy suffers from a low overall number of thread blocks. To keep all 13 SMs of the K20c GPU busy, at least $13 \times 16 = 208$ thread blocks are needed. However, even a large input domain of size 8192^2 can only be partitioned into $\frac{8192}{128} = 64$ blocks.

The number of thread blocks can be increased by using overlapped tiling along the y dimension as well. However, it is only necessary to partition the space into $\left\lceil \frac{208}{64} \right\rceil = 4$ blocks along the y dimension. Since the thread block is 1D and each tile in the input grid is 2D, streaming is performed within a tile.

This version does incur redundant computation and data access along the y dimension, but this is negligible compared to the reduction in redundant computation and data access along the x dimension due to increase in block size from 32 to 128. The lower access latency of registers, compared to shared memory is an additional benefit.

B. Overlapped Tiling for 3D Stencils

Fig. 6 depicts different possible choices of grid dimensionality for 3D stencils over a 3D domain: 3D, 2D, or 1D grid. A 3D grid implies a halo region on all six sides of a 3D interior data cube handled by a thread block, and suffers from significant redundant computation and data movement, as discussed earlier. At the other end, a 1D grid incurs the lowest overhead from redundant computation and data movement, but the total number of thread blocks may be insufficient to keep

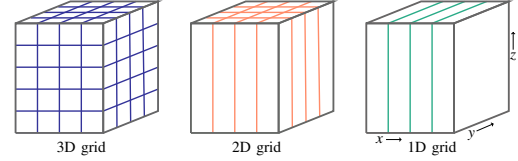


Fig. 6: Various grid dimensionalities for a 3D input

all SMs of the GPU busy. A 2D grid represents a good trade-off and offers moderate overhead from redundant computation and data movement along with adequate number of thread blocks for the SMs on the GPU.

Streamed time-tiling of a 3D stencil requires $T(2k + 1)$ planes to be cached in shared memory. For $T = 4$ and $k = 1$, a 32×32 thread block needs 48KB of shared memory, the maximum per SM on a K20 GPU. This means that only one thread block can be scheduled at a time on an SM, limiting occupancy to 50%. For $k = 2$, the required shared memory exceeds the hardware limit, forcing a smaller thread block size. Micikevicius [17] and Nguyen et al. [16] use registers to offload the caching of some planes from shared memory to registers for spatial and time tiling. Details regarding implementation of time tiling code with a combination of shared memory and registers are discussed next.

1) *Streaming and using registers for storage*: In the scenario above with $k = 1$ and block size 32×32 , if the stencil is diagonal-access free along the streaming dimension, then some storage can be offset to registers, increasing the per thread register pressure by $2Tk$, and simultaneously reducing the shared memory requirement to 16KB. With this trade off, an SM can have two active blocks, achieving maximum occupancy. If there are register spills due to the increased register pressure, then the value of T can be reduced to alleviate the register pressure.

Fig. 7 illustrates this scheme applied to time tile the 7-point 3D Jacobi stencil. Each output point at time step t needs to read data from three input planes at time step $t - 1$. The invariant maintained in this scheme is that the data needed to compute plane p'_v at time step t comes from registers that hold the values of planes p'_{v-1} and p'_{v+1} , and a shared memory buffer that holds the value of plane p_k^{v-1} at time step $t - 1$. This is illustrated in Fig. 7a, where the red (black) edges indicate that the data is being read from register (shared memory). The number on the edges indicates the sequence in which the contributions are made. Before the plane p'_{v+1} makes contribution to the plane p'_{v+1} , we perform a *shift*, where the data from the shared memory buffer holding the value of plane p'_v is moved into the register that held the value of plane p'_{v-1} , and the data from the register holding the value of p'_{v+1} is moved into the shared memory buffer. The shifts are represented as the blue edges in Fig. 7b. The shift helps maintain the invariant as the computation streams through the z dimension. In the steady state, we can compute one plane at each time step, as shown in Fig. 7c.

An implementation sketch of this scheme is presented in Algorithm 1. The algorithm is presented independently of the tiling scheme for the other two dimensions. The initializations

² $a \bmod b \equiv a \& (b-1)$ if b is a power of 2

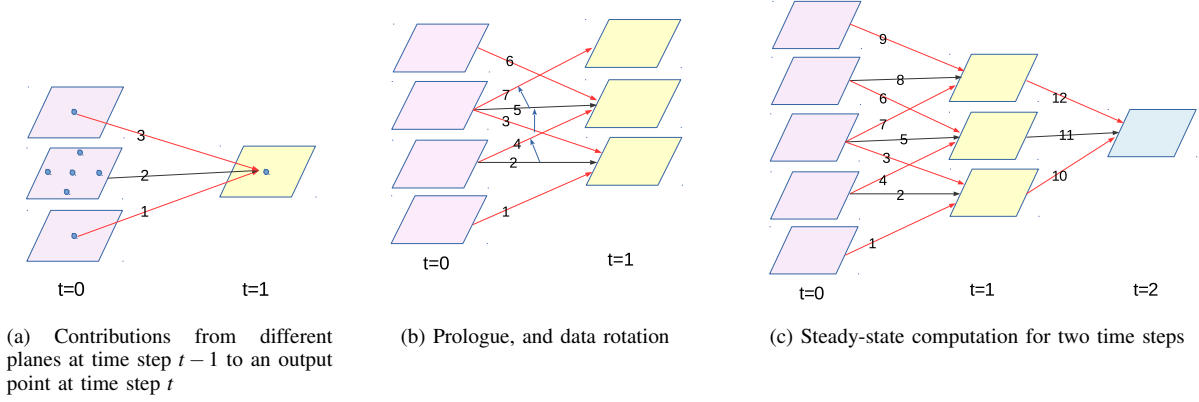


Fig. 7: Streaming for 2D overlapped blocks of a 3D 7-point Jacobi stencil. Red (black) arrows indicate that the value from the previous time step comes from register (shared memory).

Algorithm 1: Streaming and using registers for storage for a 3D order-1 stencil

Input : in : input array, T : time tile size
Output: out : output array

- 1 A_v^t : shared memory buffer for plane p_v at time step t ;
- 2 r_{v-1}^t, r_{v+1}^t : registers for planes p_{v+1}, p_{v-1} at time step t ;
// Initial reads
- 3 $r_{v-1}^0 \leftarrow load_plane(in[0][...][...]);$
- 4 $A_v^0 \leftarrow load_plane(in[1][...][...]);$
- 5 **for** each z from 1 to $N-1$ **do**
- 6 $r_{v+1}^0 \leftarrow load_plane(in[z+1][...][...]);$
- 7 $__syncthreads()$;
// Perform the computation per time step
- 8 **for** s from 1 to T **do**
- 9 $r_{v+1}^s = compute_stencil(r_{v-1}^{s-1}, A_v^{s-1}, r_{v+1}^{s-1});$
- 10 $__syncthreads()$;
- 11 $out[z-kT][...][...] = r_{v+1}^T$;
// Shift data per time step
- 12 **for** s from 0 to T **do**
- 13 $r_{v-1}^s \leftarrow A_v^s \leftarrow r_{v+1}^s$;

at line 1 depend on the tiling scheme used for the non-streaming dimensions. Algorithm 1 can be generalized to all stencils where each plane accesses only one point per plane from other planes along z dimension. If such a stencil accesses a diagonal point (x_0, y_0, z_0) , then appropriate shuffle intrinsic will have to be inserted in the code, so that each thread maintains a correct central value in registers.

2) *Storage optimization for associative stencils*: For stencils that access more than one point per plane from other planes along z dimension (i.e., stencils that are not diagonal-access free along the streaming dimension), the optimization scheme described above may require too many registers to be beneficial. For example, a 27-point 3D stencil with $T = 2$ will require 36 explicit registers per thread just for storage. If such a stencil is associative, we can use an optimization strategy that leverages contribution reordering to reduce the number of registers required at each time step to just $2k + 1$. A similar optimization strategy is used by Stock et al. [31] in the context

Listing 2: A 7-point 2D associative stencil that accesses 2 points from planes along y dimension

```

for (i=1; i<N-1; i++) {
  for (j=1; j<N-1; j++) {
    b[i][j] = c0*(a[i-1][j-2] + a[i-1][j+2]) +
              c1*(a[i][j-1] + a[i][j] + a[i][j+1]) +
              c2*(a[i+1][j-2] + a[i+1][j+2]);
  }
}

```

Listing 3: Reading one plane from input at a time, and accumulating its contribution at different output points

```

for (i=0; i<N; i++) {
  for (j=0; j<N; j++) {
    b[i+1][j] += c0*(a[i][j-2] + a[i][j+2]);
    b[i][j] += c1*(a[i][j-1] + a[i][j] + a[i][j+1]);
    b[i-1][j] += c2*(a[i][j-2] + a[i][j+2]);
  }
}

```

of high-order stencils on multi-core CPUs. Listing 2 shows an example of a stencil that is associative, and hence is amenable to our optimization. We leverage the associativity of addition and multiplication to rewrite it as an accumulation stencil of Listing 3.

Fig. 8 shows the application of this scheme to time tile the 7-point 3D Jacobi stencil. An input plane at time step $t-1$ contributes to three points belonging to distinct output planes at time step t . There are three invariants maintained: (1) the plane p_v^{t-1} making the contribution from time step $t-1$ is cached in shared memory; (2) contributions to the planes p_{v+1}^t, p_{v-1}^t and p_v^t at time step t are accumulated in registers; and (3) after the accumulation, the contribution to the plane p_{v-1}^t is written out into a shared memory buffer. This is illustrated in Fig. 8a, where the red (black) edges indicate that after accumulating the contribution, the final result is in register (shared memory). The number on the edges indicates the sequence in which the contributions are made. One can make the following observations from Fig. 8b: (a) At $t = 0$, all the contributions from a plane are made before starting with the next plane, and therefore one can reuse the same

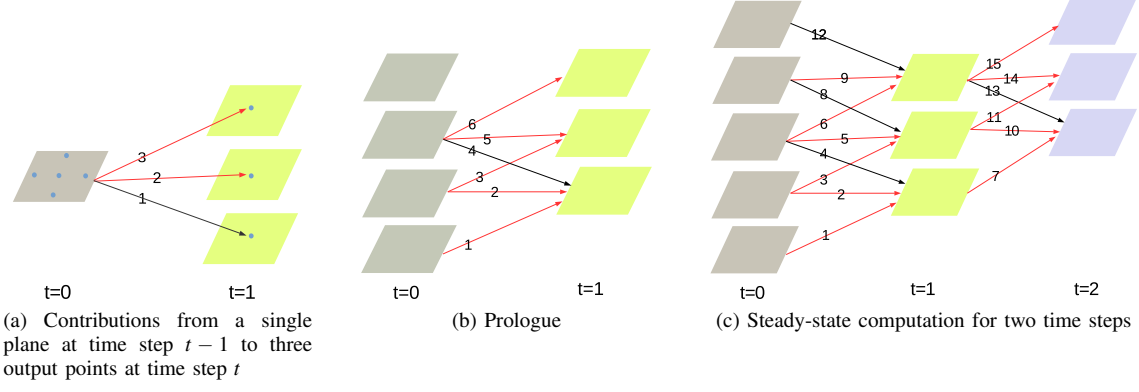


Fig. 8: Streaming with associative reordering for 2D overlapped blocks of a 3D 7-point Jacobi stencil. Red (black) arrows indicate that the contributions from the previous time step are accumulated and written into a register (shared memory buffer).

Algorithm 2: Storage optimization for an associative 3D order-1 stencil

Input : in : input array, T : time tile size
Output: out : output array

```

1  $A_{v-1}^t$  : shared memory for plane  $p_{v-1}$  at time step  $t$ ;
2  $r_{v-1}^t, r_v^t, r_{v+1}^t$  : registers for planes  $p_{v-1}, p_v, p_{v+1}$  at time
  step  $t$ ;
3 for each  $z$  from 0 to  $N-1$  do
4    $A_{v-1}^0 \leftarrow load\_plane(in[z][...][...]);$ 
5    $syncthreads();$ 
   // Perform the computation per time step
6   for  $s$  from 1 to  $T$  do
7      $r_{v+1}^s \leftarrow bottom\_plane\_contribution(A_{v-1}^{s-1});$ 
8      $r_v^s += mid\_plane\_contribution(A_{v-1}^{s-1});$ 
9      $r_{v-1}^s += top\_plane\_contribution(A_{v-1}^{s-1});$ 
10     $A_{v-1}^s = r_{v-1}^s;$ 
11     $syncthreads();$ 
12     $out[z-kT][...][...] = A_{v-1}^T;$ 
   // Shift data per time step
13    for  $s$  from 1 to  $T$  do
14       $r_{v-1}^s \leftarrow r_v^s \leftarrow r_{v+1}^s;$ 

```

shared memory buffer for all the planes; (b) At any time step, at most one shared memory plane is required, and at most $2k+1$ registers simultaneously receive contributions. Thus, we only need $2k+1$ registers, and one shared memory buffer per time step. In a steady state (Fig. 8c), a plane that is in shared memory at time step t can start making contributions to the accumulation registers at time step $t+1$. Algorithm 2 presents an implementation sketch for this approach.

VI. FUSION HEURISTICS

More complex stencil computations, as arising in image processing pipelines and multi-statement stencils, can be represented by a directed acyclic graph (DAG) where each vertex represents a single stencil operation identified by the stencil statements on a set of input domains [3], [32], and edges represent data produced/consumed by these operations. Stencils that are iterated a fixed number of times can also be

written in a DAG form, by unrolling the time loop to expose different nodes for different time steps.

Fusion of stencil operators in a DAG can be essential to improve reuse of data between operators. For example, fusing together nodes which read the same data can enable better temporal locality, and reduce memory traffic. But on the other hand, fusing nodes increase resource pressure (e.g., registers) as the kernel would implement multiple stencil operators with the same per-kernel resources budget. In this section we present a simple resource-aware greedy heuristic to fuse stencil operators in a DAG.

For the fusion of two nodes in the DAG to be valid, it must preserve data dependences. For Jacobi-like stencils in a DAG, such fusion is always valid if the fused node is atomic, i.e., there is no dependence cycle after fusion. There are many valid fusion schemes for the nodes in the DAG, but their performance may vary significantly depending on the profitability of fusion (gains in temporal locality and memory traffic versus reduction in available per-operator resources in the fused kernel). The space of all valid fusion structures can be very large [33], and exhaustively exploring this space to find the optimal solution can be prohibitively expensive. To avoid such exploration, we instead propose a greedy algorithm to determine which nodes shall be fused together, to optimize a dedicated objective function. Starting from the input DAG $D_g = (V, E)$, the objective is build the *fused* graph $D_f = (V_f, E_f)$ such that each $v_f \in V_f$ is a convex partition of nodes from V . A convex partition is an atomic “macro-node”, comprising nodes from V , such that there is no dependence path going out and back in the macro-node. To form a macro-node v_f , we fuse node(s) $v \in V$ following a profitability function. Eventually a single GPU kernel is generated for each node $v_f \in V_f$. Note that if a temporary array is produced in v_f by design, and all its uses are contained in v_f , then no global memory transaction is needed for such temporary array.

A. Roadmap of the Greedy Fusion Algorithm

The greedy resource-driven fusion algorithm consists of the following steps:

Step 1: For each stencil operator, compute the amount of shared memory and registers that will be needed to cache the data (Sec. VI-B).

Step 2: Identify pairs of stencil operators that can be fused without violating dependences (Sec. VI-C).

Step 3: For each pair of stencil operators, compute the shared memory and registers used when fusing them in a single kernel, based on the individual resource usage computed in Step 1 (Sec. VI-D).

Step 4: Based on the resource usage of the fused node computed in Step 3, create a profitability metric that encodes the impact of fusion on the GPU resources, and the data movement volume (Sec. VI-E).

Step 5: Define a custom sort to order the profitability metrics of various stencil operator pairs, and choose to fuse the operators of the most profitable pair into a macro-node (Sec. VI-F). Update the stencil DAG and the dependence graph, and repeat again from Step 1, until no more stencil operator pairs can be fused.

Once a fused stencil DAG D_f is produced by the algorithm, an optimized CUDA kernel can be generated for each node in D_f using the algorithms described in Sec. V.

B. Computing Resource Requirements of a Stencil Operator

Without loss of generality, we assume that streaming, as presented in prior sections, is done on the outermost dimension of a 3D domain. As we present in Sec. V, to determine the resource requirements of an order- k stencil operator along the streaming dimension, we must determine the storage properties for the $2k+1$ accessed planes. We proceed as follows.

- (a) If computing an output value at (z_0, y_0, x_0) only requires a single input value at (z_r, y_0, x_0) from an input plane z_r , then z_r is stored as an explicit register (that is, z_r has register storage type). If not, then z_r is cached into shared memory
- (b) If an output element is to be (over-)written again, e.g., because it is an accumulator, then it is stored in an explicit register. Conversely, the last assignment to a copy-out array uses global memory for its storage type.

An explicit register above is implemented with a (scalar) temporary variable, which contains an array element. The compiler shall place such temporary variables in registers. These explicit registers are distinguished from *implicit* registers that are used to store arrays elements and the intermediate results of computing expressions, as produced by the compiler during traditional register allocation. We use this distinction for explanation purpose only, as in the generated code the explicit registers are thread scalars.

Listing 4 shows two consecutive time steps of a 3D 7-point Jacobi stencil defined in the STENCILGEN language. From rule (a) above, the statement line 7 should use a shared memory buffer, storing the input plane $A[0,*,*]$, since five different values are read from that plane and contribute to different output points. From rule (b), three explicit registers must be used to store the output values written to $B[1,*,*]$, $B[0,*,*]$ and $B[-1,*,*]$ at lines 2, 3, and 5 respectively.

We represent the resource requirement of a stencil operator using three attributes: N_{reg} and N_{shm} capture how many explicit

Listing 4: Two time steps of a 7-point 3D Jacobi stencil

```

1 stencil j3d7pt (A, B, C) {
2   B[k+1][j][i] = (A[k][j][i]);
3   B[k][j][i] += (A[k][j-1][i] + A[k][j][i-1] + A[k][j][i+1]
4               + A[k][j+1][i] + A[k][j][i+1] + A[k][j+1][i]);
5   B[k-1][j][i] += (A[k][j][i]);
6
7   C[k][j][i] = (B[k-1][j][i]);
8   C[k-1][j][i] += (B[k-1][j-1][i] + B[k-1][j][i-1] +
9                   B[k-1][j][i+1] + B[k-1][j+1][i]);
10  C[k-2][j][i] += (B[k-1][j][i]);
11 }
```

registers and shared memory buffers are used by the operator, and $M_{acc \rightarrow res}$ is a map from each accessed data plane to the storage type used for that plane.

C. Identifying Fusion Candidates

The dependence graph $G_{dep} = (V, E_{dep})$ captures the dependences between stencil operators in the DAG. From G_{dep} , consider a transitive dependence graph $G_{trans} = (V, E_{trans})$, such that

$$s_i \rightarrow s_j \in E_{trans} \text{ iff } \{ \exists s_k \mid s_k \in V \wedge s_k \neq s_i \wedge s_k \neq s_j \wedge \\ \text{there exists a path from } s_i \text{ to } s_k \text{ (denoted as } s_i \rightsquigarrow s_k) \wedge \\ \text{there exists a path from } s_k \text{ to } s_j \text{ i.e., } s_k \rightsquigarrow s_j \}$$

G_{trans} is used to ensure the validity of fusion, i.e., that all macro-nodes generated by fusion are convex/atomic nodes. If $s_i \rightarrow s_j \in E_{trans}$, then fusing s_i and s_j would break convexity, unless all nodes along any path $s_i \rightsquigarrow s_j$ are also included in the fused node. This definition allows the fusion of two stencils in producer-consumer relationship, as long as there are no other operators along any path from s_i to s_j .

For each unordered pairs of stencil operators (s_i, s_j) such that $s_i \rightarrow s_j \notin E_{trans}$, we first compute the total GPU resources required if s_i and s_j are fused, and then compute a *fusion profitability* metric that will quantify the benefit of fusing s_i and s_j in terms of GPU resource usage and data movement reduction.

D. Computing the Post-Fusion Resource Map

We first determine the resource usage of the fused macro-node made of s_i and s_j . Let M_i and M_j be the maps from accessed array planes to GPU resources ($M_{acc \rightarrow res}$) for the operators s_i and s_j . Let M_{fused} be the resource map for the fused node. If there is no dependence between s_i and s_j , then M_{fused} is simply the union of the resource maps of the individual nodes, where the rules of union are as defined in this section.

But if there is a true dependence between s_i and s_j , to accurately compute resource usage it is necessary to first build the interleaving of the two operators domain that captures a fused execution schedule that preserve the dependence(s). Once this schedule is built, the resource map of the dependence sink may need adjustments. For example, in Listing 4, after the contribution at line 5 the values of plane $B[-1,*,*]$ can be used as input for the next time step (lines 7–10). There

is a true dependence from the write to plane $B[-1,*,*]$ at line 5 to the subsequent reads at lines 7–10. Observe that the accesses along the streaming dimension are shifted by the dependence distance in lines 7–10 to ensure that after fusing lines 2–10, dependences are preserved. It is sufficient to shift by the dependence distance along the streaming dimension.

Once the schedules are adjusted to represent a valid fused schedule, the set of planes accessed by s_i and s_j may intersect. For example, in Listing 4, $B[-1,*,*]$ is mapped to a register in line 5, and to shared memory in lines 8–9. In the presence of such storage conflict for a plane, we always choose the storage that is lower in the memory hierarchy, i.e., shared memory in this example.

E. Computing the Profitability Metric

Once the resource map for the fusion of s_i, s_j is computed, we determine the profitability of fusing these nodes using 5 metrics as follows.

D_m represents the savings, in terms of data movements, after fusing s_i and s_j . For every array leading to a true dependence between s_i and s_j , data movements are reduced after fusion, as this array does not require as much global memory transfers after fusion versus without fusion. To capture this, when D_m is incremented by 2 this models a saving in a store and then subsequent load of that array. There is no explicit data movement saving for arrays carrying write-after-read dependence. If there is a write-after-write (or read-after-read) dependence between s_i and s_j but no read-after-write dependence, then D_m is incremented by 1 to represent the reduction in multiple writes (reads) to (from) the same location in global memory.

S_{reg} represents the total number of accesses that are mapped to the same explicit registers in s_i and s_j after fusion, i.e., $S_{reg} = \text{num_registers}(M_i) + \text{num_registers}(M_j) - \text{num_registers}(M_{fused})$. Similarly, S_{shm} represents the total number of accesses that are mapped to the same shared memory buffer in s_i and s_j after fusion, i.e., $S_{shm} = \text{num_shared}(M_i) + \text{num_shared}(M_j) - \text{num_shared}(M_{fused})$.

Finally, T_{reg} represents the total number of explicit registers in the fused node, and T_{shm} the total amount of shared memory used.

F. Constructing the Objective Function

Once all metrics above has been computed for all valid stencil operator pairs, they are sorted based on their profitability. We use a customized sort operation \prec , which models a total order of the list of candidate fusion pairs. The set of rules that are used to order two pairs $c_i, c_j \in L_{tuple}$ is as follows:

- a. $(D_m)_{c_i} < (D_m)_{c_j} \Rightarrow c_i \prec c_j$
- b. $(T_{reg} + T_{shm})_{c_j} < (T_{reg} + T_{shm})_{c_i} \Rightarrow c_i \prec c_j$
- c. $(T_{shm})_{c_j} < (T_{shm})_{c_i} \Rightarrow c_i \prec c_j$
- d. $(T_{reg})_{c_j} < (T_{reg})_{c_i} \Rightarrow c_i \prec c_j$
- e. $(S_{reg} + S_{shm})_{c_i} < (S_{reg} + S_{shm})_{c_j} \Rightarrow c_i \prec c_j$
- f. $(S_{shm})_{c_i} < (S_{shm})_{c_j} \Rightarrow c_i \prec c_j$
- g. $(S_{reg})_{c_i} < (S_{reg})_{c_j} \Rightarrow c_i \prec c_j$
- h. $i < j \Rightarrow c_i \prec c_j$

where rule h is only used to ensure that a total order can be found in case of tuples with strictly identical metrics.

It is possible for two or more candidate pairs to have the same metric, and the order in which the sorting rules are applied determines how such case is handled. For example, if the objective is to minimize data movements, \prec follows the rule sequence $a \rightarrow b \rightarrow c \rightarrow d \rightarrow e \rightarrow f \rightarrow g \rightarrow h$, where a is the primary sorting rule, and ties are further sorted by implementing the remaining sorting rules (e.g., b, c , etc.).

Once sorted, the topmost tuple of L_{tuple} is the most profitable candidate for fusion. If the fusion (1) does not exceed the hardware limit on shared memory, and (2) does not result in a prohibitively high stencil order (and consequently high quantity of redundant computations) if s_i and s_j are in producer-consumer relationship, then these nodes are fused. The stencil DAG is modified by replacing the original operators by the node modeling their fusion, and the dependence graph is updated as needed. The process is then repeated on this updated DAG, till no further fusion is feasible or profitable.

VII. RELATED WORK

Automatic high-performance GPU code generation for stencils has been a topic of active research for both CPUs [2], [9]–[11], [20], [34] and GPUs [4]–[8], [12], [13], [18], [19]. The main issues discussed in this paper about developing high-performance code for stencil computations have also been addressed in similar or equivalent ways in these efforts. In this section, we provide a brief discussion of the related work.

PPCG [19] is a polyhedral source-to-source compiler that generates classically time tiled OpenCL and CUDA code from an annotated sequential program. Patus [2] is a code generation and autotuning framework for stencil computations that can generate spatially tiled CUDA code without shared memory usage from the input DSL stencil specification. Mint [6] is a pragma-based source-to-source translator implemented in the ROSE compiler [35] that generates a spatially tiled CUDA code from traditional C code. Unlike these approaches that generate code for a single GPU device, Physis [7] translates user-written structured grid code into CUDA+MPI code for GPU-equipped clusters. Zhang and Mueller [21] develop an autotuning strategy for 3D stencils on GPUs. Their framework uses thread coarsening along different dimensions with judicious use of shared memory and registers to improve performance of spatially-tiled 3D stencils.

Overtile [4] and Forma [3] are DSL compilers that generate time-tiled CUDA code from an input stencil DSL specification. PolyMage [32] is a DSL based code generator for automatic optimization of image processing pipelines. All these approaches use overlapped tiling to fuse the stencil operators in the computation, with the intermediate arrays stored in shared memory. The code generated for 3D stencils by Overtile and Forma uses overlapped tiling along all three spatial dimensions and therefore incurs higher data movement and redundant computation overhead than streamed tiling [15]. Grosser et al. [5] implement the split tiling approach of [10], and hexagonal tiling [12] for temporal tiling in GPUs. ChiLL [36] is a composable loop transformation framework which allows the user to script loop transformations for stencil computations.

Micikevicius [17] presents a CUDA implementation of 3D finite difference computation that performs spatial tiling, and

Benchmark	N	T	k	FPP	Benchmark	N	T	k	FPP	Benchmark	N	T	k	FPP
j2d5pt	8192 ²	4	1	10	j3d7pt	512 ³	4	1	13	heat	512 ³	4	1	15
j2d9pt-gol	8192 ²	4	1	18	j3d13pt	512 ³	4	2	25	poisson	512 ³	4	1	21
j2d9pt	8192 ²	4	2	18	j3d17pt	512 ³	4	1	28	cheby	512 ³	4	1	39
gaussian	8192 ²	4	2	50	j3d27pt	512 ³	4	1	54	denoise	512 ³	4	2	62
gradient	8192 ²	4	1	18	curl	450 ³	1	1	36	hypterm	300 ³	1	4	358

N: Domain Size, T: Time Tile Size, k: Stencil Order, FPP: FLOPs per Point

TABLE I: Benchmark Characteristics

uses registers to alleviate shared memory pressure. Nguyen et al. [16] extend this implementation to a 3.5D-blocking algorithm for 3D stencils: streaming along one dimension, and temporal tiling along the other two dimensions. For CPUs, their approach reduces the cache footprint, and for GPUs, it reduces the shared memory required to time-tile a class of cross stencils.

In context of CPUs, the Pochoir [9] stencil compiler uses a DSL approach to generate high-performance Cilk code that uses an efficient parallel cache-oblivious algorithm to exploit reuse. The Pochoir system provides a C++ template library that allows the stencil specification to be executed directly in C++ without the Pochoir compiler, which aids debugging. Henretty et al. [10] propose hybrid split-tiling and nested split-tiling, to achieve parallelism without incurring redundant computations. Halide [20] is a DSL language for image processing pipelines. It decouples algorithm specification from its execution schedule. The advantage of this separation is that one can write multiple schedules for the same computation without rewriting the entire computation. Halide schedules can express loop nesting, parallelization, loop unrolling and vector instruction. The performance of the optimized code depends on the efficacy of the schedule. The schedule can be either written manually by a domain expert, or generated by extensive autotuning (e.g., with OpenTuner [37]); these approaches either require some degree of expertise with Halide DSL, or are time consuming. Recently, Mullapudi et al. [38] extended the scheduling strategy of PolyMage [32] to automatically generate schedules for Halide. Bandishti et al. [11] propose diamond tiling to ensure concurrent start-up as well as perfect load-balance whenever possible. Yount et al. describe the YASK [34] framework to simplify the task of defining stencil functions, generating high-performance code targeted for Intel Xeon Phi architecture. In Sec. VIII, we present experimental results with several of these stencil optimizers. Olschanowsky et al. [39] focus on large-scale PDE benchmarks as those written using Chombo [40] framework. Chombo parallelizes the application across the nodes using MPI, where each MPI process operates over a set of boxes, and each box applies a sequence of stencil operations over its domain. They evaluate the benefits of inter-loop stencil optimizations like loop shifting, loop fusion, wavefront tiling, and overlapped tiling on various box sizes. Davis et al. [41] use modified macro dataflow graphs to represent the stencil computation, and then apply fusion, rescheduling, and tiling optimizations to reduce communication volume and storage requirements. The stencil operations in the optimized graph

Resource	Details
Quadro GP100 GPU	3584 CUDA cores, 717 GB/s peak global memory bandwidth, 10.6 TFLOPS peak single-precision performance, 5.3 TFLOPS peak double-precision performance
Quadro GV100 GPU	5120 CUDA cores, 870 GB/s peak global memory bandwidth, 14.8 TFLOPS peak single-precision performance, 7.4 TFLOPS peak double-precision performance
Intel Skylake CPU	Intel core i7-6700K (4 cores, 2 threads/core 4.00 GHz, 512 GFLOPS peak single-precision performance)
Intel Xeon CPU	Intel Xeon E5-2680 (14 cores, 2 threads/core 2.40 GHz, 1.075 TFLOPS peak single-precision performance)
KNL phi	Intel Xeon Phi 7250 (68 cores, 4 threads/core 1.40GHz, 6.092 TFLOPS peak single-precision performance)

TABLE II: Benchmarking Hardware

Compiler	Flags
<i>nvcc</i>	-O3 -maxrregcount= <i>n</i> -ccbin=g++ -std=c++11 -Xcompiler "-fPIC -fopenmp -O3 -fno-strict-aliasing" -use_fast_math -Xptxas "-v -dlcm=ca" -arch=sm_60/70
<i>gcc</i>	-Ofast -fopenmp -ffast-math -fstrict-aliasing -march=core-avx2/{ <i>kn</i> l} -mavx512f -mavx512er -mavx512cd -mavx512pf} -mfma -ffp-contract=fast
<i>llvm</i>	-Ofast -fopenmp=libomp -mfma -ffast-math -fstrict-aliasing -march=core-avx2/{ <i>kn</i> l} -mllvm -ffp-contract=fast
<i>icc</i>	-Ofast -qopenmp -no-prec-div -march=native -fma -xMIC-AVX512 -fp-model fast=2 -complex-limited-range -override-limits

TABLE III: Compilation flags for GPU, multi-core CPU, and Xeon Phi

can then be reinterpreted as relations bounded by affine constraints, to enable automatic code generation using ISL [42]. They also propose a cost model to compare schedules based on the memory traffic.

Wahib and Maruyama [13] pose kernel fusion as an optimization problem, and use a code-less performance model to choose a near-optimal fusion configuration amongst other possible variants. The space of feasible solutions is pruned using a search heuristic based on a hybrid grouping genetic algorithm. Gysi et al. [8] also propose a model-driven stencil optimization approach. Like Wahib and Maruyama [13], they use a code-less performance model to find the best fusion configuration among *all* valid topological sorts of the stencil DAG. The model has been used to guide kernel fusion in the Stella library [43]. Prajapati et al. [44] propose an analytical model that predicts the execution time of the code generated with hexagonal tiling.

VIII. EXPERIMENTAL EVALUATION

A. Evaluation on GPUs

Experimental setup: We have implemented the tiling schemes and fusion heuristics described in Secs. V and VI

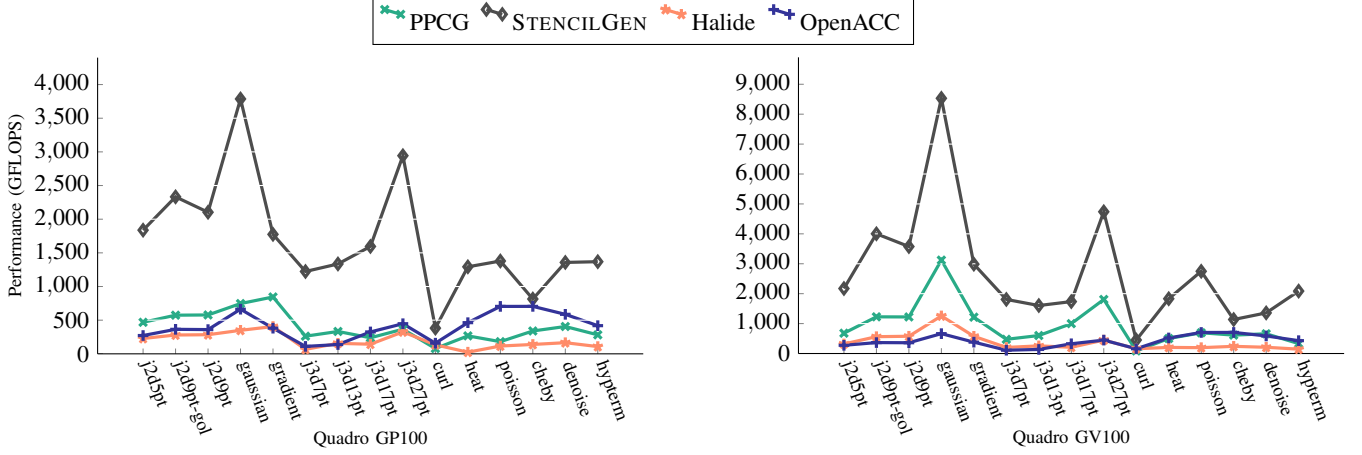


Fig. 1: Performance of single-precision benchmarks on Pascal and Volta GPU devices

Fig. 9: Performance of single-precision benchmarks on Pascal and Volta GPU devices

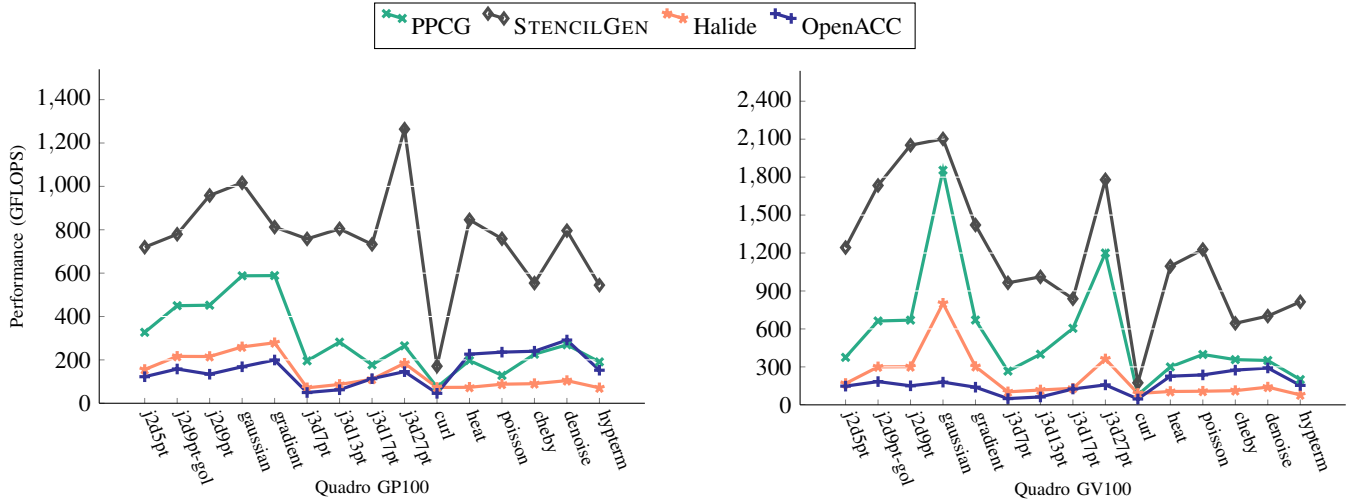


Fig. 1: Performance of double-precision benchmarks on Pascal and Volta GPU devices

Fig. 10: Performance of double-precision benchmarks on Pascal and Volta GPU devices

into the STENCILGEN code generator, and we compare below the performance of the STENCILGEN-generated code against PPCG-0.08 [19], OpenACC-17.4 [45], and the auto-scheduler branch of Halide [38]. The benchmarks used in evaluation are listed in Table I. We evaluate the performance of both single- and double-precision versions of the benchmarks on Pascal GP100 and Volta GV100 devices; the hardware is detailed in Table II. For PPCG and STENCILGEN, the generated code was compiled using NVCC 9.1 [46]. The Halide generated code was compiled with LLVM 3.7. The compilation flags are listed in Table III.

Code generation: PPCG performs classical time tiling along with thread coarsening. Mapping multiple iterations to a thread exposes instruction level parallelism. Coarsen-

ing within the sustainable per thread register pressure aids register level reuse, and helps hide memory access latency by exposing instruction-level parallelism [47]. We created multiple versions by tuning the block size and unrolling factors for PPCG, and report results for the version with best performance. STENCILGEN does not have the support for thread coarsening at present. For each benchmark, we leverage operator associativity to alleviate pressure on GPU resources. GPUs allow compile-time flexibility in assigning the number of registers per thread. The performance of the same code with varying registers per thread can vary dramatically. This flexibility particularly affects the fusion heuristics for 3D benchmarks. We use STENCILGEN to generate multiple versions for the 3D benchmarks with different number of

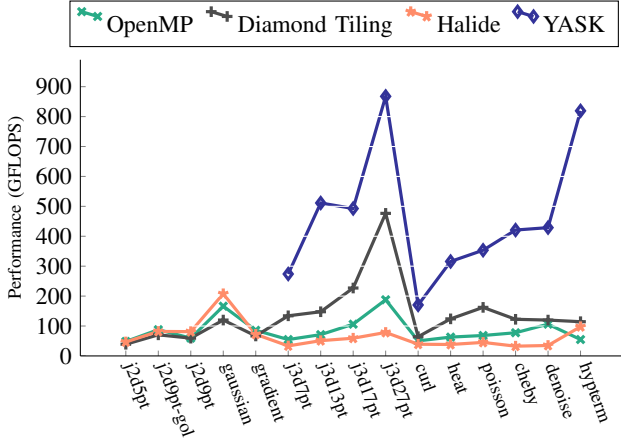


Fig. 12: Performance of benchmarks on Xeon Phi 7250

explicit registers. The versions with higher limit of explicit registers have a greater degree of fusion, and are compiled with higher registers per thread (e.g., `maxregcount=64` for *heat* and *poisson*, `maxregcount=255` for *hypterm*). Pascal and Volta devices can load the read-only data through the cache used by texture pipeline. To enable this feature, the read-only data in STENCILGEN-generated code is automatically annotated with the `__restrict__` keyword. Both PPCG and STENCILGEN-generated codes are compiled with different registers-per-thread settings to find the best configuration. To ensure coalescence, the fastest-varying dimension of the thread block is never chosen as a streaming dimension.

GPU performance results: Figs. 9 and 10 plot the performance of STENCILGEN-generated single- and double-precision codes, respectively, against different code generators on the two GPU devices. STENCILGEN systematically outperforms the other code generators in our experiments, by a factor up to $9\times$. Our optimization scheme for associative stencils helps reduce shared memory requirement for stencils like *gaussian*, *j3d27pt*, and *chebyshev*, allowing us to achieve higher occupancy, and consequently higher performance. We achieve 1.0 TFLOPS (2.1 TFLOPS) for double-precision *gaussian* stencil, 1.2 TFLOPS (1.7 TFLOPS) for double-precision *j3d27pt* stencil, and 2.2 TFLOPS (2.6 TFLOPS) for double-precision *chebyshev* stencil on GP100 (GV100) device. For the 2D stencils, the best performance is achieved by overlapped tiling along both x and y dimensions, and concurrent-streaming along y dimension. For the 3D stencils, overlapped tiling along x and y dimensions, and serial-streaming along z dimension gives best performance. For 3D order-2 time-iterated stencils, restricting the fusion to two time steps yields better performance due to lower recomputation volume.

Despite applying a plethora of optimizations, we achieve only 37% and 58% (44% and 36%) of the single-precision (double-precision) peak performance on GP100 and GV100, respectively. This serves to highlight the bandwidth-bound nature of stencil computations. For theoretically compute-bound stencils like *hypterm* and *denoise*, high register pressure lowers the achieved occupancy, making it difficult to hide the memory access latency. One can further reduce the bandwidth,

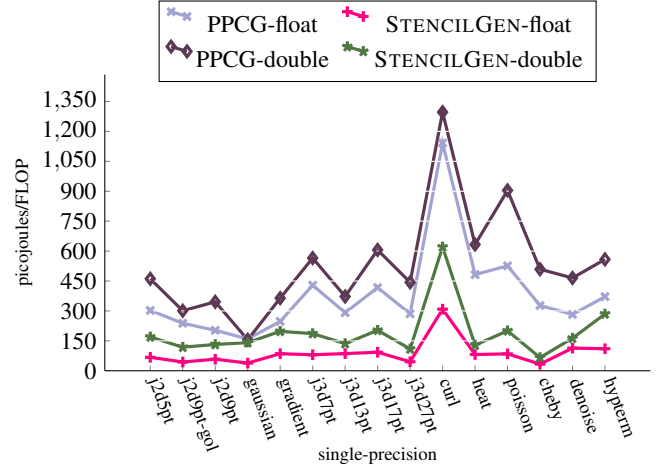


Fig. 13: Energy expended per FLOP on Pascal device

and improve the ILP by exploring optimizations like thread-coarsening, vectorized loads, and hiding access latency via data prefetching.

B. Stencil Computations on CPUs

Experimental setup: Many frameworks explicitly target stencil optimizations on CPU [2], [10], [11], [20], [34]. We evaluate the CPU code generated by diamond tiling [11] and Halide with auto-scheduling support [38] against a baseline OpenMP code on two different multi-core CPU processors, and a Xeon Phi processor. We also evaluate Intel’s latest stencil optimization framework, YASK [34], on the Xeon Phi processor. Due to the limitations of the framework, we could only generate YASK code for 3D stencils. We use three different compilers, namely ICC-17.0 [48], LLVM-5.0 [49], and GCC-7.2 [50] to compile the generated C/C++ stencil code, and report the highest performance.

CPU performance results: The performance is plotted in Figs. 11 and 12. One can observe that diamond tiling outperforms both Halide and the baseline code with its time tiling and concurrent start optimizations. On Xeon Phi, the performance of YASK is far superior to that of other frameworks. YASK performs optimizations like vector folding, cache blocking, and temporal wave-front tiling, that are specifically tuned for the Xeon Phi architecture.

For the multi-core CPUs, only diamond tiling is able to achieve nearly 30%-45% of the machine peak for *j3d27pt* stencil. Despite the high performance, YASK is able to achieve only 15% of the machine peak.

C. Energy Efficiency

Figs. 13, 14, and 15 plot the energy expended in a single floating-point operation for GP100, multi-core CPUs, and Xeon Phi respectively. For the stencil codes compiled with NVCC, we use the NVIDIA Management Library (NVML³) to measure the Joules/FLOP ratio. For multi-core CPUs and

³<https://developer.nvidia.com/nvidia-management-library-nvml>

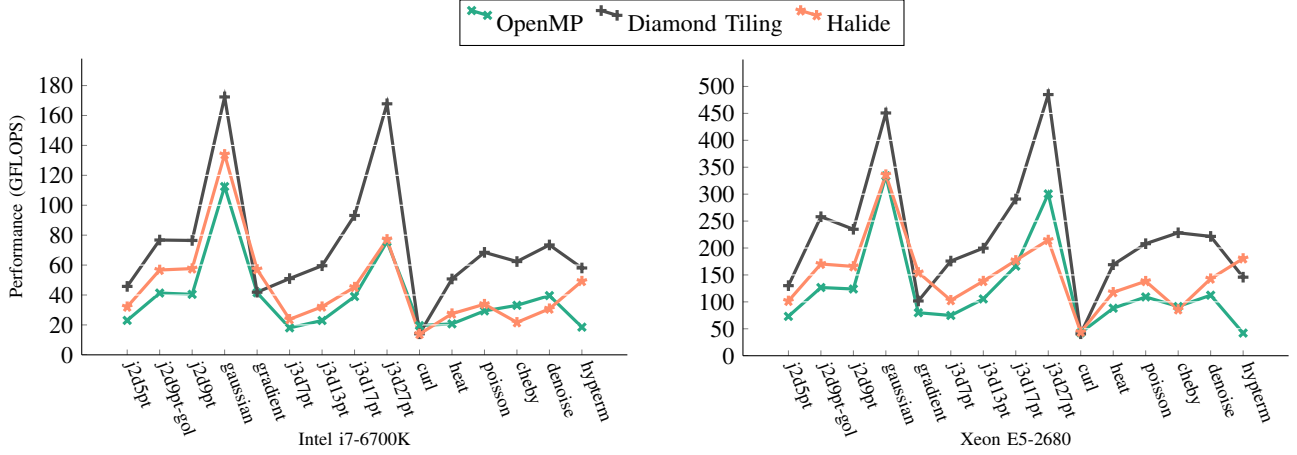


Fig. 11: Performance of benchmarks on Skylake and Xeon multi-core CPU

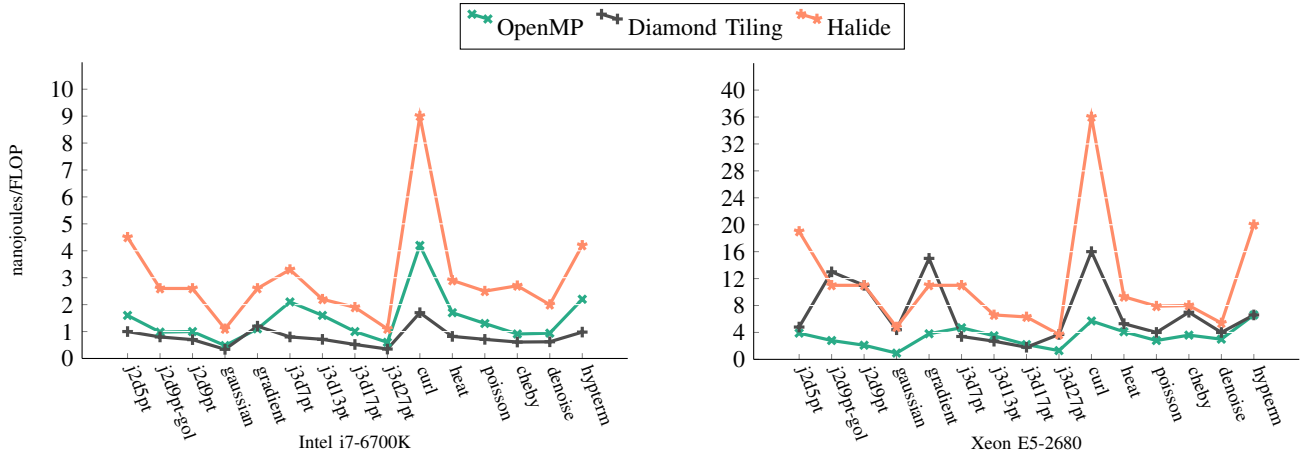


Fig. 14: Energy expended per FLOP on Skylake and Xeon multi-core CPU

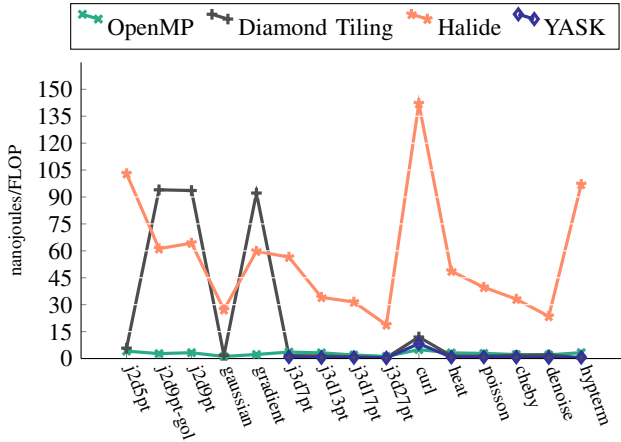


Fig. 15: Energy expended per FLOP on Xeon Phi 7250

Xeon Phi, we use RAPL⁴ to obtain the measurement. Although typically measurements indicate that faster code versions have a lower Joules per FLOP ratio, data shows situations where a faster code has a lower energy efficiency as when comparing

Halide and OpenMP versions of *j2d9pt-gol* on Xeon Phi. We conjecture the use of more power-hungry instructions (e.g., FMA) and additional in-processor data traffic as possible causes for this slight decrease in energy efficiency, as unfortunately obtaining fine-grain power measurements to determine the root cause is not feasible with our energy measurement setup.

Note that the benchmarks running on GP100 have a much lower ($1 \times - 0.1 \times$) Joules/FLOP ratio that the multi-core CPUs and Xeon Phi, while also achieving slightly higher peak performance. It indicates that the current generation of GPU devices appear more energy efficient than these CPU processors in executing such bandwidth-bound stencil computations.

IX. CONCLUSION

Stencil computations are at the computational core for many applications. Unlike dense linear algebra computations, where efficient libraries are widely available, the variety of manifestations of stencil patterns makes it infeasible to create libraries for stencils. However, the fundamental characteristics of the data access patterns for stencil computations can be used to devise domain-specific and target-specific optimizations in a source-to-source transformer and code generator

⁴<http://web.eece.maine.edu/~vweaver/projects/rapl/>

for stencils. This paper has presented an analysis of the fundamental considerations in achieving high performance for stencil computations on GPUs, focusing on tiling strategies that make effective use of key resources like shared memory and registers. Experimental results demonstrate the significant benefits from use of domain-specific optimization over state-of-the-art general-purpose optimizers.

REFERENCES

- [1] “Advanced Stencil-Code Engineering (ExaStencils),” <http://www.exastencils.org/>.
- [2] M. Christen, O. Schenk, and H. Burkhardt, “PATUS: A code generation and autotuning framework for parallel iterative stencil computations on modern microarchitectures,” in *Proceedings of the 2011 IEEE International Parallel & Distributed Processing Symposium*, ser. IPDPS ’11. IEEE Computer Society, 2011, pp. 676–687.
- [3] M. Ravishankar, J. Holewinski, and V. Grover, “Forma: A DSL for image processing applications to target GPUs and multi-core CPUs,” in *Proceedings of the 8th Workshop on General Purpose Processing Using GPUs*, ser. GPGPU-8. ACM, 2015, pp. 109–120.
- [4] J. Holewinski, L.-N. Pouchet, and P. Sadayappan, “High-performance code generation for stencil computations on GPU architectures,” in *Proceedings of the 26th ACM International Conference on Supercomputing*, ser. ICS ’12. ACM, 2012, pp. 311–320.
- [5] T. Gresser, A. Cohen, P. H. J. Kelly, J. Ramanujam, P. Sadayappan, and S. Verdoolaege, “Split tiling for GPUs: Automatic parallelization using trapezoidal tiles,” in *Proceedings of the 6th Workshop on General Purpose Processor Using Graphics Processing Units*, ser. GPGPU-6. ACM, 2013, pp. 24–31.
- [6] D. Unat, X. Cai, and S. B. Baden, “Mint: Realizing CUDA performance in 3D stencil methods with annotated C,” in *Proceedings of the International Conference on Supercomputing*, ser. ICS ’11. New York, NY, USA: ACM, 2011, pp. 214–224.
- [7] N. Maruyama, T. Nomura, K. Sato, and S. Matsuoka, “Physis: An implicitly parallel programming model for stencil computations on large-scale GPU-accelerated supercomputers,” in *Proceedings of 2011 International Conference for High Performance Computing, Networking, Storage and Analysis*, ser. SC ’11. New York, NY, USA: ACM, 2011, pp. 11:1–11:12.
- [8] T. Gysi, T. Gresser, and T. Hoefer, “MODESTO: Data-centric analytic optimization of complex stencil programs on heterogeneous architectures,” in *Proceedings of the 29th ACM on International Conference on Supercomputing*, ser. ICS ’15. New York, NY, USA: ACM, 2015, pp. 177–186.
- [9] Y. Tang, R. A. Chowdhury, B. C. Kuszmaul, C.-K. Luk, and C. E. Leiserson, “The Pochoir stencil compiler,” in *Proceedings of the Twenty-third Annual ACM Symposium on Parallelism in Algorithms and Architectures*, ser. SPAA ’11. New York, NY, USA: ACM, 2011, pp. 117–128.
- [10] T. Henretty, R. Veras, F. Franchetti, L.-N. Pouchet, J. Ramanujam, and P. Sadayappan, “A stencil compiler for short-vector SIMD architectures,” in *Proceedings of the 27th International ACM Conference on International Conference on Supercomputing*, ser. ICS ’13. New York, NY, USA: ACM, 2013, pp. 13–24.
- [11] V. Bandishti, I. Pananilath, and U. Bondhugula, “Tiling stencil computations to maximize parallelism,” in *High Performance Computing, Networking, Storage and Analysis (SC), 2012 International Conference for*, Nov 2012, pp. 1–11.
- [12] T. Gresser, A. Cohen, J. Holewinski, P. Sadayappan, and S. Verdoolaege, “Hybrid hexagonal/classical tiling for GPUs,” in *Proceedings of Annual IEEE/ACM International Symposium on Code Generation and Optimization*, ser. CGO ’14. ACM, 2014, pp. 66:66–66:75.
- [13] M. Wahib and N. Maruyama, “Scalable kernel fusion for memory-bound GPU applications,” in *Proceedings of the International Conference for High Performance Computing, Networking, Storage and Analysis*, ser. SC ’14. IEEE Press, 2014, pp. 191–202.
- [14] P. S. Rawat, C. Hong, M. Ravishankar, V. Grover, L.-N. Pouchet, A. Rountev, and P. Sadayappan, “Resource conscious reuse-driven tiling for GPUs,” in *Proceedings of the 2016 International Conference on Parallel Architectures and Compilation*, ser. PACT ’16. ACM, 2016, pp. 99–111.
- [15] P. S. Rawat, C. Hong, M. Ravishankar, V. Grover, L.-N. Pouchet, and P. Sadayappan, “Effective resource management for enhancing performance of 2D and 3D stencils on GPUs,” in *Proceedings of the 9th Annual Workshop on General Purpose Processing Using Graphics Processing Unit*, ser. GPGPU ’16. ACM, 2016, pp. 92–102.
- [16] A. Nguyen, N. Satish, J. Chhugani, C. Kim, and P. Dubey, “3.5-D blocking optimization for stencil computations on modern CPUs and GPUs,” in *Proceedings of the 2010 ACM/IEEE International Conference for High Performance Computing, Networking, Storage and Analysis*, ser. SC ’10. IEEE Computer Society, 2010, pp. 1–13.
- [17] P. Micikevicius, “3D finite difference computation on GPUs using CUDA,” in *Proceedings of 2nd Workshop on General Purpose Processing on Graphics Processing Units*, ser. GPGPU-2. ACM, 2009, pp. 79–84.
- [18] M. Ravishankar, P. Micikevicius, and V. Grover, “Fusing convolution kernels through tiling,” in *Proceedings of the 2nd ACM SIGPLAN International Workshop on Libraries, Languages, and Compilers for Array Programming*, ser. ARRAY 2015. ACM, 2015, pp. 43–48.
- [19] S. Verdoolaege, J. C. Juega, A. Cohen, J. I. Gómez, C. Tenllado, and F. Catthoor, “Polyhedral parallel code generation for CUDA,” *ACM TACO*, vol. 9, no. 4, pp. 54:1–54:23, Jan. 2013.
- [20] J. Ragan-Kelley, C. Barnes, A. Adams, S. Paris, F. Durand, and S. Amarasinghe, “Halide: A language and compiler for optimizing parallelism, locality, and recomputation in image processing pipelines,” in *Proceedings of the 34th ACM SIGPLAN Conference on Programming Language Design and Implementation*, ser. PLDI ’13. ACM, 2013, pp. 519–530.
- [21] Y. Zhang and F. Mueller, “Auto-generation and auto-tuning of 3D stencil codes on GPU clusters,” in *Proceedings of the Tenth International Symposium on Code Generation and Optimization*, ser. CGO ’12. New York, NY, USA: ACM, 2012, pp. 155–164.
- [22] K. Datta, M. Murphy, V. Volkov, S. Williams, J. Carter, L. Oliker, D. Patterson, J. Shalf, and K. Yelick, “Stencil computation optimization and auto-tuning on state-of-the-art multicore architectures,” in *Proceedings of the 2008 ACM/IEEE Conference on Supercomputing*, ser. SC ’08. IEEE Press, 2008, pp. 4:1–4:12.
- [23] J. Meng and K. Skadron, “Performance modeling and automatic ghost zone optimization for iterative stencil loops on GPUs,” in *Proceedings of the 23rd International Conference on Supercomputing*, ser. ICS ’09. New York, NY, USA: ACM, 2009, pp. 256–265.
- [24] A. Vizitiu, L. Itu, C. Ni, and C. Suciu, “Optimized three-dimensional stencil computation on fermi and kepler GPUs,” in *2014 IEEE High Performance Extreme Computing Conference (HPEC)*, Sept 2014, pp. 1–6.
- [25] N. Maruyama and T. Aoki, “Optimizing stencil computations for NVIDIA Kepler GPUs,” 2014.
- [26] M. Christen, O. Schenk, E. Neufeld, P. Messmer, and H. Burkhardt, “Parallel data-locality aware stencil computations on modern micro-architectures,” in *2009 IEEE International Symposium on Parallel Distributed Processing*, May 2009, pp. 1–10.
- [27] A. D. Pereira, M. Castro, M. A. R. Dantas, R. C. O. Rocha, and L. F. W. Ges, “Extending OpenACC for efficient stencil code generation and execution by skeleton frameworks,” in *2017 International Conference on High Performance Computing Simulation (HPCS)*, July 2017, pp. 719–726.
- [28] X. Zhou, J.-P. Giacalone, M. J. Garzarán, R. H. Kuhn, Y. Ni, and D. Padua, “Hierarchical overlapped tiling,” in *Proceedings of the Tenth International Symposium on Code Generation and Optimization*, ser. CGO ’12. New York, NY, USA: ACM, 2012, pp. 207–218.
- [29] P. Basu, S. Williams, B. V. Straalen, L. Oliker, P. Colella, and M. Hall, “Compiler-based code generation and autotuning for geometric multigrid on GPU-accelerated supercomputers,” *Parallel Computing*, vol. 64, pp. 50 – 64, 2017, high-End Computing for Next-Generation Scientific Discovery.
- [30] P. S. Rawat, “Optimization of stencil computations on GPUs,” in *Electronic Thesis or Dissertation. Ohio State University*, 2018, pp. 1–200.
- [31] K. Stock, M. Kong, T. Gresser, L.-N. Pouchet, F. Rastello, J. Ramanujam, and P. Sadayappan, “A framework for enhancing data reuse via associative reordering,” in *Proceedings of the 35th ACM SIGPLAN Conference on Programming Language Design and Implementation*, ser. PLDI ’14. New York, NY, USA: ACM, 2014, pp. 65–76.
- [32] R. T. Mullapudi, V. Vasista, and U. Bondhugula, “PolyMage: Automatic optimization for image processing pipelines,” in *Proceedings of the Twentieth International Conference on Architectural Support for Programming Languages and Operating Systems*, ser. ASPLOS ’15. ACM, 2015, pp. 429–443.

- [33] L.-N. Pouchet, U. Bondhugula, C. Bastoul, A. Cohen, J. Ramanujam, P. Sadayappan, and N. Vasilache, "Loop transformations: Convexity, pruning and optimization," in *38th Annual ACM SIGPLAN-SIGACT Symposium on Principles of Programming Languages (POPL'11)*. ACM, 2011, pp. 549–562.
- [34] C. Yount, J. Tobin, A. Breuer, and A. Duran, "YASK-yet another stencil kernel: A framework for HPC stencil code-generation and tuning," in *Proceedings of the Sixth International Workshop on Domain-Specific Languages and High-Level Frameworks for HPC*, ser. WOLFHP '16. Piscataway, NJ, USA: IEEE Press, 2016, pp. 30–39.
- [35] D. J. Quinlan, "Rose: Compiler support for object-oriented frameworks." vol. 10, pp. 215–226, 06 2000.
- [36] M. Khan, P. Basu, G. Rudy, M. Hall, C. Chen, and J. Chame, "A script-based autotuning compiler system to generate high-performance CUDA code," *ACM Trans. Archit. Code Optim.*, vol. 9, no. 4, pp. 31:1–31:25, Jan. 2013.
- [37] J. Ansel, S. Kamil, K. Veeramachaneni, J. Ragan-Kelley, J. Bosboom, U.-M. O'Reilly, and S. Amarasinghe, "OpenTuner: An extensible framework for program autotuning," in *Proceedings of the 23rd International Conference on Parallel Architectures and Compilation*, ser. PACT '14. New York, NY, USA: ACM, 2014, pp. 303–316.
- [38] R. T. Mullanpudi, A. Adams, D. Sharlet, J. Ragan-Kelley, and K. Fatahalian, "Automatically scheduling Halide image processing pipelines," *ACM Trans. Graph.*, vol. 35, no. 4, pp. 83:1–83:11, Jul. 2016.
- [39] C. Olschanowsky, M. M. Strout, S. Guzik, J. Loffeld, and J. Hittinger, "A study on balancing parallelism, data locality, and recomputation in existing PDE solvers," in *SC14: International Conference for High Performance Computing, Networking, Storage and Analysis*, Nov 2014, pp. 793–804.
- [40] M. Adams, P. Colella, D. T. Graves, J. Johnson, N. Keen, T. J. L. D. F. M. P. McCorquodale, D. M. P. Schwartz, T. Sternberg, and B. V. Straalen, in *Chombo Software Package for AMR Applications - Design Document*, Lawrence Berkeley National Laboratory Technical Report LBNL-6616E, 2009.
- [41] E. C. Davis, M. M. Strout, and C. Olschanowsky, "Transforming loop chains via macro dataflow graphs," in *Proceedings of the 2018 International Symposium on Code Generation and Optimization*, ser. CGO 2018. New York, NY, USA: ACM, 2018, pp. 265–277.
- [42] S. Verdoolaege, "isl: An integer set library for the polyhedral model," in *Mathematical Software - ICMS 2010*, ser. Lecture Notes in Computer Science. Springer, 2010, vol. 6327, pp. 299–302.
- [43] 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*, ser. SC '15. ACM, 2015, pp. 41:1–41:12.
- [44] N. Prajapati, W. Ranasinghe, S. Rajopadhye, R. Andonov, H. Djidjev, and T. Grosser, "Simple, accurate, analytical time modeling and optimal tile size selection for gpgpu stencils," in *Proceedings of the 22Nd ACM SIGPLAN Symposium on Principles and Practice of Parallel Programming*, ser. PPOPP '17. New York, NY, USA: ACM, 2017, pp. 163–177.
- [45] "OpenACC compiler," https://www.pgroup.com/doc/openacc_gs.pdf, 2017.
- [46] "NVCC," docs.nvidia.com/cuda/cuda-compiler-driver-nvcc, 2017.
- [47] V. Volkov and J. W. Demmel, "Benchmarking GPUs to tune dense linear algebra," in *Proceedings of the 2008 ACM/IEEE Conference on Supercomputing*, ser. SC '08. IEEE Press, 2008, pp. 31:1–31:11.
- [48] "Intel C++ compiler," <https://software.intel.com/en-us/>, 2017.
- [49] C. Lattner and V. Adve, "LLVM: A compilation framework for lifelong program analysis & transformation," in *Proceedings of the International Symposium on Code Generation and Optimization: Feedback-directed and Runtime Optimization*, ser. CGO '04. Washington, DC, USA: IEEE Computer Society, 2004, pp. 75–.
- [50] R. M. Stallman and G. DeveloperCommunity, *Using The GNU Compiler Collection: A GNU Manual For GCC Version 4.3.3*. Paramount, CA: CreateSpace, 2009.