

# Re-Write Rules for Flattening in Fast Memory the Parallelism of Hull-White Trinomial Pricing

COSMIN E. OANCEA, University of Copenhagen, Denmark

WOJCIECH M. PAWLAK, SimCorp, Denmark

This document sketches a formalization of the re-write rules used for optimizing the parallel GPU execution of an implementation of the Hull-White One-Factor Short-Rate pricing model, which uses the trinomial-tree lattice-based method for numerically solving the underlying stochastic differential equation.

The proposed code transformation is a specialization of Blelloch's flattening transformation, which maps irregular nested parallelism into a sequence of flat-parallel constructs, such that the work-depth asymptotics of the original program is respected.

Our transformation applies on code with two levels of parallelism: the first level corresponding to a map operation, and the second level corresponding to multiple parallel operations—potentially nested within sequential loops—but which have the same length for the same iteration of the outer map.

The main difference with respect to the classical approach is that our approach does not replicate the variables that are bound in the outer map but free in the inner parallel constructs. Instead, such variables are found by indirectly accessing some auxiliary arrays which are reused across equal-shape arrays and across sequential recurrences (loops). This technique reduces the footprint of scratchpad/fast memory, whose efficient use is paramount to achieving efficient GPU execution.

Additional Key Words and Phrases: parallelism, compilers, flattening transformation, dependence analysis, dynamic analysis, GPU, computational finance, derivative pricing.

## 1 INTRODUCTION

This document sketches a formalization by means of re-write rules of a specialized instance of Blelloch's flattening code transformation [Blelloch and Greiner 1996; Blelloch et al. 1994].

The transformation was motivated by the work [Pawlak et al. 2021] on optimizing for GPU execution of a standard pricing algorithm, namely the Hull-White One-Factor Short-Rate (*HW1F*) model [Hull and White 1994]. The latter defines the value of a financial instrument by means of a stochastic differential equation, that represents the random changes in the interest rate over time, which is solved numerically by means of a *Trinomial Tree* lattice-based numerical method [Hull and White 1996].

In this algorithm, pricing *one* financial instrument has two main stages. The forward stage builds a tree of bounded width, representing the propagation of the interest rate until the maturity of the underlying bond is reached.<sup>1</sup> The backward stage then performs the instrument valuation from maturity back to the current time. The computational structure thus consists of two sequential time-series loops, whose counts are the height of the tree, in which each iteration performs several semantically-parallel operations, which have the same length as the bounded width of the tree. For example, the forward step computes each node at the current level (time step) from three nodes at the previous level (time step) in the tree.

---

<sup>1</sup> The model uses a tree of bounded width, because in practice the short-interest rate tends to revert to the mean value over time.

---

Authors' addresses: Cosmin E. Oancea, University of Copenhagen, Copenhagen, 2100, Denmark, cosmin.oancea@diku.dk; Wojciech M. Pawlak, SimCorp, Copenhagen, 2100, Denmark, Wojciech.M.Pawlak@simcorp.com.

---

2021. /2021/5-ART1 \$15.00

<https://doi.org/>

Practical use cases require pricing of a (large) portfolio of instruments, which gives rise to an outer level of parallelism, as instruments can be priced independently of each other.

*The main challenge*, however, is that in realistic scenarios *the width and height of the trees is highly variant* across the portfolio instruments. This gives rise to a case of irregular-nested parallelism that is difficult to map efficiently to GPU hardware. In particular, related parallelization approaches either assume the homogeneous case [Gerbessiotis 2004; Huang and Thulasiram 2005] in which all trees have the same height and width, or acknowledge the problem, but do not offer a solution [Grauer-Gray et al. 2013].

This document considers the challenging case of a heterogeneous portfolio, in which the heights and widths of various instrument trees may differ significantly. At a very high-level, (at least) two parallelization strategies make sense.

The first strategy follows the common wisdom that states that outer levels of parallelism are more profitable to exploit than inner ones. As such, given a big enough portfolio, one should sequentialize the inner parallelism and perform one instrument valuation per thread. The problem with this approach is that the high variance of widths and heights across a portfolio instruments introduces two levels of thread divergence on GPU. These levels correspond

- (1) to the sequential time loop of variant height and
- (2) to the inner width-parallel operations, which are sequentialized.

This strategy allows to optimise *one level of divergence, but not both*; for example, by precomputing the widths (or heights) of the instrument' trees, and by sorting the portfolio in decreasing order of their widths (or heights).

The second strategy is to exploit both levels of parallelism, even when the outer parallelism is enough to fully utilize hardware. This allows to better exploit temporal locality by maintaining most of the intermediate arrays in the fast (shared + register) memory, and to optimize *both* levels of divergence. The height-level of divergence is optimized by sorting as before. To optimize the width-level of divergence the idea is

- (1) to *pack instruments into bins*, such that their summed widths fits the size of the thread-block (bin), and then
- (2) to “flatten” (merge) the available two-level parallelism—the first level corresponds to the instruments in a bin, and the second level corresponds to the width-length parallel operations that appear in the implementation of each instrument.

The flattening step is non-trivial and is the main focus of this document: In Section 2.2, we document it by presenting an initial simplified nested-parallel specification that uses operators such as map and reduce. In Section 3, we give the gist of the technique by demonstrating how the flat-parallel program is obtained, named GPU-FLAT.

Section 4 is the core contribution of this document that corresponds to sketching a formalization of the proposed transformation by a set of inference (rewrite) rules, which can possibly be integrated in the repertoire of a data-parallel compiler.

Please be advised that the material other than Section 4 has been already presented elsewhere [Pawlak et al. 2021]; we recount it here for convenience, i.e., to provide the intuition for the rewrite rules and to make this document self contained.

## 2 NOTATION AND NESTED-PARALLEL SPECIFICATION

This chapter introduces first a functional notation that is then used, in Section 2.2, to present the nested-parallel structure of the pricing algorithm. In comparison to a lower-level loop-based notation, such as CUDA, the functional notation has the advantage that it (i) enables a concise specification of all available (nested) parallelism in terms of well-known data-parallel operators

such as `map`, `reduce`, and `scan`, and (ii) allows to demonstrate at a high level the rewrite rules, that are applied for deriving GPU-FLAT.

## 2.1 Functional Notation

We use `i32` and `real` to denote the type of a 32-bit integer and (single- or double-precision) floating-point numbers, respectively,  $[N]\alpha$  to denote the type of an array, whose elements have type  $\alpha$ ,  $[a_1, \dots, a_n]$  to denote an array literal, and  $(a, b)$  to denote a tuple (record) value. Function  $f$  applied to two arguments  $a$  and  $b$  is written as  $f\ a\ b$  (without parenthesis or commas).

The notation supports the usual unary/binary operators and (normalized) `let` bindings, which have the form `let a = e1 let b = e2... in en` and are similar to a block of statements followed by a return denoted by keyword `in`. In-place updates to array elements are allowed and are written as `let arr[i] = x`.<sup>2</sup> The notation supports (i) conditional if statement `if c then e1 else e2` with semantics similar to the C ternary operator `c? e1 : e2`, and (ii) sequential loop expression:

$$\mathbf{loop}\ (x^1, \dots, x^m) = (x_0^1, \dots, x_0^m)\ \mathbf{for}\ i < n\ \mathbf{do}\ e.$$

Here,  $x^{1\dots m}$  are loop-variant variables that are initialized for the first iteration with in-scope variables  $x_0^{1\dots m}$ . The loop executes iterations  $i$  from  $0$  to  $n-1$ , and the result of the loop-body expression  $e$  provides the values of  $x^{1\dots m}$  for the next iteration. The initialization part may be syntactically omitted—i.e., `loop (x1, ..., xm) for i < n do e` is legal—in which case the initialization refers to in-scope variables bearing the same name ( $x^{1\dots m}$ ).

Most importantly, the notation supports several key-parallel operators, whose types and semantics are shown in Listing 1: `iota` applied to integer  $n$  creates the array with elements from  $0$  to  $n-1$  (i.e., an iteration space), and `replicate n v` creates an array of length  $n$  whose elements are all  $v$ . A `map` operation applies its function argument  $f$  to each element of the input array, resulting in an array of same length. The function can be declared in the program or can be an anonymous ( $\lambda$ ) function—e.g., `map ( $\lambda x \rightarrow x+1$ ) arr` adds one to each element of `arr`. Similarly, `map2` applies its function argument to corresponding elements from its two input arrays. `reduce` successively applies a binary-*associative* operator  $\odot$  to all elements of its input array, where  $e_0$  denotes the neutral element of  $\odot$ . `scan` [Blelloch 1989] (a.k.a., parallel-prefix sum) is similar to `reduce`, except that it produces an array of the same length ( $n$ ) containing all prefix sums of its input array. The *inclusive* scan (`scaninc`) starts with the first element of the array, and the *exclusive* scan (`scanexc`) starts with the neutral element.

`Segmented scan` (`sgmscan`) has the semantics of a scan applied to each subarray of an irregular array of subarrays. The latter has a flat representation consisting of (i) a flag array composed of 0s and 1s, where 1 denotes the start position of a subarray, and (ii) by a length-matching flat array containing in order all elements of all subarrays. For example, `flag = [1, 0, 1, 0, 0, 0, 1]` denotes an array with three rows, having two, four and one elements, respectively, and `sgmscaninc (+) 0 flag [1, 2, 3, 4, 5, 6, 7]` results in array `[1, 3, 3, 7, 12, 18, 7]`. `Segmented scan` can be implemented in terms of a scan with a modified operator [Blelloch 1989], e.g., for the *inclusive* one:

$$\lambda(f1, v1)\ (f2, v2) \rightarrow \mathbf{if}\ f2 \neq 0\ \mathbf{then}\ (f1\ |f2, v2)\ \mathbf{else}\ (f1\ |f2, v1\ \odot\ v2)$$

The last operator `scatter x is vs` updates in place the array  $x$  at indices contained in array  $is$  with the values contained in array  $vs$ , except that out-of-bounds indices are ignored (not updated). For example, in Listing 1, value  $b_1$  is not written in the result because its index  $-1$  is out of bounds.

<sup>2</sup> In place updates can be supported without affecting language purity by means of a uniqueness type mechanism [Henriksen et al. 2017].

```

1  iota : (n: i32) → [n]i32
2  iota n = [0, ..., n-1]
3
4  replicate : (n: i32) → α → [n]α
5  replicate n v = [v, ..., v]
6
7  map : ∀ n. (α → γ) → [n]α → [n]γ
8  map f [a1, ..., an] = [f a1, ..., f an]
9
10 map2: ∀ n. (α → β → γ) → [n]α → [n]β → [n]γ
11 map2 g [a1, ..., an] [b1, ..., bn] =
12   [g a1 b1, ..., g an bn]
13
14 reduce: ∀ n. (α → α → α) → α → [n]α → α
15 reduce ⊙ e⊙ [a1, ..., an] = a1 ⊙ ... ⊙ an
16
17 scan: ∀ n. (α → α → α) → α → [n]α → [n]α
18 scaninc ⊙ e⊙ [a1, ..., an] =
19   [a1, a1⊙a2, ..., a1⊙...⊙an]
20 scanexc ⊙ e⊙ [a1, ..., an] =
21   [e⊙, a1, ..., a1⊙...⊙an-1]
22
23 sgmscan : ∀ n. (α → α → α) → α →
24   [n]i32 → [n]α → [n]α
25 sgmscaninc ⊙ e⊙
26   [..., 1, 0, ..., 0, 1, ...]
27   [..., a1k, a2k, ..., ank, a1k+1, ...] =
28   [..., a1k, ..., a1k⊙...⊙ank, a1k+1, ...]
29
30 scatter: ∀ n, m. [n]α → [m]i32 → [m]α → [n]α
31 scatter [a0, a1, a2, a3, ..., an-1]
32   [2, -1, 0, 3]
33   [b0, b1, b2, b3] =
34   [b2, a1, b0, b3, ..., an-1]

```

Listing 1. Data-Parallel Operators Semantics

```

1  let valuate (ins : Instrument) : real =
2  let (w, h) = f1(ins)
3  let Qs = replicate w 0.0
4  let Qs[w/2+1] = f2(ins)
5  let αs = replicate h 0.0
6  let αi = f3(ins)
7  let αs[0] = αi
8  let (_, αs) =
9    loop (Qs, αi, αs) for i < h-1 do
10     let Qs' = map (λ j →
11       let q0 = Qs[j]
12         let q1 = if j > 0 then Qs[j-1] else 1.
13         let q2 = if j < w-1 then Qs[j+1] else 1.
14         in g1(i, j, αi, q0, q1, q2)
15       ) (iota w)
16     let αv = reduce (+) 0.0 Qs'
17     let αi' = g2(αv, ins)
18     let αs[i+1] = αi'
19     in (Qs', αi', αs)
20 let Cs = replicate w 100.0
21 let Cs =
22   loop (Cs) for ii < h-1 do
23     let i = h - 2 - ii
24     let αi = αs[i]
25     in map (λ j →
26       let c0 = Cs[j]
27       let c1 = if j > 0 then Cs[j-1] else 1.
28       let c2 = if j < w-1 then Cs[j+1] else 1.
29       in g3(i, j, αi, c0, c1, c2)
30     ) (iota w)
31 in Cs[w/2+1]
32
33 let main(portfolio:[] Instrument) =
34   map valuate portfolio

```

Listing 2. Nested-Parallel Implementation.

## 2.2 Simplified Nested Data-Parallel Specification

Listing 2 sketches the (simplified) implementation of the pricing algorithm, which nevertheless accurately captures the nested-parallel structure. The `main` function (at the bottom) receives a portfolio of instruments and performs a valuation of each by an embarrassingly-parallel `map` operation, that can be easily distributed across threads, GPUs or nodes.

Function `valuate` receives an instrument data as an argument and computes its price approximation. Computation starts by determining the width `w` and height `h` of the trinomial tree (at line 2), and by initializing arrays `Qs` of size `w` (lines 3-4) and `αs` of size `h` (lines 5-7). The two sequential loops of indices `i` and `ii` implement the forward and backward tree propagation.

The first loop fills in the values of an array `αs`: First, the `map` operation of length `w` (at lines 10-15) computes each element across width at the current time step (height) level in the tree, i.e., `Qs'[j]`, by aggregating the three different values belonging to the previous time step level, i.e., `Qs[j-1]`, `Qs[j]`, `Qs[j+1]`. Note that only the current and previous time step elements across width levels—rather than the entire tree—are manifested in memory by means of arrays `Qs'` and `Qs`.) The newly created array `Qs'` is then summed up—by means of the `reduce` operation at line 16—and provides the value of `αs[i+1]`. Note that both parallel operations occur inside the outer `map` operation, which is applied to the whole portfolio, thus giving raise to nested parallelism. Finally, the resulted

values of  $Qs'$ ,  $\alpha_i'$  and  $\alpha s$  are bound to the loop-variant variables  $Qs$ ,  $\alpha_i$  and  $\alpha s$  for the next iteration.

The second loop traverses the tree backwards, from the maturity to the present date, and at each step it computes the prices associated to the current width level by a similar map operation. The price of the instrument today is at the root of a tree, corresponding to  $Cs[w/2+1]$  (after the loop).

### 3 GPU-FLAT: FLATTENING TWO-LEVEL PARALLELISM

This section demonstrates how the GPU-FLAT implementation is derived from the nested-parallel code of Listing 2. While we keep the discussion here intuitive and specific to the trinomial-pricing algorithm, Section 4 formalizes the transformation by means of a set of inference (rewrite) rules, which can be integrated in the repertoire of a compiler.

Essentially, GPU-FLAT utilizes both levels of parallelism, which allows to simultaneously optimize (i) the temporal locality and (ii) both levels of divergence. The idea is to first sort the options in decreasing order of their heights—thus optimizing the divergence of the forward/backward traversal loops—and then to (bin-)pack the input (instruments) into bins, such that the summed widths of their trees do not exceed the thread-block size—we choose the maximal size 1024—which is considered the capacity of the bin. The two parallel levels—of the instruments in a bin, and of the inner parallelism inside `valuate` function—are then combined (flattened) and mapped at the thread-block level, while the parallelism across bins is mapped on the CUDA grid.

On the one hand, this implicitly optimizes the width-level of thread divergence, because the flatten parallelism has roughly the size of the thread-block. On the other hand, temporal locality is also optimized because the data created by inner-parallel operations (inside `valuate`)—such as the arrays  $Qs$  and  $Cs$ —is maintained in fast (scratchpad/shared) memory.<sup>3</sup> The downsides are that the flattening transformation introduces instructional overhead, and shared-memory/register pressure.

#### 3.1 Flat-Parallel Version in Fast Memory

Listing 3 shows the code that is obtained from applying the flattening transformation: the `bin` array corresponds to a batch of  $q$  instruments—whose summed widths is less than the thread-block size—and the result is an array of length  $q$  of real numbers denoting the prices of those instruments at current time. The flat code is obtained by distributing the (outer) `map` operation—i.e., over the  $q$  instruments of the bin—around each `let` statement of the original `valuate` function shown in Listing 2. In essence, distributing the `map` (i) across a scalar statement results in a `map` of size  $q$ , and (ii) across a parallel operation of size  $width$  results in a parallel operation of size  $\sum_{k=0}^{q-1} width_k$ , which is padded to thread-block size. For brevity and clarity of exposition, our discussion ignores the complications related to (i) padding parallel arrays to thread-block size and to (ii) using the non-trivial offsets in the arrays of all instruments and of  $\alpha s$  corresponding to the sub-arrays of the current block—these are tedious but straightforward to add.

Listing 3 starts by computing the widths and heights,<sup>4</sup> of the trees of the  $q$  instruments (line 2). For demonstration, we assume that  $q=2$ , and the widths and heights are  $ws=[2, 4]$  and  $hs=[4, 3]$ . Lines 3-11 compute three helper arrays (`flag`, `outinds` and `inninds`) that are used in the code transformation. The first array `flag` is the flag component in the flat-representation of an irregular array of shape  $ws$ , such as  $Qss$ . We recall that the flag arrays is required by the segmented scan operations. An irregular array of shape  $[2, 4]$  has two rows of lengths 2 and 4, respectively, and its flag array marks with 1 the start of each subarray and has otherwise 0 elements. It follows that

<sup>3</sup> The array  $\alpha s$  is maintained as before in global memory—because it is not guaranteed to fit in shared memory—as well as it is padded and transposed at block level to optimize coalescing and memory footprint, as before.

<sup>4</sup> In practice the widths and heights are precomputed by an inspector which is sliced out of the original code. This is because, in the preliminary steps, the instruments are first sorted after their heights in order to optimize the divergence of the sequential loops, and then their widths, that is necessary for packing instruments into bins.

```

1  let valuatebin(q: i32, bin: [q]Instrument) : [q]real =
2    let (ws,hs) = map f1 bin
3    let Bw = scanexc (+) 0 ws
4    let lenflat = Bw[q-1] + ws[q-1]
5    let tmp = map2 (λs b → if s == 0 then -1 else b) ws Bw
6    let flag = scatter (replicate lenflat 0) tmp (replicate q 1)
7    let tmp = scaninc (+) 0 flag
8    let outinds = map (λx → x-1) tmp
9    -- map (λw → iota w) ws
10   let tmp = map (λf → 1-f) flag
11   let inninds = sgmscaninc (+) 0 flag tmp
12   -- map(λw → replicate w 0) ws
13   let Qss = replicate lenflat 0.0
14   -- map2 (λ Qs w → Qs[w/2+1] = f2( ins ) ) Qss ws
15   let tmp_i = map2(λ b w → b + w/2 + 1) Bw ws
16   let tmp_v = map f2 bin
17   let Qss = scatter Qss tmp_i tmp_v
18   -- init regular (transposed) hmax×q matrix assT
19   let hmax = reduce max 0 hs
20   let αis = map f3 bin
21   let assT = scatter (replicate (hmax+q) 0.0) (iota q) αis
22
23   -- map-loop interchange; loop count padded to hmax-1
24   let (_,_,assT) = loop(Qss,αis,assT)
25     for i < hmax-1 do
26       -- map2 (λis αi → map (...) is) inninds αis
27       let Qss' = map2 (λ j oi → let (b,h) = (Bw[oi], hs[oi])
28         in if i ≥ h-1 then Qss[b+j]
29         else let q0 = Qss[b+j]
30           let q1 = if j > 0 then Qss[b+j-1] else 1.
31           let q2 = if j < w-1 then Qss[b+j+1] else 1.
32           in g1( i, j, αis[oi], q0, q1, q2)
33         ) inninds outinds
34       -- map (reduce (+) 0) Qss'
35       let scQs = sgmscaninc (+) 0.0 flag Qss'
36       let αvs = map2 (λ b w → scQs[b+w-1]) Bw ws
37       -- map(λα → α[i+1] = g2(...)) ass
38       let tmp_i = map2 (λ h k → if i ≥ h-1 then -1 else (i+1)*q + k) hs (iota q)
39       let αis' = map2 g2 αvs bin
40       let assT = scatter assT tmp_i αis'
41       in (Qss', αis', assT)
42   let Css = replicate lenflat 100. ... -- second loop is not shown (similar)

```

Listing 3. Flat-Parallel Implementation.

we expect `flag` to be equal to  $[1, 0, 1, 0, 0, 0]$ . This is computed by applying an exclusive scan on `ws`, resulting in  $B_w = [0, 2]$ , then by computing the total number of elements  $len_{flat} = 2 + 4 = 6$ , and finally, by the `scatter` operation at line 6 that writes 1s at the indices in  $B_w = [0, 2]$  into an array of  $len_{flat} = 6$  0s; hence  $flag = [1, 0, 1, 0, 0, 0]$ , as expected.

The second array `outinds` records, for each of the width entries associated with an instrument, the index of that instrument in the current bin. Thus, we expect  $out_{inds} = [0, 0, 1, 1, 1, 1]$ . This is achieved by (inclusive) scanning the flag array, which results in  $[1, 1, 2, 2, 2, 2]$ , and by subtracting 1 from each obtained element (lines 7-8).

The final array `inninds` is the expansion of `iota w` across the `q` widths, hence we expect  $inn_{inds} = [0, 1, 0, 1, 2, 3]$ . This is achieved at lines 10-11 by negating the flag array, resulting in  $[0, 1, 0, 1, 1, 1]$ , and applying an inclusive segmented scan on the result under the flag array `flag`, i.e., scanning independently the two logical rows of two and four elements, respectively.



Lines 12-17 in listing 3 are the flattening across  $q$  instruments of the lines 3-4 in listing 2—which initializes  $Q_s$  elements to zeroes and sets index  $w/2+1$  to value  $f2(\text{opt})$ . The zero-initialization of  $Q_s$  is translated to a replicate of length  $\text{len}_{flat}$ —i.e., the summed widths of the  $q$  instruments—and the update at index  $w/2+1$  is translated to a scatter on the expanded  $Q_{ss}$  in which

- the updated indices are computed at line 15 by summing the offset in  $Q_{ss}$  of each instrument, denoted by  $b \in B_w$ , with  $w/2+1$ , where  $w \in ws$ , and
- the updated values are the result of mapping  $f2$  on the  $q$  instruments at line 16.

The initialization of  $\alpha_{ss}$ —the expansion of  $\alpha_s$ —is simply obtained by padding each row to the maximal height  $h_{max}$  of the  $q$  instruments—hence total length is  $q \times h_{max}$ —and by using a scatter to overwrite the first entry in every row with the result of calling  $f3$ . This is implemented in lines 19-21, except that  $\alpha_{ss}^T$ —the transpose of  $\alpha_{ss}$ —is used in order to achieve coalesced access to global memory. Next, the forward loop is padded to count  $h_{max}$ , the outer map of length  $q$  is interchanged inside the loop, and the outer-map distribution continues on the loop-body statements.

Lines 27-33 correspond to flattening the map that is applied to  $\text{iota } w$  to compute array  $Q_s'$  in Listing 2, lines 10-15. Since the flattened code corresponds to applying the original map simultaneously to all entries of all  $q$  instruments, it is translated to a map2 over  $\text{inn}_{inds}$  and  $\text{out}_{inds}$ :

- $\text{inn}_{inds}$  is precisely the expansion of  $(\text{iota } w)$  across the  $q$  instruments, hence  $j$  takes the same values as in Listing 2;
- $\text{out}_{inds}$  is used to access values that are the same across the width threads assigned to process the current instrument, but are needed by each thread—we recall that the  $\text{out}_{inds}$  values record the index of each instrument in each of the width entries associated with it. For example,  $\text{out}_{inds}$  is used to (indirectly) index into length- $q$  arrays  $B_w$ ,  $hs$  and  $\alpha\_is$  in order to compute the start offset  $b$  into array  $Q_{ss}$ , the height  $h$  and the  $\alpha$  value corresponding to the current instrument, respectively.
- The body of the mapped function is protected by an `if` condition ( $i \geq h-1$ ) that checks that the tree traversal has not already terminated for the current instrument, because the loop count is padded to maximal value  $h_{max}$ . If so, then the input value of  $Q_{ss}$  is directly returned.

The code between lines 35-36 is the flattening across all  $q$  instruments of the (original) reduce at line 16 in Listing 2, which sums up the values of array  $Q_s'$ . This is implemented by first performing an inclusive segmented scan on the expanded array  $Q_{ss}'$ , which by definition, computes the  $q$  corresponding sums in the last entry of each logical subarray of the result  $scQ_s$ . Then these last entries are extracted by a map operation of length  $q$ ; the index of the last entry of the  $i^{th}$  subarray is  $B_w[i] + ws[i] - 1$ , because  $B_w$  and  $ws$  record the offset and the size of each subarray, respectively.

Finally, lines 38-40 implement the expansion of the update to  $\alpha_s[i+1]$  at line 18 in Listing 2. This is translated to a scatter that updates  $\alpha_{ss}^T$  at the  $q$  flat-indices belonging to row  $i+1$  (stored in  $\text{tmp}_i$ ) with the values  $\text{tmp}_v$  obtained by applying  $g2$  to all  $\alpha\_vs$  and batched instruments. Note that if the loop index  $i$  is greater or equal than the logical loop count  $h-1$  then the return index is  $-1$ , hence the update is ignored.

Similar ideas apply for the translation of the backward loop, which is not shown. Our CUDA implementation of GPU-FLAT fuses aggressively the inner-parallel operations and reuses shared memory buffers whenever possible: e.g.,  $Q_{ss}$ ,  $Q_{ss}'$ ,  $C_{ss}$ ,  $C_{ss}'$  use the same buffer. Arrays of size  $q$  are typically stored in the shared memory (since they save register space), and arrays  $\text{out}_{inds}$  and  $\text{inn}_{inds}$  are held in registers. The shortcomings of GPU-FLAT are that (i) it introduces instructional overhead—i.e., the code is more complex than the nested-parallel one, (ii) it introduces significant

register pressure <sup>5</sup> and that (iii) the parallel operations of size  $q$  underutilize the block-level parallelism, which is typically much larger than  $q$ .

#### 4 INFERENCE RULES FOR TWO-LEVEL FLATTENING

This section provides (the gist of) the formal inference rules that implement the two-level flattening transformation, which is used in Section 3 to derive the GPU-FLAT code version. The transformation is applied on code exhibiting exactly two-levels of parallelism, in which the outer level is a map of length  $q$ , whose (lambda) body respects the following restrictions:

- all inner-parallel operations—such as `map`, `reduce`, `scan`, `scatter`—have the same length, denoted  $w$ , for the same outer-map iteration, albeit  $w$  may vary across iterations;
- there exists at most one (inner) sequential loop whose count is variant to the outer map, and such a loop is not contained in an inner parallel operation;
- most computation is carried out by the inner-parallel operations, i.e., all sequential loops should contain inner parallelism;
- for simplicity of exposition, we assume that all inner arrays are one dimensional—this can be achieved by a pre-processing step that flattens the array indexing.

We also assume that the lengths of the parallel operations for each of the  $q$  iterations of the outer map, have been pre-computed by slicing out the computation of  $w$  into a simpler map construct. Figure 4 shows the code for computing the helper arrays, discussed in Section 3, where  $q$  denotes the number of outer-map iterations packed in the current bin, and  $shp$  denotes the parallel lengths in the current bin. We assume that this code has been already inserted, at the beginning of the translation. (To note, the code allows empty parallel lengths, i.e., elements of  $shp$  are allowed to be zero.) We briefly recount the rationale for each helper variable:

- $len_{flat}$  is the sum of the inner-parallel lengths of the  $q$  packed outer-map iterations, i.e., the size of the flat-parallel operation;
- $B_{shp}$  is an array of length  $q$ , which records the start index in the flat-parallel array representation of each of the  $q$  logical sub-arrays (segments);
- $inn_{inds}$  is an array of length  $len_{flat}$  that records the inner iteration space of each inner-parallel construct. For example, in figure 4, the first sub-array has length 3, therefore the first three elements of  $inn_{inds}$  are  $0, 1, 2$ ; the last sub-array has length 2, therefore the last two elements of  $inn_{inds}$  are  $0, 1$ .
- $out_{inds}$  is an array of length  $len_{flat}$ , which records the index of the outer-map iterations for each of the inner-parallel elements. For example, the first sub-array has length 3, therefore the first three elements of  $out_{inds}$  are  $0, 0, 0$ ; the third sub-array has length 2, therefore the last two elements of  $out_{inds}$  are  $2, 2$ . The rationale is that scalar variables that are variant to the outer map are expanded in the transformed code to length- $q$  arrays, and accesses to these scalars from inner map operations is translated by using  $out_{inds}$  to indirectly index into these expanded arrays.
- $flag$  is the flag array (also of length  $len_{flat}$ ) that is used by segmented scan, and which records with an one the start of each sub-array and has zero elements otherwise.

##### 4.1 The Translation Context $\Sigma$

The context of the translation—denoted  $\Sigma$  and summarized in Figure 1—is represented by a record containing the following fields:

<sup>5</sup> Nvidia nvcc compiler reports that 74 – 76 registers per thread are used by default and a speedup of up to 1.66× is achieved by limiting the number of registers to 32.



```

1  let helper (q:int) (shp: [q]int) :
2      (int,[q]int,[ ]int,[ ]int,[ ]int) =
3      -- Assume shp = [3,1,2]
4      let inds = scanexc (+) 0 shp
5      let Bshp = map2(λs k→ if (s <= 0)
6                      then -1 else k
7                      ) shp inds
8      let lenflat = Bshp[q-1] + shp[q-1]
9      -- Bshp : i32 = [0,3,4], lenflat = 6
10     let flag = scatter (replicate len 0)
11                    Bshp (replicate q 1)
12     -- flag : [lenflat]i32 = [1,0,0,1,1,0]
13
14     let inninds = sgmscanexc (+) 0 flag
15                    (replicate lenflat 1)
16     -- inninds : [lenflat]i32 = [0,1,2,0,0,1]
17     let tmps = scaninc (+) 0 flag
18     let outinds = map (-1) tmps
19     -- outinds : [lenflat]i32 = [0,0,0,1,2,2]
20
21     in (lenflat, Bshp, inninds, outinds, flag)
    
```

Listing 4. Helper Function for Flattening

$\mathcal{V}^o, \mathcal{V}^t$ : set of variable names of the original and target program, respectively.

$\Sigma = \langle \mathcal{H}, \mathcal{V}_{inv}^{map}, \mathcal{A}_{var}^{par}, \mathcal{S}_{var}, \mathcal{A}_{var}^{seq}, \mathcal{L} \rangle$

$\mathcal{H} = \langle w, q, i, o, shp, len_{flat}, B_{shp}, inn_{inds}, out_{inds}, flags \rangle$

$\mathcal{H} \in \mathcal{V}^o \times \mathcal{V}^t \times \dots \times \mathcal{V}^t$

$\mathcal{V}_{inv}^{map} \in Set(\mathcal{V}^{o/t})$

$\mathcal{A}_{var}^{par} \in \mathcal{V}^o \mapsto \mathcal{V}^t$

$\mathcal{S}_{var} \in \mathcal{V}^o \mapsto \mathcal{V}^t$

$\mathcal{A}_{var}^{seq} \in \mathcal{V}^o \mapsto \mathcal{V}^t \times \mathcal{V}^t$

$\mathcal{L} = \langle j, n^o, n_{max}^t \rangle \in \mathcal{V}^{o/t} \times \mathcal{V}^o \times \mathcal{V}^t$

 Fig. 1. The Structure of the Transformation Context  $\Sigma$ .

- $\mathcal{H}$  is a record containing helper variables for translation:
  - w is the original-program scalar variable w containing the length of the inner-parallel operation;
  - q is the number of outer-map iterations packed in the current bin;
  - i, o are the scalar variables used by the translation as formal arguments of the lambda body of a flat map, and that take values from  $inn_{inds}$  and  $out_{inds}$ , respectively.
- The other variables—shp,  $len_{flat}$ ,  $B_{shp}$ ,  $out_{inds}$ ,  $inn_{inds}$  and flag—have been discussed before, and are assumed available in the translation process.
- $\mathcal{V}_{inv}^{map}$  represents the set of original program variables that are invariant to the outer map (free variables). The translation directly reuses these names and may also insert variables of the transformed program in this set.
- $\mathcal{A}_{var}^{par}$  is a finite map that binds each (inner) parallel array—of symbolic length w—in the original code to their corresponding expansion (across q iterations) in the transformed code; these array are stored in fast (shared) memory.
- $\mathcal{A}_{var}^{seq}$  is a finite map that binds each sequential array—i.e., that are computed sequentially and are assumed to be of length other than w—in the original code to their expansion in the transformed code. In this case, the array expansion is performed by padding to the maximal size of the q subarray, and the expanded arrays are stored in global memory. The lookup returns the expanded array name together with its maximal (padded) row/segment length.
- $\mathcal{S}_{var}$  is a finite map that binds each scalar variable in the original program that is variant within the outer map to its corresponding expanded-array variable in the translation.
- $\mathcal{L}$  is a record containing the information of the sequential loop that has irregular count across the q-packed iterations of the outer map:
  - field j denotes the original-loop index, whose name is reused in the translation;
  - field  $n^o$  denotes the original-program variable storing the loop count. The lookup  $ns^t = \mathcal{S}_{var}(n^o)$  must succeed inside such a loop, and the result  $ns^t$  is the array expansion of the scalar  $n^o$  in the translated program;

- field  $n_{max}^t$  denotes the transformed-program variable recording the maximal count across the  $q$  packed iterations of the outer map—i.e.,  $n_{max}^t = \text{reduce max } 0 \text{ ns}^t$ ;
- Whenever the translation encounters such a loop, this information is updated; otherwise, outside such loops, field  $n^o$  is set to a dummy value that fails the  $S_{var}$  lookup.

The translation process assumes that the initial context have the fields  $\mathcal{H}$  and  $\mathcal{V}_{inv}^{map}$  already filled, while the other are unset. The translation requires that the input program is normalized to A-normal form: **let** bindings can be seen as a block of statements followed by a sequence of result variables, where the statement can be:

- unary/binary operations in the form of three-address code;
- a parallel operation, whose lambda body is also normalized;
- sequential loops whose body is also normalized and the loop initializers are variables.

For readability and clarity of exposition:

- The resulted program is not normalized and variables are not necessarily uniquely named.
- The translation rules use the fields of record  $\Sigma.\mathcal{H}$  freely—i.e., omitting the  $\Sigma.\mathcal{H}$  prefix.
- The discussion ignores the complications related to padding and to the non-trivial offsets at which the sequential arrays of the current bin are located inside flat (global-memory) arrays that extends across all bins; these details are tedious but straightforward to add.

## 4.2 Formalizing the Translation By Inference Rules

The core inference (rewrite) rules of our translation are formalized in figures 2 and 3. The rules allow inferences of the form

- $\Sigma \vdash_{out} e \Rightarrow e'$ , which are read “in context  $\Sigma$ , the source expression  $e$  appearing outside the lambda function of an inner-parallel map operation can be translated into the target expression  $e'$ .”
- $\Sigma \vdash_{inn} e \Rightarrow e'$ , which are read “in context  $\Sigma$ , the source expression  $e$  appearing inside the lambda function of an inner-parallel map operation can be translated into the target expression  $e'$ .” In this case, the variables  $\mathcal{H}.o$  and  $\mathcal{H}.i$  are assumed available in the declaration of the enclosing map lambda function, and are taking values from  $\mathcal{H}.out_{inds}$  and  $\mathcal{H}.inn_{inds}$ .

In an inference rule, the part below the line specifies the translation (the conclusion), i.e., this source expression is translated to this target expression. The part above the line contains the premises necessary for the translation to fire, including (i) generation of fresh names, (ii) lookup operations into finite maps (that need to succeed for the rule to fire), (iii) recursive application of inference rules, (iv) creation of new contexts, and (v) abbreviating sub-expression so that they fit in the space of the conclusion.

We discuss the inference rules  $\vdash_{out}$  from Figure 2:

**G0** refers to a binary operation  $\odot$  between two scalar variables, in which one  $v_1^o$  is invariant to the outer map and the other  $v_2^o$  is variant. The array-expanded translation of the latter,  $vs_2^t$ , is obtained after a successful lookup in  $\Sigma.S_{var}$ . The bound expression is translated to a **map** that applies  $\odot$  to the invariant  $v_1^o$  and each element of  $vs_2^t$ . The result is stored in fresh variable  $vs^t$ , a new context  $\Sigma'$  is created by extending  $\Sigma.V_{var}$  with the new binding  $v^o \mapsto vs^t$ , and the let-body expression  $e^o$  is recursively translated in the new context  $\Sigma'$ . If both operands are invariant, then the original binding is left unmodified. If both are variant, **map2** is used to combine elements from both expanded arrays.

**G1** refers to translating an array-indexing expression  $Q^o[ind^o]$ , where  $ind^o$  is outer-map variant and  $Q^o$  is a parallel array. The expanded arrays of the translation  $Qs^t$  and  $inds^t$  are looked up, and the indexing expression is translated to a map that takes values  $k \in \{0, \dots, q-1\}$  and in

Specialized-Flattening Rules for code outside inner-map constructs

$$\Sigma \vdash_{out} e^o \Rightarrow e^t$$

$$\frac{v_1^o \in \Sigma.\mathcal{V}_{inv}^{map}, \quad v_2^o \notin \Sigma.\mathcal{V}_{inv}^{map}, \quad vs_2^t = \Sigma.\mathcal{S}_{var}(v_2^o), \quad vs^t \text{ fresh name}, \quad \Sigma' = \Sigma \text{ with } \{\mathcal{S}_{var} = \Sigma.\mathcal{S}_{var} \cup \{v^o \mapsto vs^t\}\}}{\Sigma \vdash_{out} \text{let } v^o = v_1^o \odot v_2^o \text{ in } e^o \Rightarrow \text{let } vs^t = \text{map } (\lambda v_2^t \rightarrow v_1^o \odot v_2^t) \text{ vs}_2^t \text{ in } e^t} \quad (G0)$$

$$\frac{Qs^t = \Sigma.\mathcal{A}_{var}^{par}(Q^o), \quad inds^t = \Sigma.\mathcal{S}_{var}(ind^o) \quad (\text{i.e., both lookups succeed}), \quad k, vs^t \text{ fresh name}, \quad e_{map}^t = \text{map } (\lambda k \rightarrow Qs^t[B_{shp}[k] + inds^t[k]]) \text{ (iota } q), \quad \Sigma' = \Sigma \text{ with } \{\mathcal{S}_{var} = \Sigma.\mathcal{S}_{var} \cup \{v^o \mapsto vs^t\}\}}{\Sigma \vdash_{out} \text{let } v^o = Q^o[ind^o] \text{ in } e^o \Rightarrow \text{let } vs^t = e_{map}^t \text{ in } e^t} \quad (G1)$$

$$\frac{Qs^t = \Sigma.\mathcal{A}_{var}^{par}(Q^o), \quad inds^t = \Sigma.\mathcal{S}_{var}(ind^o), \quad vs^t = \Sigma.\mathcal{S}_{var}(v^o), \quad \Sigma' \vdash_{out} e^o \Rightarrow e^t, \quad j = \Sigma.\mathcal{L}.j, \quad ns^t = \Sigma.\mathcal{S}_{var}(\Sigma.\mathcal{L}.n^o), \quad (\text{i.e., inside a loop}), \quad k^t \text{ fresh name}, \quad e_{inds}^t = \text{map2 } (\lambda k \rightarrow \text{if } j < ns^t[k] \text{ then } inds^t[k] + B_{shp}[k] \text{ else } -1) \text{ (iota } q)}{\Sigma \vdash_{out} \text{let } Q^o[ind^o] = v^o \text{ in } e^o \Rightarrow \text{let } Qs^t = \text{scatter } Qs^t \ e_{inds}^t \text{ vs}^t \text{ in } e^t} \quad (G2)$$

$$\frac{n = \Sigma.\mathcal{H}.w, \quad v^o \notin \mathcal{V}_{inv}^{map}, \quad vs^t = \Sigma.\mathcal{S}_{var}(v^o), \quad e_{rep}^t = \text{map}(\lambda o \rightarrow vs^t[o]) \text{ out}_{inds}, \quad Qs^t \text{ fresh name}, \quad \Sigma' = \Sigma \text{ with } \{\mathcal{A}_{var}^{par} = \Sigma.\mathcal{A}_{var}^{par} \cup \{Q^o \mapsto Qs^t\}\}}{\Sigma \vdash_{out} \text{let } Q^o = \text{replicate } n \ v^o \text{ in } e^o \Rightarrow \text{let } Qs^t = e_{rep}^t \text{ in } e^t} \quad (G3)$$

$$\frac{m^o \notin \Sigma.\mathcal{V}_{inv}^{map}, \quad m^o \neq \Sigma.\mathcal{H}.w, \quad ms^t = \Sigma.\mathcal{S}_{var}(m^o), \quad v^o \notin \Sigma.\mathcal{V}_{inv}^{map}, \quad vs^t = \Sigma.\mathcal{S}_{var}(v^o), \quad \Sigma' = \Sigma \text{ with } \{\mathcal{V}_{inv}^{map} = \Sigma.\mathcal{V}_{inv}^{map} \cup \{m_{max}\}, \quad \mathcal{A}_{var}^{seq} = \Sigma.\mathcal{A}_{var}^{seq} \cup \{X^o \mapsto (Xs^t, m_{max})\}\}}{\begin{aligned} k, m_{max} \text{ fresh names, } \quad e_{red}^t &= \text{reduce max } 0 \ ms^t, \quad \Sigma' \vdash_{out} e^o \Rightarrow e^t, \\ e_{map}^t &= \text{map } (\lambda k \rightarrow vs^t[k/m_{max}]) \text{ (iota } (q * m_{max})) \end{aligned}} \quad (G4)$$

$$\Sigma \vdash_{out} \text{let } X^o = \text{replicate } m^o \ v^o \text{ in } e^o \Rightarrow \text{let } m_{max} = e_{red}^t \text{ in let } Xs^t = e_{map}^t \text{ in } e^t$$

$$\frac{\Sigma.\mathcal{H}.w = \text{length } Q^o, \quad Qs^t = \Sigma.\mathcal{A}_{var}^{par}(Q^o), \quad vs^t, Qs_{scn} \text{ fresh name}, \quad \Sigma' = \Sigma \text{ with } \{\mathcal{S}_{var} = \Sigma.\mathcal{S}_{var} \cup \{v^o \mapsto vs^t\}\}}{\begin{aligned} e_{scn}^{sgm} &= \text{segscan } (\odot) \ \theta_{\odot} \ \text{flag } Qs^t, \\ e_{map}^{pck} &= \text{map2 } (\lambda b \ s \rightarrow \text{if } s == 0 \ \text{then } \theta_{\odot} \ \text{else } Qs_{scn}[b + s - 1]) \ B_{shp} \ \text{shp} \end{aligned}} \quad (G5)$$

$$\Sigma \vdash_{out} \text{let } v^o = \text{reduce } (\odot) \ \theta_{\odot} \ Q^o \text{ in } e^o \Rightarrow \text{let } Qs_{scn} = e_{scn}^{sgm} \text{ in let } vs^t = e_{map}^{pck} \text{ in } e^t$$

$$\frac{\Sigma.\mathcal{H}.w = \text{length } Q^o, \quad Qs^t = \Sigma.\mathcal{A}_{var}^{par}(Q^o), \quad Xs^t \text{ fresh name}, \quad \Sigma' = \Sigma \text{ with } \{\mathcal{A}_{var}^{par} = \Sigma.\mathcal{A}_{var}^{par} \cup \{X^o \mapsto Xs^t\}\}}{\begin{aligned} e_{scn}^{sgm} &= \text{segscan } (\odot) \ \theta_{\odot} \ \text{flag } Qs^t \\ \Sigma \vdash_{out} \text{let } X^o = \text{scan } (\odot) \ \theta_{\odot} \ Q^o \text{ in } e^o &\Rightarrow \text{let } Xs^t = e_{scn}^{sgm} \text{ in } e^t \end{aligned}} \quad (G6)$$

$$\frac{\Sigma.\mathcal{H}.w = \text{length } Q^o, \quad Qs^t = \Sigma.\mathcal{A}_{var}^{par}(Q^o), \quad flgs_Q^t = \Sigma.\mathcal{A}_{var}^{par}(flgs_Q^o), \quad Xs^t, \text{ fresh name}, \quad \Sigma' = \Sigma \text{ with } \{\mathcal{A}_{var}^{par} = \Sigma.\mathcal{A}_{var}^{par} \cup \{X^o \mapsto Xs^t\}\}}{\begin{aligned} e_{scn}^{sgm} &= \text{segscan } (\odot) \ \theta_{\odot} \ flgs_Q^t \ Qs^t, \quad \Sigma' \vdash_{out} e^o \Rightarrow e^t \\ \Sigma \vdash_{out} \text{let } X^o = \text{segscan } (\odot) \ \theta_{\odot} \ flgs_Q^o \ Q^o \text{ in } e^o &\Rightarrow \text{let } Xs^t = e_{scn}^{sgm} \text{ in } e^t \end{aligned}} \quad (G7)$$

Fig. 2. Flattening rules for the code outside inner-parallel constructs (assumed to have all the same length).

which  $Qs^t$  is indexed by  $inds^t[k]$  to which we add  $B_{shp}[k]$ —the start offset of each parallel sub-array. Finally, the context is extended with the result array, and the let-body expression is recursively translated (into  $e^t$ ). Indexing a sequential array is similar: the translation and its row-padded size are looked up  $(Ys^t, m_{max}) = \Sigma.\mathcal{A}_{var}^{seq}(Y^o)$  and the indexing operation inside the map becomes:  $Ys^t[m_{max}*k + inds^t[k]]$ .

- G2** refers to an in-place update to an element of a parallel array  $Q^o[ind^o]$ , where the index is variant. The update occurs inside a loop of index  $j$  and variant count  $n^o$ , which is translated to array  $ns^t$ —i.e., please note that the rule fires only if the lookup  $\Sigma.S_{var}(\Sigma.L.n^o)$  succeeds. This case is translated by a **scatter** operation in which the to-be-updated indices into  $Qs^t$  are computed similarly to **G1**, except that they are guarded by condition  $if\ j < ns^t[k]$ . If this condition does not hold, then the current sub-array has logically finished the execution of the loop, and index  $-1$  is returned instead, which prompts the **scatter** to ignore the updates to such sub-arrays. For updating in-place a sequential array, the indexing is computed similarly to **G1**, i.e.,  $m_{max}*k + inds^t[k]$ , and the guard is also inserted.
- G3** refers to a **replicate**-based initialization of a parallel array—i.e., the length of the original array is  $w$ —where the replicated variable  $v^o$  is variant, and has the array-expanded translation  $vs^t$ . This is translated by a **map** in which the corresponding value for the elements of a sub-array is taken by indirectly accessing  $vs^t$  with the sub-array indices stored in helper array  $out_{inds}$ , i.e.,  $vs^t[o]$ , where  $o \in out_{inds}$ . Finally, the context is extended with the result array ( $Q^o \mapsto Qs^t$ ), and the let-body expression  $e^o$  is recursively translated to  $e^t$ .
- G4** is similar to **G3** except that the **replicate** initializes a sequential array, i.e., of length  $m^o$  different than  $w$ . The translation of  $m^o$  is the array  $ms^t$ , and its maximal element  $m_{max}$  is computed by a **reduce**. The result array  $Xs^t$  is padded to have logically  $q$  rows of length  $m_{max}$ , and the corresponding values of sub-array elements are selected from  $vs^t$  by using the regular indexing  $k / m_{max}$  where  $k \in \{0, \dots, q*m_{max}-1\}$ .
- G5** refers to flattening a reduce operation—on the original parallel array  $Q^o$  of size  $w$ —nested inside the outer map. This is translated by (i) performing a segmented scan on the translated array  $Qs^t$ , and (ii) by selecting the last element of each sub-array (segment) by means of a map (see  $e_{map}^{pck}$ ). The index of the last element is  $b+s-1$ , where  $b \in B_{shp}$  represents the start position of each sub-array, and  $s \in shp$  represents the length of each sub-array. The  $B_{shp}$ ,  $shp$ , and flag array (for segmented scan) are taken from  $\Sigma.H$ .
- G6** refers to flattening a scan operation on an original parallel array  $Q^o$  of size  $w$ . This is translated (by definition) to a segmented scan where the flag array is taken from  $\Sigma.H$ .
- G7** refers to a segmented scan on an original parallel array  $Q^o$  of size  $w$ . This is translated to a segmented scan, where the flag array is the translation of the original flag array—i.e.,  $flgs_Q^t = \Sigma.\mathcal{A}_{var}^{par}(flgs_Q^o)$ —because the original flag array is necessarily a parallel array (length  $w$ ).

We now discuss the inference rules  $\vdash_{out}$  from Figure 3:

- G8** refers to an inner-map operation, necessarily of length  $w$ , which is applied directly to parallel arrays. The rule normalizes/rewrites the input map into one that is applied to one array: the iteration space, which is given by  $iota\ w = [0, \dots, w-1]$ . Then, essentially it relies on rule **G9** to perform the translation.
- G9** refers to a map operation applied (only) to  $iota\ w$ , that occurs inside a sequential loop of index  $j$  and variant loop count  $n^o$ , which is translated to array  $ns^t$ . The original map is translated to a **map2** operating on helper arrays  $inn_{inds}$  and  $out_{inds}$ , and the original lambda expression  $e_k^o$  is translated by the  $\vdash_{inn}$  inference rules to  $e_k^t$ . These rules require that  $\mathcal{H}.i$  and  $\mathcal{H}.o$  are available. To this extent  $\mathcal{H}.i$  is set to the original formal argument of the lambda  $k$ —which takes values from  $iota\ w$ , hence it preserves semantics since in the

Specialized-Flattening Rules for code outside inner **map** constructs

$$\Sigma \vdash_{out} e^o \Rightarrow e^t$$

$$\frac{\Sigma.\mathcal{H}.w = \text{length } A^o, \quad k \text{ fresh name,}}{\Sigma \vdash_{out} \text{let } Q^o = \text{map } (\lambda k \rightarrow \text{let } a = A[k] \text{ in } e_a^o) (\text{iota } w) \text{ in } e_{let}^o \Rightarrow e_{all}^t} \quad (G8)$$

$$\Sigma \vdash_{out} \text{let } Q^o = \text{map } (\lambda a \rightarrow e_a^o) A^o \text{ in } e_{let}^o \Rightarrow e_{all}^t$$

$$\frac{\begin{array}{l} \Sigma' = \Sigma \text{ with } \{\mathcal{H}.i = k\}, \quad \Sigma' \vdash_{inn} e_k^o \Rightarrow e_k^t, \quad ns^t = \Sigma.\mathcal{S}_{var}(\Sigma.\mathcal{L}.n^o), \quad j = \Sigma.\mathcal{L}.j, \\ e_{map}^t = \text{map2 } (\lambda k \ o \rightarrow \text{if } j < ns^t[o] \text{ then } e_k^t \text{ else dummy}) \text{ inn}_{inds} \text{ out}_{inds}, \\ Qs^t \text{ fresh name, } \quad \Sigma'' = \Sigma \text{ with } \{\mathcal{A}_{var}^{par} = \Sigma.\mathcal{A}_{var}^{par} \cup \{Q^o \mapsto Qs^t\}\}, \quad \Sigma'' \vdash_{out} e^o \Rightarrow e^t \end{array}}{\Sigma \vdash_{out} \text{let } Q^o = \text{map } (\lambda k \rightarrow e_k^o) (\text{iota } w) \text{ in } e^o \Rightarrow \text{let } Qs^t = e_{map}^t \text{ in } e^t} \quad (G9)$$

$$\frac{\begin{array}{l} \Sigma.\mathcal{H}.w = \text{length } Q^o, \quad Qs_0^t = \Sigma.\mathcal{A}_{var}^{par}(Q_0^o), \quad ns^t = \Sigma.\mathcal{S}_{var}(n^o), \\ Qs^t, n_{max} \text{ fresh names, } \quad \Sigma' = \Sigma \text{ with } \{\mathcal{A}_{var}^{par} = \Sigma.\mathcal{A}_{var}^{par} \cup \{Q^o \mapsto Qs^t\}\} \\ \Sigma'' = \Sigma' \text{ with } \{\mathcal{L} = \langle j, n^o, n_{max} \rangle, \mathcal{V}_{inv}^{map} = \Sigma.\mathcal{V}_{inv}^{map} \cup \{n_{max}\}\} \\ e_{red}^t = \text{reduce max } 0 \ ns^t, \quad \Sigma'' \vdash_{out} e_{body}^o \Rightarrow e_0^t, \quad e_0^t = \text{let } Qs_r^t = e_r^t \text{ in } Qs_r^t \\ e_{inds}^t = \text{map2 } (\lambda \ o \ i \rightarrow \text{if } j < ns^t[o] \text{ then } i \text{ else } -1) \text{ out}_{inds} (\text{iota } len_{flat}) \\ e_{body}^t = \text{let } Qs_r^t = e_r^t \text{ in scatter } Qs^t \ e_{inds}^t \ Qs_r^t \end{array}}{\Sigma \vdash_{out} \text{loop } (Q^o) = (Q_0^o) \text{ for } j < n^o \text{ do } e_{body}^o \Rightarrow \\ \text{let } n_{max} = e_{red}^t \text{ in loop } (Qs^t) = (Qs_0^t) \text{ for } j < n_{max} \text{ do } e_{body}^t} \quad (G10)$$

Specialized-Flattening Rules for code inside inner-map constructs

$$\Sigma \vdash_{inn} e^o \Rightarrow e^t$$

$$\frac{v^o \text{ has scalar type, } \quad e_k^t = \begin{cases} vs^t[o], & \text{if } vs^t = \Sigma.\mathcal{S}_{var}(v^o) \\ v^o, & \text{otherwise} \end{cases}}{\Sigma \vdash_{inn} v^o \Rightarrow e_k^t} \quad (G11)$$

$$\frac{Qs^t = \Sigma.\mathcal{A}_{var}^{par}(Q^o), \quad \Sigma \vdash_{inn} v_{ind}^o \Rightarrow e_{ind}^t}{\Sigma \vdash_{inn} Q^o[v_{ind}^o] \Rightarrow Qs^t[B_{shp}[o] + e_{ind}^t]} \quad (G12)$$

$$\frac{(\mathcal{X}^t, h_{max}) = \Sigma.\mathcal{A}_{var}^{seq}(\mathcal{X}^o), \quad \Sigma \vdash_{inn} v_{ind}^o \Rightarrow e_{ind}^t}{\Sigma \vdash_{inn} \mathcal{X}^o[v_{ind}^o] \Rightarrow \mathcal{X}^t[h_{max} * o + e_{ind}^t]} \quad (G13)$$

Fig. 3. Flattening rules for the code outside (G8-G10) and inside (G11-13) inner-map constructs.

translation it feeds from  $inn_{inds}$ —and  $\mathcal{H}.o$  is used as the second formal argument of **map2**'s lambda, feeding as expected from  $out_{inds}$ . Finally, an **if** guard is inserted to return a dummy value for the sub-arrays that have logically finished executing their loops (similar to **G2**).

**G10** refers to the treatment of a loop whose count  $n^o$  is variant to the outer map. The (original) loop-variant variable is assumed to be a parallel array  $Q^o$  whose size is assumed invariant to the loop. The count of the translated loop is padded to  $n_{max}$ : the maximal value of  $ns^t$ , which is the array-expansion translation of scalar  $n^o$ . A fresh variable is created for the translation of  $Q^o$ , denoted  $Qs^t$ , and the context is updated with the new binding  $Q^o \mapsto Qs^t$  and the loop information. The body of the original loop  $e_{body}^o$  is recursively translated in the newly created context  $\Sigma''$  and its result is assumed to be variable  $Qs_r^t$ . This variable contains

the correct loop result for the sub-arrays that have not yet finished the loop execution, but contains dummy results for the ones that have already finished execution—see for example rule **G9**. The latter values are stored in  $Qs^t$ , and an extra step is necessary to put together the correct values for all subarrays. The **map2** operation computes the indices of all elements belonging to loop-active sub-arrays and  $-1$  for the rest. These indices are fed to a **scatter** operation which updates  $Qs^t$  at the positive indices to the values of  $Qs^t$  (and preserves the other values). The **scatter** operation adds negligible overhead in practice because it can be fused in most cases with the parallel computation of  $Qs^t$  and with the computation of indices, hence neither needs to be manifested in memory. The case when the loop-variant variable is a scalar requires a similar scatter. The case when the loop variant symbols is a sequential array, the scatter operation is not necessary, because such array can only be updated in place, and rule **G2** already ensures that loop-inactive sub-arrays are not updated.

We finally discuss the inference rules  $\vdash_{inn}$  from Figure 3, which assume that the code of interest is inside the lambda-expression of an inner **map** operation:

- G11** refers to the translation of a scalar variable  $v^o$ . If the scalar is variant to the outer map—i.e., lookup  $vs^t = \Sigma.S_{var}(v^o)$  succeeds—than the translation is  $vs^t[o]$ , where  $o$  takes values from  $out_{inds}$  in the enclosing map. Otherwise,  $v^o$  is used as it is because it is either invariant to the outer map or it is a variable local to the lambda body of the inner map.
- G12** refers to indexing into an original parallel array, where the index is stored in scalar variable  $v_{ind}^o$ . This is translated by selecting from the translated array the element at index  $B_{shp}[o] + ind^t$ , where  $ind^t$  is the translation of  $v_{ind}^o$  by rule **G11**, and  $B_{shp}[o]$  stores the start offset of the current sub-array—both  $B_{shp}$  and  $o$  are taken from  $\Sigma.H$ .
- G13** refers to indexing into a sequential array  $X^o$  at index  $v_{ind}^o$ . The translation  $Xs^t$  and the maximal row length  $h_{max}$  are obtained by a lookup in  $\Sigma.A_{var}^{seq}$ , and the start offset of the current sub-array is computed by  $h_{max} * o$ .

## 5 COMPILER RELATED WORK AND CONCLUSIONS

Our implementation draws inspiration from a number of compiler techniques. The GPU-FLAT version builds on the flattening transformation [Blelloch 1990; Blelloch and Greiner 1996; Blelloch et al. 1994], which maps irregular nested parallelism into a sequence of flat-parallel ones, and has been also implemented for GPU execution [Bergstrom and Reppy 2012].

There are two key differences here: The first one is that flattening pushes all sequential recurrences outside the parallel code, and it introduces many prefix-sum operations that are executed in global memory and thus limit performance gains. Instead, we bin-pack inner parallelism at Cuda-block level so that temporarily locality can be efficiently exploited by maintaining and reusing arrays in/from shared (scratchpad) memory.

The second difference is more subtle and refers to the fact that the traditional flattening transformation replicates variables that bound in the outer map operation but are free in the inner parallel constructs. This may lead to memory explosion, which might prevent the use of shared memory, which is a scarce resource. In comparison, our procedure does not expand such variables, but instead indirectly accesses them by means of auxiliary arrays, such as  $out_{inds}$ ,  $inn_{inds}$ ,  $B_w$  in Listing 3. The latter can be seen as part of the shape representation of an irregular array of arrays, which is reused between similarly-shaped arrays, such as  $Qss$  and  $Css$ . Furthermore, such auxiliary arrays are invariant to the sequential loop, hence they are created once and the overhead is amortized across the execution of many loop iterations.

The pricing algorithm that motivated our flattening approach is part of a class of applications, whose best parallelization strategy goes against the common-wisdom of always sequentializing



the parallelism in excess of what the hardware can support, because this reduces the per-thread memory footprint and allows to better exploit locality of reference. Another example that benefits from a similar optimization strategy has been studied in the remote-sensing field, where GPU acceleration enabled to apply at continental scale an algorithm that analyses satellite data to detect landscape changes such as deforestation [Gieseke et al. 2020]—but there, the flattening step is straightforward.

Our techniques for optimizing the two-level divergence were inspired by data reordering transformations aimed at improving locality of reference by dynamically (and systematically) reorganizing the array elements in the order in which they are accessed in loop iterations [Ding and Kennedy 1999; Strout et al. 2003, 2018]. Such transformation typically use an inspector-executor model, which has also been used (i) to compute sufficient conditions for loop parallelization [Oancea and Rauchwerger 2013, 2015; Rus et al. 2003] or (ii) to restructure a partially-parallel loop into a sequence of parallel waves [Rauchwerger et al. 1995]. In the context of software thread-level speculation [Dang et al. 2002; Rauchwerger and Padua 1999]—which is another technique that extracts partial loop parallelism at runtime—similar work was aimed at optimizing at runtime the layout of the speculative storage [Oancea and Mycroft 2008], or at optimizing communication overheads by optimistic execution of remote calls [Oancea et al. 2005] in a distributed-component framework [Chicha et al. 2004; Oancea and Watt 2005]. Finally, other work optimizes locality by a memory-centric approach that prefers executions of work-items from the current memory partition and delays the others [Oancea et al. 2009].

However, we are not aware of any compiler framework that is able to derive the GPU-FLAT version of the code. For example the code transformations of the polyhedral Pluto compiler [Bondhugula et al. 2008; Verdoolaege et al. 2013] prioritize optimization of locality of references over the degree of parallelism, but it does not support the flattening transformation (based on scans). The data-parallel language Futhark [Elsman et al. 2018; Henriksen et al. 2020, 2017] supports flattening, but only of regular parallelism, i.e., the parallel sizes of inner operations cannot vary across different iterations of an outer map. Flattening of irregular can be performed manually in Futhark, but it results in code that maintains all intermediate arrays in global memory, and is significantly slower than GPU-FLAT. While Futhark’s incremental-flattening transformation [Henriksen et al. 2019] can create code versions that utilize shared memory, it does so only when all inner-parallel constructs have the same size, and this is not the case of GPU-FLAT, which has some parallel operation of size  $len_{flat}$  and others of size  $q$ .<sup>6</sup>

The work presented in this document may provide useful insights into how to integrate such a technique into the repertoire of a compiler.

## REFERENCES

- Christian Andreetta, Vivien Bégot, Jost Berthold, Martin Elsman, Fritz Henglein, Troels Henriksen, Maj-Britt Nordfang, and Cosmin E. Oancea. 2016. FinPar: A Parallel Financial Benchmark. *ACM Trans. Archit. Code Optim.* 13, 2, Article 18 (2016).
- Lars Bergstrom and John Reppy. 2012. Nested Data-parallelism on the GPU. *ICFP '12: Proceedings of the 17th ACM SIGPLAN International Conference on Functional Programming* 47, 9 (2012), 247–258.
- Guy E. Blelloch. 1989. Scans as Primitive Parallel Operations. *Computers, IEEE Transactions* 38, 11 (1989), 1526–1538.
- Guy E Blelloch. 1990. *Vector models for data-parallel computing*. Vol. 75. MIT press Cambridge.
- Guy E. Blelloch and John Greiner. 1996. A Provable Time and Space Efficient Implementation of NESL. In *Procs. of ACM SIGPLAN International Conference on Functional Programming (ICFP '96)*. ACM, 213–225.
- Guy E Blelloch, Jonathan C Hardwick, Jay Sipelstein, Marco Zagha, and Siddhartha Chatterjee. 1994. Implementation of a Portable Nested Data-Parallel Language. *Journal of parallel and distributed computing* 21, 1 (1994), 4–14.

<sup>6</sup> In Futhark defense, a set of other financial algorithms have been reported to have GPU-efficient Futhark implementations [Andreetta et al. 2016; Henriksen et al. 2018; Oancea et al. 2012; Pawlak et al. 2019].

- Uday Bondhugula, Albert Hartono, J. Ramanujam, and P. Sadayappan. 2008. A Practical Automatic Polyhedral Parallelizer and Locality Optimizer. In *Proceedings of the 29th ACM SIGPLAN Conference on Programming Language Design and Implementation (PLDI '08)*. ACM, New York, NY, USA, 101–113. <http://doi.acm.org/10.1145/1375581.1375595>
- Y. Chicha, M. Lloyd, C. Oancea, and S. M. Watt. 2004. Parametric Polymorphism for Computer Algebra Software Components. In *Proc. 6th International Symposium on Symbolic and Numeric Algorithms for Scientific Computing (SYNASC '04)*. Mirton Publishing House, 119–130.
- Francis Dang, Hao Yu, and Lawrence Rauchwerger. 2002. The R-LRPD Test: Speculative Parallelization of Partially Parallel Loops. In *International Parallel and Distributed Processing Symposium (PDPS)*. 20–29.
- Chen Ding and Ken Kennedy. 1999. Improving Cache Performance in Dynamic Applications Through Data and Computation Reorganization at Run Time. In *Procs. of the ACM SIGPLAN Conference on Programming Language Design and Implementation (PLDI '99)*. ACM, 229–241.
- Martin Elsmann, Troels Henriksen, Danil Annenkov, and Cosmin E. Oancea. 2018. Static Interpretation of Higher-order Modules in Futhark: Functional GPU Programming in the Large. *Proc. ACM Program. Lang.* 2, ICFP, Article 97 (July 2018), 30 pages. <https://doi.org/10.1145/3236792>
- Alexandros V. Gerbessiotis. 2004. Architecture independent parallel binomial tree option price valuations. *Parallel Comput.* 30, 2 (2004).
- Fabian Gieseke, Sabina Rosca, Troels Henriksen, Jan Verbesselt, and Cosmin E. Oancea. 2020. Massively-Parallel Change Detection for Satellite Time Series Data with Missing Values. In *2020 IEEE 36th International Conference on Data Engineering (ICDE)*. 385–396.
- Scott Grauer-Gray, William Killian, Robert Searles, and John Cavazos. 2013. Accelerating financial applications on the GPU. In *Procs. of GPGPU-6*. ACM, New York, NY, USA, 127–136.
- Troels Henriksen, Martin Elsmann, and Cosmin E. Oancea. 2018. Modular Acceleration: Tricky Cases of Functional High-performance Computing. In *Proceedings of the 7th ACM SIGPLAN International Workshop on Functional High-Performance Computing (FHPC 2018)*. ACM, New York, NY, USA, 10–21. <https://doi.org/10.1145/3264738.3264740>
- Troels Henriksen, Sune Hellfritsch, Ponnuswamy Sadayappan, and Cosmin Oancea. 2020. Compiling Generalized Histograms for GPU. In *Proceedings of the International Conference for High Performance Computing, Networking, Storage and Analysis (SC '20)*. IEEE Press, Article 97, 14 pages.
- Troels Henriksen, Niels G. W. Serup, Martin Elsmann, Fritz Henglein, and Cosmin E. Oancea. 2017. Futhark: Purely Functional GPU-programming with Nested Parallelism and In-place Array Updates. In *Procs. 38th ACM SIGPLAN Conference on Programming Language Design and Implementation (PLDI 2017)*. 556–571.
- Troels Henriksen, Frederik Thorøe, Martin Elsmann, and Cosmin E. Oancea. 2019. Incremental Flattening for Nested Data Parallelism. In *Procs. Symposium on Principles and Practice of Parallel Programming (PPoPP '19)*. ACM, New York, NY, USA, 53–67.
- Kai Huang and Ruppa K. Thulasiram. 2005. Parallel algorithm for pricing American Asian options with multi-dimensional assets. In *19th International Symposium on High Performance Computing Systems and Applications (HPCS'05)*. IEEE, Washington, DC, USA, 177–185.
- John Hull and Alan White. 1994. Numerical Procedures for Implementing Term Structure Models I: Single-Factor Models. *The Journal of Derivatives* 2, 1 (1994), 7–16.
- John Hull and Alan White. 1996. Using Hull-White Interest Rate Trees. *The Journal of Derivatives* 3, 3 (1996), 26–36.
- Cosmin E. Oancea, Christian Andreaetta, Jost Berthold, Alain Frisch, and Fritz Henglein. 2012. Financial Software on GPUs: Between Haskell and Fortran. In *Procs. of Workshop on Functional High-Performance Computing (FHPC '12)*. ACM, 61–72. <https://doi.org/10.1145/2364474.2364484>
- Cosmin E. Oancea and Alan Mycroft. 2008. Set-Congruence Dynamic Analysis for Thread-Level Speculation (TLS). In *Languages and Compilers for Parallel Computing*, José Nelson Amaral (Ed.). Springer, Berlin, Heidelberg, 156–171.
- Cosmin E. Oancea, Alan Mycroft, and Stephen M. Watt. 2009. A New Approach to Parallelising Tracing Algorithms. In *Proceedings of the 2009 International Symposium on Memory Management (ISMM '09)*. Association for Computing Machinery, New York, NY, USA, 10–19. <https://doi.org/10.1145/1542431.1542434>
- Cosmin E. Oancea and Lawrence Rauchwerger. 2013. A Hybrid Approach to Proving Memory Reference Monotonicity. In *Languages and Compilers for Parallel Computing (LCPC'11) (LNCS)*, Vol. 7146. Springer, Berlin, Heidelberg, 61–75.
- Cosmin E. Oancea and Lawrence Rauchwerger. 2015. Scalable Conditional Induction Variables (CIV) Analysis. In *Proceedings of the 13th Annual IEEE/ACM International Symposium on Code Generation and Optimization (CGO '15)*. IEEE Computer Society, Washington, DC, USA, 213–224. <http://dl.acm.org/citation.cfm?id=2738600.2738627>
- Cosmin E. Oancea, Jason W. A. Selby, Mark Giesbrecht, and Stephen M. Watt. 2005. Distributed Models of Thread-Level Speculation. In *Procs. of International Conference on Parallel and Distributed Processing Techniques and Applications (PDPTA'05)*. 920–927.
- Cosmin E. Oancea and Stephen M. Watt. 2005. Parametric Polymorphism for Software Component Architectures. In *Proceedings of the 20th Annual ACM SIGPLAN Conference on Object-Oriented Programming, Systems, Languages, and*

- Applications (OOPSLA '05)*. Association for Computing Machinery, New York, NY, USA, 147–166. <https://doi.org/10.1145/1094811.1094823>
- Wojciech Michal Pawlak, Martin Elsmann, and Cosmin Eugen Oancea. 2019. A Functional Approach to Accelerating Monte Carlo based American Option Pricing. In *Procs. of 31st Symposium on Implementation and Application of Functional Languages (IFL '19)*. ACM, New York, NY, USA, 12.
- Wojciech Michal Pawlak, Marek Hlava, Martin Metaksov, and Cosmin Eugen Oancea. 2021. Acceleration of Lattice Models for Pricing Portfolios of Fixed-Income Derivatives. In *Proceedings of the 8th ACM SIGPLAN International Workshop on Libraries, Languages and Compilers for Array Programming (ARRAY '21)*. ACM, 12. <https://doi.org/10.1145/3460944.3464309>
- Lawrence Rauchwerger, Nancy Amato, and David Padua. 1995. A Scalable Method for Run Time Loop Parallelization. *Int. Journal of Par. Prog* 26 (1995), 26–6.
- L. Rauchwerger and D. Padua. 1999. The LRPD Test: Speculative Run-Time Parallelization of Loops with Privatization and Reduction Parallelization. *IEEE Trans. Parallel Distrib. System* 10(2) (1999), 160–199.
- Silvius Rus, Jay Hoeflinger, and Lawrence Rauchwerger. 2003. Hybrid Analysis: Static & Dynamic Memory Reference Analysis. *Int. Journal of Par. Prog* 31(3) (2003), 251–283.
- Michelle Mills Strout, Larry Carter, and Jeanne Ferrante. 2003. Compile-time Composition of Run-time Data and Iteration Reorderings. In *Procs. Conference on Programming Language Design and Implementation (PLDI '03)*. ACM, 91–102.
- Michelle Mills Strout, Mary Hall, and Catherine Olschanowsky. 2018. The Sparse Polyhedral Framework: Composing Compiler-Generated Inspector-Executor Code. *Proc. IEEE* 106, 11 (2018), 1921–1934.
- Sven Verdoolaege, Juan Carlos Juega, Albert Cohen, José Ignacio Gómez, Christian Tenllado, and Francky Catthoor. 2013. Polyhedral Parallel Code Generation for CUDA. *ACM Trans. Archit. Code Optim.* 9, 4, Article 54 (2013), 54:1–54:23 pages. <http://doi.acm.org/10.1145/2400682.2400713>