

Modular GPU Programming with Typed Perspectives

ANONYMOUS AUTHOR(S)

Abstract

1 Introduction

CUDA [22] is a low-level, imperative programming language for NVIDIA GPUs. These GPUs are organized into a hierarchy of compute resources. *Threads* are the basic unit of sequential execution; *blocks* are groups of threads that can cooperate through shared scratchpad memory; and the *grid* is the full collection of blocks launched for a GPU kernel. A CUDA program executes in parallel across this hierarchy, but is written from the *perspective of a single thread*.

While operations are specified per-thread, some operations are only semantically valid when executed *collectively* by a group of threads. The `__syncthreads()` primitive, for instance, synchronizes all threads within a block, and will cause a deadlock if executed by only some of the threads within that block. There are many such collective operations, including tensor core instructions, warp shuffle operations [10], and other kinds of synchronization primitives. These operations, therefore, require programmers to carefully marshal compute resources to coordinate which threads execute which lines of code. As a result, such collective operations *break the illusion* of threads executing independently.

The conceptual clash between regular statements (executed by a single thread), and collective operations (executed by a groups of threads) hinders compositional reasoning. In CUDA, function calls, like all other statements, are per-thread operations. However, a function may contain code that executes a collective operation, creating a gap between the per-thread syntax for its invocation and the cooperative semantics of its execution.

This gap creates tension between abstraction and correctness, a tension apparent even in widely used CUDA libraries that package common functionality via function interfaces. Consider the following snippet of documentation taken directly from CUB [21], a library of parallel primitives. The documentation describes the **BlockReduce** function [6], in which all threads in a block collaboratively apply a reduction operation, such as a maximum or prefix sum, over an array:

Computes a block-wide reduction for `thread0` using the specified binary reduction functor. The first `num_valid` threads each contribute one input element.

- The return value is undefined in threads other than `thread0`.
- A subsequent `__syncthreads()` threadblock barrier should be invoked after calling this method if the collective's temporary storage (e.g., `temp_storage`) is to be reused or repurposed.

Through informal documentation, CUB attempts to convey several assumptions implicit in **BlockReduce**'s implementation. First, the function is only well-defined when invoked by all threads in a block. Second, because it accesses memory shared among threads, any subsequent reuse of that storage requires synchronization to ensure all threads have completed their accesses.

In effect, CUB attempts to retrofit CUDA with information about the compute and memory requirements of collective operations. In fact, by prefixing functions with identifiers such as **Block**, CUB creates ad-hoc "namespaces" for different functions that assume similar invariants. However, without a type system to statically enforce these requirements in CUDA, correct usage of this function depends on the user carefully reading and interpreting the documentation.

In this work, we ask the following question: can we provide low-level access to collective operations while statically guaranteeing that they execute with the necessary configuration of compute resources? By reifying these configurations as type-level *perspectives*—so named because they describe the view of GPU resources against which a program is written—we find that we can. Our key insight is that GPU programmers naturally map computations onto different compute resources, and unlike CUDA, which obscures this mapping, we can expose it in a type system that tracks perspectives. This tracking enforces that collective operations are executed with the necessary resources, while still allowing users to access low-level interfaces. Further, it allows users to define and compose their own collective operations by specifying the perspective with which their code must be run.

Our approach is a departure from previous work, which attempts to resolve the mismatch between per-thread syntax and collective execution by restricting access to low-level operations. Triton [?], a tile-based, multi-dimensional array language, limits the user’s perspective to the block level, which prevents programmers from writing the highest-performance kernels. A variety of functional languages, such as Futhark [18], provide compile-time guarantees through their type systems but are machine-agnostic: they cannot expose low-level control over hardware. Some other efforts, like Descend [20], aim to provide a low-level, safe GPU systems programming language, but lack support for collective operations like tensor cores. As a result, CUDA remains the de-facto standard for writing high-performance kernels on modern GPUs.

We introduce Parallax, a new low-level GPU programming language that guarantees correct usage of hardware resources by construction. Inspired by dependency calculi [1], Parallax statically tracks the configurations of compute and memory resources with type-level perspectives. We also develop Prism, a core calculus underpinning Parallax, provide formal rules for manipulating the perspectives of both code and data, and prove *type-and-perspective safety*. This ensures that Prism is sound, and that code always has the right perspective to execute operations at run time.

A parallel goal, alongside to correctness, is performance. We incorporate modern GPU features such as tensor cores and asynchronous data movement into Parallax, and demonstrate that Parallax can achieve the same performance as hand-written, highly optimized code on an H100 and a 4070 Ti Super. Our contributions are:

- (1) Parallax, a low-level GPU programming language with type-level perspectives that track the logical grouping of compute and memory resources (Section 3);
- (2) Prism, a formal model of Parallax that tracks perspectives in both its type system and operational semantics, along with a soundness theorem guaranteeing that programs execute operations only when they have been statically proven to possess the necessary perspective.
- (3) An implementation of Parallax (Section 5) that demonstrates that it can support modern GPU features like tensor cores and asynchronous data movement, achieving performance comparable to hand-optimized CUDA implementations (Section 6).

Section 7 discusses related work, and Section 8 discusses the limitations of our approach and outlines future work.

2 Background & Motivation

Before diving into the design of Parallax, we begin with an overview of the GPU’s compute and memory hierarchies and outline the challenges posed by reasoning about them collectively.

2.1 Compute Hierarchy

In CUDA, programmers launch computations that run on thousands of threads. These threads are organized into a *compute hierarchy* that defines how work is distributed and scheduled on the GPU. At the top of this hierarchy is the *grid*, representing all threads launched as part of a single

```

99 1 int tid = threadIdx.x;
100 2 if (tid >= 0 && tid < 31){
101 3     float A[4];
102 4     float B[2];
103 5     float C[4] = { 0 };
104 6     # Populate A with unique values
105 7     for (int i = 0; i < 4; i++)
106 8         A[i] = tid * 4 + 1;
107 9     # Populate B with unique values
108 10    for (int i = 0; i < 2; i++)
109 11        B[i] = tid * 4 + 1;
110 12    # Issue a warp-level tensor-core
111 13    # operation: D = A * B + C
112 14    ↪ asm("mma.sync.aligned.m16n8k8.
15    { \"%0, %1, %2, %3\", \" /*D*/
16    { \"%4, %5, %6, %7\", \" /*A*/
17    { \"%8, %9\", \" /*B*/
18    { \"%10, %11, %12, %13\";\" /*C*/
19    : \"r\"(C[0]), \"r\"(C[1]), ...
20    : \"r\"(A[0]), \"r\"(A[1]), ...
21    : \"r\"(B[0]), \"r\"(B[1]),
22    : \"r\"(C[0]), \"r\"(C[1]), ...);}

```

Fig. 1. Invoking a warp-level tensor-core instruction in CUDA. We elide some typecasts necessary to invoke the assembly operation.

```

1 tid : int @ thread[1] = id();
2 with group(thread[32]):
3 A : float[4] @ thread[4]
4 B : float[2] @ thread[1]
5 C : float[4] @ thread[1]
6
7 # Populate A with unique values
8 for i in range(0, 4, 1):
9 A[i] = tid * 4 + 1
10 C[i] = 0
11 # Populate B with unique values
12 for i in range(0, 2, 1):
13 B[i] = tid * 4 + 1
14
15 # Issue tensor-core op.
16 intrinsic.mma(
17 A[0], A[1], A[2], A[3],
18 B[0], B[1],
19 C[0], C[1], C[2], C[3],
20 out=[A[0], A[1], A[2],
    ↪ A[3]])

```

Fig. 2. Invoking a warp-level tensor-core instruction in Parallax.

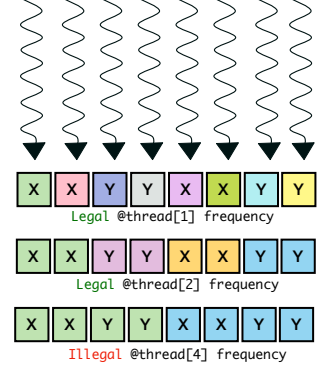


Fig. 3. The @ syntax represents the frequency at which a variable varies across threads **T1–T8**. Each @ denotes a coloring of the variable; threads that “view” the variable with the same color must observe the same value.

computation. The grid is divided into *blocks*, each containing a user-specified number of threads. A *thread*, in turn, is the finest unit of execution, and it is the machine’s basic unit of sequential control.

Users define the behavior of their computation by writing a single program which is replicated across all threads. This program controls the behavior of individual threads by reading built-in identifiers like `threadIdx.x`—to determine a thread’s position within a block—and `blockIdx.x`—to determine the block’s position within the grid—at runtime. By making control flow decisions based on these two variables, users can assign each thread to its share of the full computation.

A natural and tempting way to interpret these built-in variables is to think of them as the indices of implicit “parallel **for**-loops” surrounding the program, where each iteration is executing simultaneously. While this view is sufficient to understand CUDA’s programming model when threads do independent work, it quickly breaks down in the presence of collective operations.

Unlike most instructions, which are executed by a single thread, collective operations must be executed collaboratively by a group of threads. For example, on line 14 of Figure 1 we perform a “warp-level” tensor core operation [?], which is only *meaningful* when invoked by a collection of 32 threads—a *warp*—acting together. For the call, each thread sets up its portion of the operands, and the operation is performed *once* for the entire warp, with the results scattered across participating threads. The first 32 threads in a block are tasked with this operation; programmers mentally group these 32 threads into a collective unit, and then reflect that grouping in the program via the **if**-statement on line 2. If the condition on line 2 were instead `tid >= 0 && tid < 30`, making fewer than 32 threads reach line 14, the instruction would be undefined. To make matters worse, the restriction is not only on the number of threads executing the operation, but also on their alignment. In this case, the starting ID of the group of 32 participating threads must be aligned to a multiple of 32. So, if the condition on line 2 were instead `tid >= 1 && tid < 32`, the instruction would still be invalid, even though 32 threads would execute it.

Collective operations make reasoning about CUDA programs challenging because they force programmers to track the *convergence behavior* of threads. Specifically, programmers must reason both about how many threads reach a particular point in the program and how those threads are arranged,

accounting for alignment. This difficulty is further amplified due to two reasons. First, programs may have multiple *points* of convergence, requiring programmers to mentally track the relative ID of a thread within a logical group as that group evolves over the course of a program’s execution. Second, threads may be participating in multiple *levels* of convergence within the same program. In our example, we only considered the warp-level tensor core operation, but there are other collective operations that require convergence at different granularities like the *warp-group*-level tensor core operations, which must be issued collectively by four warps, or the block-level `__syncthreads()` synchronization primitive, which must be executed by all threads within a block.

Reasoning about collective operations is already error-prone within the context of a single function, but becomes even more difficult when reasoning interprocedurally *across* functions. A callee may assume a certain number of threads or blocks and may structure its computation, including convergence behavior, around that assumption. However, such assumptions are not visible in the callee’s function interface, which only exposes the input and output types. To invoke that function correctly, users must read documentation, or worse, read the callee’s implementation to understand its assumptions, breaking modular reasoning.

In Parallax, we materialize the programmer’s mental grouping of the compute hierarchy explicitly in the program’s source. Consider the example in Figure 2 which shows an equivalent rewrite of the CUDA program in Figure 1 using Parallax. In Figure 2, the call to the `mma` intrinsic is valid only because it is executed within a **group** of 32 threads, as made explicit on line 2. Since Parallax’s **group** construct already enforces both the size and alignment of participating threads, the validity of the `mma` operation is guaranteed by construction.

2.2 Data and the Memory Hierarchy

Data on the GPU is organized into a *memory hierarchy*, mirroring the compute hierarchy. All threads in a grid can read and write from *global memory*, where operands typically reside at the start of a computation and where results are eventually written. Each block has access to a limited amount of *shared memory*, a programmer-managed scratchpad typically used to stage repeatedly-used data. Finally, each thread maintains its own state in *registers* and *local memory*, which is used for storing the thread’s stack. Crucially, registers and local memory are per-thread data, while shared and global memory are accessible by multiple threads.

Thread-Local Data. Because registers and local memory are owned by individual threads, a variable with the same name in a program can have different values on different threads at run time. In CUDA, threads can use this to induce divergent behavior. We’ve already seen an example of this: the `tid` variable in Figure 1. Branching on this variable is an instance of data-dependent control flow; in the general case, the possibility of divergence lurks anywhere that execution branches on data.

A key challenge for Parallax is ensuring that all threads organized by other language constructs like **group** remain logically unified, even in the face of data-dependent control flow, which is necessary to compute anything interesting. To do this, Parallax must track the frequency with which values vary in space. The `@` syntax attached to each variable declaration denotes the rate at which a value changes across compute resources on the machine. For example, a `thread[1]` value can vary for every thread, while a `block[1]` can vary for every block, but not for threads within that block. Intuitively, the `@`

```

1 int x = threadIdx.x;
2 if (x >= 0 && x < 31){
3     float A[4];
4     float B[2];
5     float C[4] = { 0 };
6
7     // k is A's stride
8     A[0] = a_mem[(x/4)*k+(x%4)];
9     A[1] = a_mem[(x/4)*k+(x%4)+4];
10    A[2] = a_mem[(x/4+8)*k+(x%4)];
11    A[3] = a_mem[(x/4+8)*k+(x%4)+4];
12
13    // n is B's stride
14    B[0] = b_mem[(x%4)*n+(x/4)];
15    B[1] = b_mem[(x%4+4)*n+(x/4)];
16
17    asm("mma.sync...");
18
19    // Write back into c_mem
20    // n is C's stride
21    c_mem[(x/4)*n+2*(x%4)] = C[0];
22    c_mem[(x/4)*n+2*(x%4)+1] = C[1];
23    c_mem[(x/4+8)*n+2*(x%4)] = C[2];
24
25    ↪ c_mem[(x/4+8)*n+2*(x%4)+1]=C[3];

```

Fig. 4. Invoking a warp-level tensor core instruction in CUDA after loading data from memory.

construct “colors” a variable across threads, and any threads that share a color are required to agree on the variable’s value. Figure 3 illustrates valid and invalid colorings of a variable.

Using this information, Parallax enforces rules for reading from and writing to data. So, for example, had we attempted to introduce a condition based on `tid` within the `group`’s scope in Figure 2, Parallax would reject that program at compile time; code that diverges with low frequency cannot read from variables that change at a higher frequency. We will explain the syntax and rules of legal Parallax programs in greater detail in Section 3.

Thread-Shared Data. Shared and global memory, on the other hand, is not replicated by every thread; instead, all threads in a block see the same shared memory, while all threads on the grid see the same global memory.

CUDA does not explicitly model shared and global memory spaces, nor does it track whether allocations in a given memory space exceed device limits. Parallax, on the other hand, distinguishes these spaces and restricts shared memory to static allocations, throwing an error if an allocation exceeds device limits. We will discuss these aspects in detail in ?? For now, however, we focus on how memory converges and diverges in tandem with the compute hierarchy.

Because threads share their view of global and shared memory, these spaces are subject to convergent and divergent behavior the same way that threads are. For example, when launching a kernel on the GPU, a global memory pointer may initially belong to the grid because every thread sees the same pointer value at the start of the computation. To write to that memory, each thread computes an offset from the original base address. In doing so, the pointer effectively *diverges* across threads within the larger memory region so each thread can write independently. After these writes, the memory must *reconverge*, ensuring that all threads have completed their updates before the memory can safely be returned to its original logical owner.

Let us reconsider the tensor core example in Figure 1, this time initializing the operands of the tensor core operation from pointers into memory. In Figure 4, `A` and `B` are now populated from data in global memory, `a_mem` and `b_mem`. After the tensor core operation completes, the result is written back to `c_mem`, also in global memory. This variable `c_mem` is logically accessed from two different levels of the compute hierarchy. When entering this code, each thread in a group of 32 threads sees the same value for `c_mem`. Then, each thread within these 32 locates an offset within `c_mem` that it can write to (lines 21–24). To restore `c_mem` back to its 32 thread-level ownership, all 32 threads must first synchronize to ensure all per-thread writes have completed. This is similar to the requirement we saw documented in CUB in Section 1.

Similarly to the compute hierarchy (Section 2.1), CUDA programmers are responsible for tracking the logical owner of a piece of memory as it evolves over the course of a program, and for ensuring that appropriate synchronization occurs whenever that ownership changes.

Parallax, by contrast, makes the evolution of a memory object’s logical ownership explicit in the type system and *automatically* inserts the synchronization required to restore memory to its original ownership. In Figure 5, we show how memory is lowered through the compute hierarchy for the same tensor core operation in introduced in Figure 4.

```

1 @ccuda("device")
2 @requires(thread[32])
3 def simple_mma(
4   a: ptr(const(float)) @ thread[32],
5   b: ptr(const(float)) @ thread[32],
6   c: ptr(float) @ thread[32]):
7   x: int @ thread[1] = id()
8   with group(thread[32]):
9     A: float[4] @ thread[1]
10    B: float[2] @ thread[1]
11    C: float[4] @ thread[1]
12    # Reads do not need to be lowered.
13    A[0] = a_mem[(x/4)*k+(x%4)];
14    A[1] = a_mem[(x/4)*k+(x%4)+4];
15    A[2] = a_mem[(x/4+8)*k+(x%4)];
16    A[3] = a_mem[(x/4+8)*k+(x%4)+4];
17    B[0] = b_mem[(x%4)*n+(x/4)];
18    B[1] = b_mem[(x%4)*n+(x/4)];
19    # Skipping initialize C to 0...
20    intrinsic.mma(
21      A[0], A[1], A[2], A[3],
22      B[0], B[1],
23      C[0], C[1], C[2], C[3],
24      out=[C[0], C[1], C[2], C[3]])
25    # Must be at thread[1] to write.
26    idx: int @ thread[1] =
27      lambda ro, co: (x/4+ro)*n+2*(x%4)+co
28    with parti-
29      tion(c_mem,p=thread[32],f=idx) as
30      c_thrd:
31        c_thrd[0, 0] = C[0]
32        c_thrd[0, 1] = C[1]
33        c_thrd[0, 0] = C[2]
34        c_thrd[0, 1] = C[3]
35    # <-- Sync point inferred by name

```

Fig. 5. Invoking a warp-level tensor core instruction in Parallax after loading data from memory.

```

246 1 @ccuda("global")
247 2 # Top-level perspective bounds and shared memory usage.
248 3 @requires(grid[1], block[1], thread[32], smem=1280)
249 4 def mmaTF32NaiveKernel(A: ptr(const(float)) @ grid[1],
250 5   B: ptr(const(float)) @ grid[1],
251 6   C: ptr(float) @ grid[1],
252 7   M: int @ grid[1],
253 8   N: int @ grid[1],
254 9   K: int @ grid[1]):
255 10
256 11 # Starts out with grid[1] perspective.
257 12 with group(grid[1]):
258 13   # @ grid[1] inferred from current perspective
259 14   # Each block computes an 16 x 8 tile
260 15   num_blocks_n : const(int) = (N + 8 - 1) / 8
261 16
262 17   # id() function returns the block id
263 18   # inferred from @ block[1].
264 19   blk_row : const(int) @ block[1] = (id()/num_blocks_n)*16
265 20   ↳ |label[line][blkRow]|
266 21   blk_col : const(int) @ block[1] = (id()%num_blocks_n)*8
267 22
268 23   # Give each block an offset into C
269 24   offset = lambda x: blk_row * N + blk_col + x
270 25   with partition(C, p=block[1], f=offset) as C_blk:
271 26     with group(block[1]):
272 27       # SMEM declarations are only allowed
273 28       # with a block[1] perspective.
274 29       A_smem : shared(float[16 * 8]) @ block[1]
275 30       B_smem : shared(float[8 * 8]) @ block[1]
276 31       C_smem : shared(float[16 * 8]) @ block[1]
277 32       # Now, id() returns the thread id
278 33       idx : int @ thread[1] = id() * 4
279 34       # To write to C_smem, drop to thread[1] perspective
280 35       with partition(C_smem, p=thread[1], offset=idx) as C_th:
281 36         for i in range(0, 4, 1):
282 37           with group(thread[1]):
283 38             C_th[i] = 0
284 39
285 40   for i in range(0, K_tiles, 1):
286 41     # SYNC POINT (back edge from for loop)
287 42     for j in range(0, 4, 1):
288 43       global_row : int @ thread[1] = warp_row + row
289 44       global_col : int @ thread[1] = i * 8 + col
290 45       with partition(A_smem, p=thread[1], f=...) as A_th:
291 46         with group(thread[1]):
292 47           A_thrd[0] = A[global_row * K + global_col]
293 48
294 49     # SYNC POINT (back edge from for loop)
295 50     for j in range(0, 2, 1):
296 51       # Similar to write into A_smem ...
297 52       with partition(B_smem, p=thread[1], f=...) as B_th:
298 53         with group(thread[1]):
299 54           B_th[0] = B[global_row_b * N + global_col_b]
300 55
301 56     # Give each warp an offset into C_smem
302 57     # SYNC POINT (back edge from for loop)
303 58     with claim(C_smem, p=thread[32]) as Cs_warp:
304 59       match split(thread):
305 60         case 32:
306 61           # Call tensor core function inside
307 62           # a function interface.
308 63           simple_mma(A_smem, B_smem, Cs_warp)
309 64
310 65     # Write back final result to C
311 66     for j in range(0, 4, 1):
312 67       flat_idx_c : int @ thread[1] = id() * 4 + j
313 68       row_c : int @ thread[1] = flat_idx_c / MMA_K
314 69       col_c : int @ thread[1] = flat_idx_c % MMA_K
315 70       idx = lambda x: row_c * N + col_c + x
316 71       with partition(C_smem, p=thread[1], f=idx) as C_th:
317 72         with group(thread[1]):
318 73           C_th[0] = C_smem[row_c * MMA_N + col_c]
319 74     return

```

Fig. 6. Tensor Float 32 matrix Multiplication in Parallax.

The lowering of `c_mem` is required in this program as Parallax only permits writes when data is owned by a single thread. To lower `c_mem`, we use Parallax’s `partition` operation on line 31. The operation takes the name of the memory to partition, `c_mem`, and lowers it to a single thread, and assigns it a new name, `c_thrd`; Parallax will not allow references to the old name `c_mem` within the partition’s scope. The `partition` function also takes an indexing function, and each time `c_thrd` is accessed, this indexing function is implicitly applied. Once the `partition`’s scope ends, Parallax will insert a synchronization before the first use of the original name, so that all per-thread writes have completed.

We would like to emphasize that preventing data races is not one of Parallax’s goals. In Parallax, data races *can* occur; however, since the data is eventually synchronized before it is reused, the last-writer wins. Out-of-bounds accesses are considered undefined behavior. We believe that prior work, such as Descend [20], lays the foundation for reasoning about data-race free programs, and that a Descend-like memory model can be adapted to Parallax. Parallax is instead focused on the *interaction* between the compute and memory hierarchies, and on reasoning about them simultaneously to ensure that an operation is executed only when sufficient compute resources are available, a guarantee that prior work like Descend cannot provide.

3 The Parallax Language

Parallax is an imperative, low-level programming language designed at a level of abstraction comparable to CUDA. Unlike CUDA, Parallax’s syntax materializes the mapping of a computation onto the compute and memory hierarchy explicitly in the program source. Using this information, Parallax enforces that programs only execute collective operations with sufficient resources.

To guide our discussion about Parallax’s design, we use the program in Figure 6 as a running example. It computes a matrix multiplication between two float arrays, `A` and `B`, to produce an output

matrix **C**. In this program, each block computes a 16×8 independent tile of the output. To do so, the it first locates the tile index assigned to it (line 19-20). Next, threads in each block load corresponding rows and columns from **A** and **B** (line 41-53) into shared memory. Finally, the program invokes a warp-level, tensor-core instruction to compute the output on (line 62), requiring threads in a warp to *converge* at the point. We encapsulate this tensor-core instruction in a function, demonstrating how function composition works in Parallax. The function is the same as Figure 5.

3.1 Levels

Parallax models the machine’s compute hierarchy through *levels*. There are three levels in Parallax—**grid**, **block** and **thread**—which are organized as expected: a **grid** consists of multiple units of the **block** level, which in turn consist of multiple units of the **thread** level. Levels have an order defined over them, with **thread** < **block** < **grid**.

There are two key differences between Parallax’s levels and those in CUDA. First, Parallax does not model CUDA’s three-dimensional grid or block structure. Second, there are two commonly-used “levels”, *warp* and *warpgroup*, noticeably absent from our hierarchy. On the hardware, the units of each level are arranged in a single linear order, and the three-dimensional structures of CUDA are simply *interpretations* of this ordering, not distinct hardware resources. Similarly, warps and warpgroups are organizational constructs defined in terms of existing levels. Namely, a warp is a group of 32 threads whose first thread ID is aligned to 32. A warpgroup, which only became a meaningful construct with the release of the Hopper architecture, consists of 128 threads aligned to 128.

Rather than baking these interpretations into Parallax by adding new dimensions and levels, we let users express multi-dimensional structures and define groupings of custom sizes within Parallax.

3.2 Perspectives

Perspectives are the central concept of Parallax. A perspective represents the layer of the hierarchy from which the programmer is defining the machine’s behavior. They are what allow Parallax to determine which compute resources the programmer is controlling, whether they are available in the program’s context, and, once provided, if those resources are sufficient for a given operation.

A *perspective* is a level—**grid**, **block**, or **thread**—paired with a static constant **n**, specifying the number of units at that level. A perspective is written as **level** **n**. For example, **thread**[2] denotes a perspective of two threads, **block**[4] denotes a perspective of four blocks, and so on. Perspectives also carry *alignment information*: a perspective of size **n** is aligned to **n**. In this way, a warp is simply a desugaring of **thread**[32], and a warpgroup is a desugaring of **thread**[128].

Finally, perspectives have a partial order defined on them. We say that **level**₂[**n**₂] is *wider* than **level**₁[**n**₁], or that **level**₁[**n**₁] is *narrower* than **level**₂[**n**₂], if either:

- (1) **level**₁ < **level**₂; or,
- (2) **level**₁ = **level**₂ and **n**₂ % **n**₁ = 0.

The condition **n**₂ % **n**₁ = 0 enforces alignment, and we will see why it is required in Section 3.4.

3.3 Perspectives on Code

Code is associated with a set of perspectives, called a *perspective bound*. Perspective bounds directly correspond to the compute resources a program expects to exercise control over the course of its execution. At every point in the program, the perspective bound for that point indicates which layer of the hierarchy is being programmed, and how that layer can be destructured into narrower perspectives. For example, a line of code with a {**block**[1], **thread**[4]} perspective bound tells Parallax that the current line of code is being programmed at the **block**[1] perspective and that the **block**[1] perspective can be destructured into a some number of **thread**[4] perspectives. As a

shorthand, when we refer to a code’s perspective, we mean the broadest perspective available in its perspective bounds.

Each function starts out with a top-level perspective bound defined using the notation `@requires(grid[n1], b`. In Figure 6, the perspective bound is defined on line 2. Starting from this top-level declaration, programmers can shape the program’s current perspective using two constructs **group** and **split**.

Group. The **group** construct lets programmers shift from a broader perspective to some number of narrower perspectives contained in it. Operationally, the **group** construct does this by replicating code written from the perspective of the narrower perspective across the broader one. That is, **group** forks many copies of a narrower perspective.

Let’s consider this in context of our example. In Figure 6, execution begins at `grid[1]` on line 12. At that point, a programmer logically controls the whole grid’s behavior. To produce different output tiles, the programmer needs to be able to control the behavior of individual blocks, and they shift their perspective to `block[1]` on **todo**. Inside the `group(block[1])`’s scope, the behavior of the machine can be defined with respect to a single block, and code is replicated across all blocks in the grid.

As the reader may have already guessed, not all uses of **group** are valid. The examples in Figure 7 have no meaningful interpretation on the hardware. On line 4, the program tries to broaden its perspective to `block[1]` from a `thread[1]` perspective, which is clearly illegal. On the other hand, in the second example, the program tries to narrow its perspective to a `block[5]` from a `block[6]` perspective. While $5 < 6$, Parallax cannot evenly replicate `block[5]` perspective inside a `block[6]` perspective, and so rejects this program. Recall, from Section 3.2, that the condition on whether a perspective is narrower than another is actually a condition on divisibility, not just size.

To eliminate such malformed cases, Parallax only allows an invocation of `group(level[n])` if the current perspective bounds contains a perspective wider than `level[n]`. Once `group(level[n])` is invoked, it modifies the current perspective bounds in two key ways. First, it removes all perspective wider than `level[n]` from it. Second, it sets the widest perspective within that **group**’s scope to be `level[n]`.

Split. Unlike **group**, which is used for replication into multiple narrower perspectives (all of equal size), **split** is used for sharding the current perspective unequally. For example, Figure 8 shows a **split** from `thread[4]` perspective to one branch with `thread[2]` perspective and two with `thread[1]` perspective. The three arms of the **split** execute independently in parallel¹, as a form of unordered composition. When `split(level)` is invoked, the perspectives of its branches diverge. At the end of the **split**, they reconverge, and continue execution with the original perspective.

Use of **split** is necessary to write warp-specialized code, a programming pattern used in high-performance kernels. Another important use of **split**—masking off threads—can be found in our running example. Line 63 in Figure 6 shows a **split** that requests the first warp in the block, narrowing from `block[1]` into a single `thread[32]`. Later in Figure 8, this warp is used to execute a tensor-core operation.

Because **split** corresponds to unordered composition, it must provide each of its branches their requested perspectives simultaneously. Parallax thus checks that

```

1 # Example 1
2 with group(thread[2]):
3   # Illegal because block > thread
4   with group(block[1]):
5     pass
6 # Example 2
7 with group(block[6]):
8   # Illegal because 6 % 5 != 0
9   with group([block[5]]):
10    pass

```

Fig. 7. Illegal uses of **group**.

```

1 with group(thread[4]):
2   with split(thread):
3     case 2:
4       ...
5     case 1:
6       ...
7     case 1:
8       ...

```

Fig. 8. Simple example of **split**.

¹Despite the use of the pattern-matching **case** syntax, all branches of the **split** execute.

the cumulative sum of the perspectives requested by each branch of the **split** can be satisfied. For example, the program in Figure 9 does not type check. Finally, because perspectives enforce alignment, every branch of the **split** must also be aligned; not all **splits** whose sizes sum to less than or equal to the available units are valid. Figure 10 shows an example of this constraint being violated.

Once **split** is invoked, for each branch that requests n units, all perspectives wider than **level** $[n]$ are removed from its perspective bound, and the available units for **level** are set to exactly n .

3.4 Perspectives on Data

Section 3.2 described how a programmer can control different layers of the hierarchy by changing the perspective of code through **group** and **split**. To actually respect the structure of these constructs, Parallax must ensure that all units inside a perspective remain logically bundled together within code written from that perspective, and do not diverge. That is, Parallax must be able to reason about how control flow introduced through **for**, **while** and **if** statements effects the flow of perspectives through a program. To see why, consider a group of threads seen from **block** $[1]$ perspective, branching on a variable which has different values across those threads. In this case, different threads will take different paths, and the structure of the **block** $[1]$ perspective will be lost.

To reason about data-dependent control flow, Parallax must track which perspectives a piece of data lives at. Parallax distinguishes between two kinds of data: thread-local data (or, register and stack data), and thread-shared data (or, shared memory and global memory).

3.4.1 Thread-local data. Thread-local data is data stored in registers and stack-variables that is private to each thread. In the program source in CUDA, a thread-local that has the same name may contain different values at runtime. Indeed, to compute anything useful, threads must operate on different input values, and being able to reason about this divergent state is essential.

Instead of imposing severe restrictions on data-dependent control flow, Parallax solves this problem by understanding *how* values diverge across threads. That is, users tell Parallax which perspective a thread-local variable lives at. Intuitively, the perspective that a variable lives at tells Parallax the *frequency* at which its values change in space. For example, a **grid** $[1]$ value is the same for all threads in a grid, meanwhile a **block** $[2]$ value changes every two blocks. Crucially, this frequency tells us that a **block** $[2]$ variable is *indistinguishable* by threads within 2 blocks where the starting block ID is aligned to 2. Figure 3 shows different frequencies as different colorings of a variable.

Programmers specify the perspective that each variables live at in its declaration. A variable v of type **int** is declared at **thread** $[1]$ perspective using the syntax $v : \text{int} @ \text{thread}[1]$.² Variables are only available for reading and writing from certain perspective. The rule can be summarized as follows: Parallax allows programs to “write down” to narrower perspectives and “read up” from wider ones.

Write-down. Only values that live at wider perspectives can be written to values that live at narrower perspectives. For example, a **block** $[1]$ value can be used to write into a **thread** $[1]$ value, but a **thread** $[1]$ value cannot be written into a **block** $[1]$ value

```
1 # Example 1
2 with group(thread[4]):
3   match split(thread):
4     case 4:
5       ...
6       # Illegal! 4 + 1 > 4
7     case 1:
8       ...
```

Fig. 9. Illegal split exceeds the available perspectives.

```
1 # Example 2
2 with group(thread[3]):
3   match split(thread):
4     case 2:
5       ...
6     case 1:
7       ...
8       # Illegal! 2 + 1 <= 3,
9       # but the second branch
10      # is not aligned by 2.
```

Fig. 10. Illegal split violates alignment.

²If not explicitly annotated, Parallax infers a variable’s perspective to be the perspective of the code where it was declared.

Read-up. Reads from variables are dually constrained. Wider perspective values can only be read from values that live narrower perspectives. Figure 11 gives an example of an illegal read. Because the variable `flag` varies across threads, branching on it causes a subset of threads in the block to reach the `__syncthreads()`, causing a deadlock, despite the fact that `__syncthreads()` should always be safe in `block[1]`.

Together, the “read up” and “write down” rules ensure that information only flows from high perspective variables to low perspective ones. In Figure 6, we can see rule “read up” in action on lines 45 and 46. Meanwhile, “write down” is being used to set up variables on line 18-22.

3.4.2 Thread-Shared data. Unlike thread-local data, which is literally replicated across threads and is backed by distinct physical storage, thread-shared data consists of pointers into shared memory and global that are visible to some collection of threads. As a result, the perspective such data inhabits can evolve as the program executes.

In Parallax, there are three mechanisms for obtaining a thread-shared data. The first method is to directly allocate data residing in shared memory (Parallax assumes that all global memory allocations have been made before launching the CUDA kernel, as is typically the case with CUDA programs). The second and third are to obtain offsets into existing data by using `partition` or `claim`. These constructs mirror the `group` and `split` constructs, and are used to widen the perspective associated with a pointer.

Parallax simplifies shared memory allocation by requiring all shared memory allocations have a static size. Since shared memory is only visible to all threads in a block, a shared memory allocation is only valid with if code has a `block[1]` perspective. An example of such an allocation is shown on lines 28-30 of Figure 6. During compilation, Parallax will automatically handle the necessary pointer arithmetic to assign each allocation an appropriate offset within the shared memory space. The `@requires` annotation specifies the amount of shared memory a function expects to have available, and Parallax uses this information to ensure that a function’s allocations do not exceed its declared limit, and checks whether there is enough shared memory available at its call sites. Figure 6 declares the amount of shared memory it will require on 2. We use standard techniques to statically bound memory usage with Parallax’s type system.

Now, given an allocation into global or shared memory, Parallax introduces two constructs `partition` or `claim` that help programmers change the perspective that a pointer resides at.

Partition. The `partition` construct plays the same role for memory that `group` plays for compute. It is used to narrow the perspective of a pointer by assigning each narrow perspective a distinct offset into the original pointer

The `partition` construct takes the original name of the pointer, `base_name`, an indexing function `f`, and a target perspective `level[n]` to which the memory is being narrowed, and produces a new variable, `new_name`. Once a memory object has been partitioned, the original name referring to it goes out of scope, and the memory can only be accessed through the new name. Each use of this new name applies the indexing function `f` to compute the true offset of the access. In Figure 6, for example, we partition the `C` buffer to `block[1]` perspective on line 25 first, and then narrow it further to `thread[1]` perspective on line 78 so that we can actually write to it. When the scope of the `partition` statement ends, Parallax automatically inserts the necessary synchronization to restore the memory to its original perspective. The mechanism for inserting and optimizing these synchronization points is described in ??

At this point, it is necessary to consider how different indexing functions affect the possibility of data races. If the indexing function is injective, each narrower perspective receives a distinct offset into the underlying array, the resulting partition is free of data races. When the indexing function

```

1 flag : bool @ thread[1] = ...
2 with group(block[1]):
3   if (flag)
4     __syncthreads();

```

Fig. 11. Illegal read of `thread[1]` variable.

is not injective, multiple threads may write to the same location, introducing a potential data race. Out-of-bounds accesses are undefined behavior.

As we discussed in Section 2.2, achieving data-race freedom is *not* one of Parallax’s goals, and our guarantees continue to hold even when indexing functions are not injective. Because Parallax automatically inserts the necessary synchronization before the original memory name is restored to its higher perspective, Parallax ensures that all writes have completed before accessing the original name. The subsequent uses of that pointer will observe the value written by the last writer.

Claim. The `claim` construct is the analogue of `split` for memory. Using `claim`, the entire memory region can be moved to a narrower perspective.

Unlike `partition`, `claim` does not take an indexing function. The `claim` construct takes the original memory’s name, the target perspective to narrow to, and a new name to assign to the narrowed memory. The `new_name` is accessible only within a single, explicit `split` branch; sibling branches are not permitted to read or write from this memory. In Figure 6, we use the `claim` operation on line 63 to narrow the memory to `thread[32]` perspective.

As with `partition`, a `split` construct introduces its own scope, and the original name is inaccessible within this scope.

3.5 The `id()` Function

As opposed to exposing users to special hardware `blockIdx.x`, `threadIdx.x` variables directly, Parallax provides the `id()` function instead. The `id` function returns the relative index from a given perspective. The interpretation of the `id` function depends on both the perspective of the variable it is being written to, and the perspective that the `id` function is being invoked in. In our example in Figure 6, we use a call to `id` in `grid[1]` scope at line 18 and 19 to locate the tile that each block is ultimately charge of computing. In Figure 2, however, a call to `id` in `thread[32]` scope at line 7 will not return the global thread ID of the threads launched in the kernel, but rather the *relative* thread ID within the 32 threads, ranging from 0 to 31.

3.6 Collective Operations

With perspectives on code and data in place, we can now turn our attention to Parallax’s reasoning about collective operations.

There are two ways to access collective operations in Parallax. The first is to use intrinsics provided by the compiler. These intrinsics have built-in perspective classes attached to them, and Parallax enforces that the operation is only ever executed if code is being executed from the correct perspective. In Figure 2, the tensor-core call at line 22 is calling into intrinsics provided by our compiler. The second, more flexible approach allows users to extend these operations using the `unsafe` feature. Using unsafe code, users can inline a collective operation—often a single assembly instruction—into a Parallax function and ensure that the function’s top-level requirements, specified with `@require`, faithfully describe the perspectives required to execute that instruction. Then, Parallax will check the call to the function as any other function call, and will check whether it is being called with a compatible perspective. The rules for checking function calls are described in Section 3.8.

3.7 Asynchrony

Asynchronous data-movement works similarly to other data operations.

Users can mark storage as asynchronous with the `async` construct. As with `partition` and `claim`, invoking `async` removes the old name and exposes a fresh name that is only accessible within the scope of the `async` statement. Inside this

```
1 with async(old_name) as new_name:
2   match thread:
3     case 1:
4       cp_async(smem, new_name, 16);
```

Fig. 12. Example of an `cp.async` instruction in Parallax.

scope, Parallax restricts the new name to be used solely with the asynchronous data-movement intrinsics that the language provides. These intrinsics operate on a per-thread level, so Parallax will require users to program from the `thread[1]` perspective.

Parallax models two primary forms of asynchronous data movement: async mempy memory copies and tensor-memory accelerator instructions. We show an example of the former in Figure 12. As with the data operations in ??, Parallax is responsible for inserting the necessary synchronization before the program’s first use of the original name, ensuring that the asynchronous transfer has completed.

3.8 Function Composition

To check whether a function call is valid, Parallax needs to ensure that the call site can provide the set of perspectives that a function will attempt to use over its execution, and that the arguments passed to the function also respect the perspectives declared for its input and output types. Both requirements are represented explicitly via the `@require` construct.

As mentioned in Section 3.2, each function carries a perspective class. We call this top-level perspective class the function’s *perspective signature*. The signature specifies the top-level perspective class the function *assumes* will be provided to it. This signature represents the *minimum amount* of compute and memory resources the function must be called with. In the example above in Figure 6, line 2 defines the perspective signature. Up to this point, we have described programs written *inside* functions under the assumption that the resource signature can be satisfied. In Section 3.2, for example, we used this signature to set up the top-level perspective class for the function.

Compute perspectives. At a function’s call site, then, we need to ensure that the function’s perspective signature can be satisfied. Snap compares the current compute perspective with the one required by the function, and verifies that only functions whose perspective signatures can be satisfied are called.

Per-argument checks. In Parallax, functions have pass-by-value semantics. For arguments, we distinguish between primitive data types and pointers to memory. For primitive data types, like `int`, `bool`, and stack allocations, users can pass either arguments that live at or above the data perspective specified by the function’s signature. For pointers, we allow passing pointers that live at higher data perspective if the pointer is marked as `const`, indicating that it will only be used for reading (line 4 and 5 in Figure 2). Otherwise, we require that pointers match the exact perspective of the function signature.

4 Formalization in Prism

Having introduced the full Parallax language, we now describe Prism, a core calculus that formalizes its most fundamental aspects by statically tracking perspectives on code and data. We use Prism to argue that well-typed Parallax programs are not only type-safe, but will also never get stuck trying to execute operations for which they lack the correct perspective.

In this section, we describe Prism’s type system and operational semantics—in particular how it manages compute and data perspectives—and build up to a formal proof of type-and-perspective safety. We instrument Prism’s operational semantics with runtime perspective enforcement: its rules will get stuck if they encounter an operation for which they have the wrong perspective. This runtime enforcement means that our safety theorem can guarantee that dynamically-realized perspectives match the ones inferred by the type system.

$$\begin{array}{c}
\boxed{\Gamma \vdash^\pi e : \tau} \quad (Expression\ typing) \\
\\
\frac{x :^\pi \tau \in \Gamma}{\Gamma \vdash^\pi x : \tau} \text{ T-VAR} \quad \frac{\Gamma \vdash^{\pi'} e_1 : \tau[] \quad \Gamma \vdash^\pi e_2 : \mathbf{int} \quad \pi \leq \pi'}{\Gamma \vdash^\pi e_1[e_2] : \tau} \text{ T-ARR-ACCESS} \\
\\
\boxed{\Gamma \vdash^\pi s} \quad (Statement\ typing) \\
\\
\frac{\Gamma \vdash^{(h,n_1)} s_1 \quad \Gamma \vdash^{(h,n_2)} s_2 \quad n_1, n_2 \mathbf{align\ to\ } n}{\Gamma \vdash^{(h,n)} \mathbf{split}(n_1, n_2)\{s_1\}\{s_2\}} \text{ T-SPLIT} \\
\text{where } n_1, n_2 \mathbf{align\ to\ } n ::= (n_1 + n_2 \leq n) \text{ and } (n_1 | n) \text{ and } (n_2 | n) \text{ and } (n_2 | n_1 + n) \\
\\
\frac{\Gamma \vdash^\pi s}{\Gamma \vdash^{q \cdot \pi} \mathbf{group\ } q\ s} \text{ T-GROUP} \quad \frac{\Gamma, y : \downarrow^\pi \tau[]^l \vdash^\pi s \quad l \neq \mathbf{Local}}{\Gamma, x :^\pi \tau[]^l \vdash^\pi \mathbf{lower\ } x \text{ into } y \text{ in } s} \text{ T-LOWER} \\
\\
\frac{\Gamma, y :^{(h,n/c)} \tau[]^l \vdash^{(h,n)} s \quad c | n \quad l \neq \mathbf{Local}}{\Gamma, x :^{(h,n)} \tau[]^l \vdash^{(h,n)} \mathbf{partition\ } x \text{ into } y \text{ by } c \text{ in } s} \text{ T-PARTITION} \quad \frac{\Gamma \vdash^\pi s}{\Gamma \vdash^\pi \mathbf{destruct\ in\ } s} \text{ T-DESTRUCT} \\
\\
\frac{\Gamma, y :^{(h,n')} \tau[]^l \vdash^{(h,n')} s \quad n' \leq n \quad l \neq \mathbf{Local}}{\Gamma, x :^{(h,n)} \tau[]^l \vdash^{(h,n)} \mathbf{claim\ } x \text{ into } y \text{ at } n' \text{ in } s} \text{ T-CLAIM}
\end{array}$$

Fig. 13. Core typing rules of Prism. The typing rules presented here are a simplified selection of the full rules, which can be found in Appendix A.2.

4.1 Prism Type System

The core idea in Prism is to track, at the type-level, the program’s perspective on code and data. To achieve this, we borrow techniques from the literature on *dependency tracking* [1]; in particular, the code perspective is tracked on the typing judgment, which has the form $\Gamma \vdash^\pi e : \tau$ for expressions and $\Gamma \vdash^\pi s$ for statements. The π over the \vdash is the code perspective on e and s , and is comprised of a level and a size, the same structure as a perspective in Parallax.

The typing context also tracks the perspective at which each variable lives; data can only be read or written from a variable when its perspective is compatible with that of the code interacting with it. This requirement is made manifest in the **T-VAR** rule, which can be found in Figure 13. Observe that the π in the variable rule must match exactly between data and code; principles like “read up” and “write down” are instead encoded directly in the rules for reading and writing, like **T-ARR-ACCESS**, which views the array being read with broader perspective than the current code perspective.

Figure 13 also shows other key rules, and these rules fall into two main categories: rules for managing the code perspective and rules for managing data perspective.

4.1.1 Managing Perspectives on Code. In Parallax, there are two operations which manage code perspectives: **group** and **split**. In Prism, to better model the details of how perspectives shift throughout programs, we introduce a third construct, **destruct**.

Prism’s **group** statement directly corresponds to Parallax’s, and is checked by the **T-GROUP** rule. Given some perspective that can be divided into q equally sized narrower π s, the statement **group** $q\ s$ will check with perspective $q \cdot \pi$ provided that s itself checks with perspective π . The starting perspective is of the form $q \cdot \pi$, and encodes Parallax’s alignment requirement.

The **split** statement, meanwhile, is checked by the **T-SPLIT** rule, and functions like a binary version of the n -ary **split** construct in Parallax. It enforces the same divisibility requirements to ensure that the perspectives on code and data remain properly aligned, and then checks the two sub-statements s_1 and s_2 with the divided, narrower perspectives.

The final rule for managing code perspective is **T-DESTRUCT**, which shifts perspective down the GPU hierarchy and makes explicit in Prism programs exactly where such shifts occur. The \downarrow operation on π s “destructs” the perspective into many narrower perspectives at a lower level; $\downarrow(\mathbf{Grid}, 1) = (\mathbf{Block}, B)$ and $\downarrow(\mathbf{Block}, 1) = (\mathbf{Thread}, T)$, where B and T are parameters to a particular instantiation of Prism to describe the number of blocks per grid and threads per block.³ Because \downarrow is only defined on $(\mathbf{Grid}, 1)$ and $(\mathbf{Block}, 1)$, the rule enforces that one can only **destruct** their code perspective with exactly one grid or block.

4.1.2 Managing Perspectives on Data. The mechanism for managing data perspective mirrors that of code perspective, with each operation for data corresponding to an operation for code.

The **partition** operation is analogous to **group**-ing a code perspective. The typing rule for this operation, **T-PARTITION**, requires that the data perspective on the variable to be partitioned, x , is the same as the current perspective on code. After partitioning, a fresh variable y is introduced with a new perspective π/c , which we use as a shorthand to denote a division of a perspective π into c equally-sized narrower perspectives. Within the body of the **partition**, we disallow any references to the original variable x and continue checking with the original code perspective π ; the **partition** has no effect on the current code perspective.

Unlike **partition**, which divides up a piece of data equally among narrower perspectives, the **claim** operation views the claimed data with exactly one narrower perspective. Accordingly, Prism needs to ensure that only one branch of a **split** operation, with the appropriate π , can refer to the claimed variable. To ensure that this is the case, the **T-CLAIM** rule links the data perspective of the variable to the compute perspective of the code claiming it by changing *both* at the same time. This represents a minor difference from Parallax, which uses additional static analysis to ensure that a claimed variable is only accessed in a single **split** branch.

Lastly, the **T-LOWER** rule mirrors the **T-DESTRUCT** rule; it uses the \downarrow operator to move a variable from one level of the hierarchy to another, distributing it equally among all the narrower perspectives at that level in the same manner as **T-PARTITION**.

4.2 Prism Semantics

Having explained the key rules of the type system, we can move on to discuss Prism’s operational semantics. To reflect the fact that GPU programs execute in parallel across numerous threads, we model the semantics of Prism in the style of $[\cdot]_?$, using a two-level small step judgment. We present the key rules of this semantics in Figure 14.

The top level (i.e., device-level) judgment has just one rule: **S-PROGRAM**. This rule acts as a “frame” for the lower level (i.e., thread-level) judgment, and steps a collection of thread-ID-indexed local memories (L), a collection of block-ID-indexed shared memories (S), a global memory (Σ), and a *thread pool* to an updated collection of memories and thread pool. The thread pool maps thread and block IDs to code, intuitively representing the program being executed by each thread at the current moment. The **S-PROGRAM** non-deterministically chooses a thread ID and block ID and steps it according to

³Prism is abstracted over these B and T values, so instead of tracking perspective bounds the way that Parallax does, it only tracks the top-level perspective described in Section 3.2.

$$\begin{array}{c}
\frac{L(t), S(b), \Sigma, t, b, 0 \vdash^{(\text{Grid}, 1)} s \rightsquigarrow s' \dashv \eta', \sigma', \Sigma' \quad P(t, b) = s}{L, S, \Sigma, P \rightsquigarrow L[t \mapsto \eta'], S[b \mapsto \sigma'], \Sigma', P[(t, b) \mapsto s']} \text{S-PROGRAM} \\
\\
\frac{p < n_1 \quad n_1, n_2 \text{ align to } n \quad \eta, \sigma, \Sigma, t, b, p \vdash^{(h, n_1)} s_1 \rightsquigarrow s'_1 \dashv \eta', \sigma', \Sigma'}{\eta, \sigma, \Sigma, t, b, p \vdash^{(h, n)} \text{split}(n_1, n_2) \{s_1\} \{s_2\} \rightsquigarrow \text{split}(n_1, n_2) \{s'_1\} \{s'_2\} \dashv \eta', \sigma', \Sigma'} \text{S-SPLIT-LEFT} \\
\\
\frac{p \geq n_1 \quad p < n_1 + n_2 \quad n_1, n_2 \text{ align to } n \quad \eta, \sigma, \Sigma, t, b, p - n_1 \vdash^{(h, n_2)} s_2 \rightsquigarrow s'_2 \dashv \eta', \sigma', \Sigma'}{\eta, \sigma, \Sigma, t, b, p \vdash^{(h, n)} \text{split}(n_1, n_2) \{s_1\} \{s_2\} \rightsquigarrow \text{split}(n_1, n_2) \{s_1\} \{s'_2\} \dashv \eta', \sigma', \Sigma'} \text{S-SPLIT-RIGHT} \\
\\
\frac{\eta, \sigma, \Sigma, t, b, p \text{ mod } n \vdash^{(h, n)} s \rightsquigarrow s' \dashv \eta', \sigma', \Sigma'}{\eta, \sigma, \Sigma, t, b, p \vdash^{(h, q \cdot n)} \text{group } q \ s \rightsquigarrow \text{group } q \ s'; \dashv \eta', \sigma', \Sigma'} \text{S-GROUP} \\
\\
\frac{\eta, \sigma, \Sigma, t, b, t \text{ mod } T \vdash^{(\text{Thread}, T)} s \rightsquigarrow s' \dashv \eta', \sigma', \Sigma'}{\eta, \sigma, \Sigma, t, b, 0 \vdash^{(\text{Block}, 1)} \text{destruct in } s \rightsquigarrow \text{destruct in } s' \dashv \eta', \sigma', \Sigma'} \text{S-DESTRUCT-BLOCK} \\
\\
\frac{}{\eta, \sigma, \Sigma, t, b, p \vdash^{(\text{Block}, 1)} x := \text{alloc Shared } \tau \ n \text{ in } s \rightsquigarrow s \dashv \eta, \sigma [x \mapsto^\pi \langle x, n \rangle], \Sigma} \text{S-ALLOC-SHARED}
\end{array}$$

Fig. 14. Core semantic rules of Prism. As with the typing rules, we present only a simplified selection of the full rules, which can be found in Appendix A.3.

the thread-level judgment. This allows the semantics to model the full range of non-deterministic behavior arising from the GPU's thread scheduler.⁴

The thread-level judgment has the shape $\eta, \sigma, \Sigma, t, b, p \vdash^\pi s \rightsquigarrow s' \dashv \eta', \sigma', \Sigma'$, where

- η denotes thread-local memory,
- σ denotes shared memory,
- Σ denotes global memory,
- t denotes the thread's ID,
- b denotes the ID of the block in which the thread lives, and
- p denotes the relative position of the thread within π (the perspective ID).

Critically, notice that a π also appears on the thread-level judgment just as it does on the typing judgment. This is because the thread semantics *dynamically tracks and enforces* perspectives. The same way evaluation of a program “gets stuck” if a value does not have the right type, the semantics of Prism also get stuck if code attempts to access data or invoke commands with the wrong perspective. As an example, observe the **S-ALLOC-SHARED** rule in Figure 14, which requires a **(Block, 1)** perspective and will fail to step if encountered with a different one. This runtime perspective is present in Prism, but is erased by Parallax during compilation; we will later use it to prove that well-typed programs will always execute with the same perspective that the type system viewed them with.

The semantic rules for perspectives involve manipulating p to track which threads take which code paths when perspectives are **split** or **grouped**. Notice that in **S-PROGRAM**, the thread-level judgment always begins with perspective **(Grid, 1)**: all the perspective management rules are congruences, narrowing the perspective of further evaluation as determined by the particular rule used.

⁴In reality, the GPU's thread scheduler guarantees that the threads of a warp execute in lockstep, but modeling every thread as completely independent is both simpler and a conservative overestimate of the nondeterministic behavior of the GPU.

These rules take great care to ensure that p always describes the relative position of a compute resource within its perspective; the payoff is that Prism’s semantics can later use this p value to model the way that Parallax automatically adjusts indices into data when partitioning a data perspective. Beyond these key rules for managing perspective on code, we have modeled many of the other features of Parallax, such as asynchronous operations and thread synchronization, in Prism. To handle features like these we equip the operational semantics with additional structure, including sets of semaphores [12] for synchronization and a stack of effect handlers for modeling deferred asynchronous computations inspired by Ahman and Pretnar [3]. We have elided these details here for simplicity, but interested readers can see further details in Appendix A.3.

4.3 Soundness Theorem

Together, the type system and operational semantics allow us to prove the following syntactic soundness theorem, which says that Prism programs are type safe and do not get stuck trying to execute operations for which they lack the required perspective:

THEOREM 4.1. (Type-and-Perspective Safety). *For any program s such that $\Gamma \vdash^\pi s$, either:*

- (1) s is **skip**, or
- (2) *for any well-typed environments η, σ , and Σ , there is an s', η', σ' , and Σ' such that $\eta, \sigma, \Sigma, t, b, p \vdash^\pi s \rightsquigarrow^* s' \dashv \eta', \sigma', \Sigma'$ such that $\Gamma' \vdash^\pi s'$, where Γ' is an extension of Γ , and η', σ' , and Σ' are well-typed with respect to Γ' .*

PROOF. Via the usual progress and preservation lemmas, available in Appendix A.4. \square

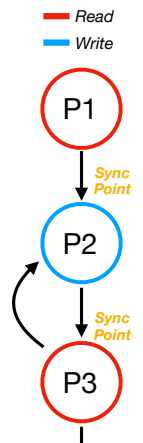
It is worth noting that this soundness theorem guarantees a syntactic safety property, not a liveness property: it does not guarantee that all threads sharing perspective π that *can* reach a program point typed with π *will* eventually do so. Indeed, in the presence of nontermination, liveness does not hold—some of the threads could **split** off and loop forever. While we believe the liveness version of this theorem holds for a terminating fragment of Prism, it is not provable with syntactic methods; the proof would require semantic techniques that are notoriously challenging and would be a research contribution [4, 15?] in and of itself. We hope to tackle this proof for Prism in future work.

5 Implementation

Parallax is implemented as an embedded language in Python. Once a program type checks, Parallax lowers it to a CUDA file. All perspective information is erased during this step, and the generated CUDA contains no run-time checks. The file can then be compiled by **nvcc**, NVIDIA’s closed-source compiler, to produce an executable. Because Parallax operates at roughly the same abstraction level as CUDA, there is a one-to-one mapping between most language constructs and their CUDA counterparts. A notable change is the addition of three parameters to each device function: the thread’s relative ID, the block’s ID, and an offset into shared memory that it can use for allocations.

Inserting Synchronization. Most of Parallax’s implementation is straightforward, but inserting synchronization points is more involved. As described in ??, once data has been partitioned, Parallax is responsible for inserting the appropriate synchronization points when the partition’s scope ends.

To determine where these points need to be inserted, Parallax constructs a *data-control-flow graph* from the program. Each node corresponds to a partition, and each edge captures program-order precedence: the parent partition must complete before the child begins. The graph can have backedges introduced through loops. In this graph, each partition is categorized as a *read* or a *write* partition. Synchronization points are inserted according to the following scheme:



- (1) If the parent partition is a write partition, a synchronization point is inserted to ensure that the current partition observes the most recent data; or
- (2) If the current partition is a write partition, a synchronization point is inserted to ensure that all preceding reads on the same memory have completed.

Figure 15 shows an example graph with the synchronization points derived from these two conditions. The inferred synchronization points in Figure 6 have also been marked on lines **TODO**.

Using this information, Parallax emits *wait* operations before a partition begins (to wait for its parent) and *arrive* operations when a partition ends (to signal completion). CUDA exposes a general *split-barrier* primitive that we use to implement these waits and arrives. For special cases, such as synchronizing an entire block or warp, Parallax instead uses primitives like `__syncthreads()` or `__syncwarp()`.

Synchronization for asynchronous data movement is handled in exactly the same way and uses the same underlying graph. CUDA allows asynchronous load operations to be associated with a split barrier, so Parallax binds each asynchronous transfer to the appropriate wait–arrive pair inferred from the graph. Certain features, such as commit-group–style synchronization, require additional reasoning, and Parallax performs further static analysis to insert the necessary synchronization around these operations.

It is worth noting that naively inserting synchronization immediately after each partition would be correct but prohibitively slow. To avoid this, Parallax applies two optimization passes. A wait-motion pass pushes waits downward toward the first use of the partitioned name, and an arrive-motion pass pulls arrives upward toward its last use, reducing unnecessary stalls while preserving correctness.

6 Evaluation

Having explained the design of Parallax, we now evaluate it in the context of four main questions:

RQ1 Can Parallax express a diversity of CUDA programs?

RQ2 Can Parallax express programs that use advanced GPU features?

RQ3 Can Parallax match the performance of existing, speed-of-light CUDA code?

Although the following question is not central to our paper’s main claims—since library design ultimately reflects choices made by programmers, and our language allows abstraction design much like any other language—we were, nevertheless, curious to explore it based on empirical and anecdotal observation:

RQ4 Can Parallax help build compositional libraries that users can use with confidence?

To perform this evaluation, we use two GPUs. The first is the NVIDIA H100 SXM5, a server-grade chip that supports tensor-core operations and a dedicated hardware copy engine, the Tensor Memory Accelerator. Notably, the H100 introduces a new *logical* level called the *warpgroup* (collection of 4-aligned warps), and we show that our programming model can accommodate this new level. Moreover, because the H100 has historically served as the primary GPU for large-scale AI training, many CUDA kernels on this hardware are already highly optimized and achieve near–speed-of-light performance, providing a rigorous baseline for comparison. To ensure our results generalize beyond the H100, we additionally test programs on a second GPU, the NVIDIA 4070 SuperTi GPU, a consumer-grade chip.

Like mentioned in Section 5, when name typechecks a program, it produces a CUDA file. We compile the CUDA file with `nvcc` version with flags . To measure performance of our benchmark, we set up the values from an input using a random number generator. We use the average running time sampled across 10 iterations, with a warm-up of 5 iterations.

6.1 RQ1: Can Parallax express a diversity of CUDA programs?

To answer this research question, we evaluate Parallax on two benchmarks: a sequence of matrix multiplication kernels that grow in complexity and a single-pass parallel prefix scan with decoupled look-back that requires programmers to think carefully about multiple points of convergence.

Matrix Multiplication. We chose matrix multiplication as our first benchmark for two main reasons. First, there exist several implementations that achieve near-peak performance, providing a strong baseline. Second, matrix multiplication allows for a range of increasingly sophisticated implementations that stress different parts of the language, making it ideal for evaluating expressiveness.

To undertake this exploration, we adapt the codebase from to implement a matrix multiplication on the `float` datatype, also known as `sgemm`, in Parallax. As this is a `float` benchmark, it *does not need* to use advanced GPU features like asynchrony or tensor cores to achieve speed-of-light performance. We will discuss these advanced features in Section 6.2.

We implement several variants of the `sgemm` benchmark, reproduced in Appendix . These include a naive baseline, a version that exploits memory coalescing, 2D blocking with shared-memory staging, 2D blocking with vectorized loads, and a warp-tiling strategy that effectively introduces an additional level of tiling from the perspective of a warp. We find that the runtime performance of our Parallax implementations closely matches that of the original versions. This is expected, since the generated code is nearly identical to the hand-written versions, aside from small differences such as hoisted expressions and a few index calculations. The downstream compiler (`nvcc`) readily inlines hoisted expressions that are used only once, so the performance-critical inner loops remain effectively unchanged.

Single-pass Parallel Prefix Scan with Decoupled Look-Back. We also implement scan, a widely used parallel primitive, in Parallax. We focus on the prefix-sum scan, which computes, for each position in an array, the sum of all elements up to that position. Prefix-sum sits in a different corner of the GPU design space from matrix multiplication: it is memory intensive, requires careful attention to the convergence behavior of different threads, and traditionally requires multiple passes over data.

We implement the single-pass parallel prefix scan with decoupled look-back, introduced by Merrill and Garland , an elegant algorithm that does not require multiple passes over the input data. In this algorithm, each block computes a local prefix sum over its assigned region of the array. Once finished, it writes the final element of its region into a global array. Each block then *looks back* to accumulate the contributions of prior blocks, allowing them to independently determine the global prefix without a full sweep over memory. This decoupled look-back mechanism lets blocks progress at different speeds while still producing correct global results.

The algorithm involves several distinct points of convergence. Within each block, work is decomposed into fine-grained thread-level and warp-level scans. After producing the local result, blocks publish their partial prefix to global memory. Finally, each block waits until enough prefix information from earlier blocks becomes available, at which point it accumulates the value and completes its section of the global scan.

We implement this full strategy in Parallax, available in Appendix . Our implementation uses the unsafe feature to implement a global-memory spinlock that lets blocks check when the preceeding block's data is ready.

6.2 RQ2: Can Parallax express programs that use advanced GPU features?

To test whether Parallax can express programs that use advanced features of modern GPUs, we write a matrix multiplication for the `bfloat16` datatype for the H100.

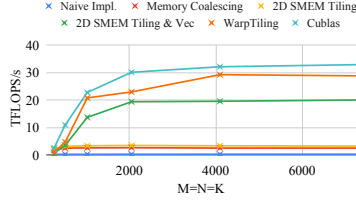


Fig. 16. Performance of **sgemm** on square matrices as matrix dimension $M=N=K$ increases.

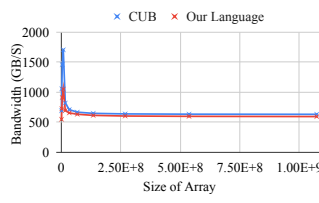


Fig. 17. Performance for prefix sum as input array size increases.

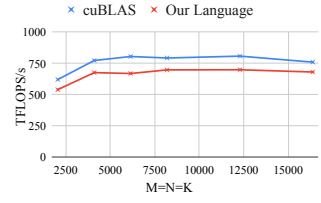


Fig. 18. Performance of **hgemm** on square matrices as matrix dimension $M=N=K$ increases.

This benchmark is an acid test of our language, as the H100 bf16 matrix multiplication pushes several language features to the extreme. To write a matrix multiplication that can hit peak throughput on an H100, we need to write a warp-specialized kernel that uses the Tensor Memory Accelerator (TMA), an asynchronous hardware copy engine that can move tiles of data at a time, and uses the warp-group-level tensor core instructions, or **wgmma**, specifically introduced for the Hopper architecture.

The implementation works as follows. First, we assign each block on our machine a logical tile of the output to compute. The block is then divided into a producer warp-group and a set of consumer warp-groups, where the producer loads data for multiple consumers, and both must signal to each other when one is done loading and the other is ready to compute. In this way, the benchmark overlaps computation with data movement by pipelining loads.

The implementation in Parallax looks different from normal CUDA code, particularly in how pipelining is expressed. Since Parallax uses *names* and subsequent **partition** or **claim** constructs to determine the synchronization each region requires, when pipelining, we cannot dynamically change the pipeline slot simply by maintaining an index variable that wraps around based on the pipeline length. Instead, each pipeline slot must be given a separate name so that Parallax can track them independently and actually overlap compute with data-movement. This leads to pipeline slots that must be individually named and forces the load logic to be effectively “unrolled,” since we can no longer iterate over a pipeline-slot index. In turn, this forces that all pipelines in Parallax be statically sized. In practice, these pipelines are statically sized in CUDA anyway because they occupy shared memory, which is a small, finite resource that must be explicitly managed.

Notably, for this benchmark, to get the **wgmma** instruction to work we did not need to add a new perspective to our language; we could simply use **thread[128]**. We did, however, need to add a special TMA asynchronous data-movement construct, since Parallax will eventually need to synchronize these transfers.

6.3 RQ3: Can Parallax match the performance of existing, speed-of-light CUDA code?

In Section 6.1 and Section 6.2, we examined programs that expressed the same computation in multiple ways, expressed operations that rely on multiple points of convergence, and used advanced GPU features. We now discuss the performance of these programs.

The performance of the matrix multiplication benchmarks is shown in , where we demonstrate that Parallax is competitive with cuBLAS [7].

For the prefix scan, we compare our performance to , the library introduced earlier in Section 1, and show that we can reach approximately % of CUB’s peak bandwidth.

Finally, we evaluate our H100 implementation and show that it is competitive with cuBLAS on square sizes, falling between %. We emphasize that this is an exacting benchmark, and achieving this level of performance requires writing large kernels with precise control over low-level features.

6.4 Case Study: Can Parallax help with library design?

Over the course of our evaluation, we found ourselves developing a small library of functions—similar in spirit to the CUB library—that we could call to execute our operations. Based on our experience, we would like to study, qualitatively, if Parallax can help programmers design libraries that they can compose with confidence.

As mentioned in Section 1, CUB occupies a unique design point in the GPU library ecosystem. Unlike many other libraries such as cuBLAS [7], cuDNN [8], and cuSPARSE [9], which provide global interfaces that users can call and configure, CUB provides a device-side library organized into different levels. It makes these levels apparent by prefixing each of its functions with **Device**, **Block**, and **Thread**. The single-scan prefix sum in Section 6.1 used these functions in Parallax already.

Let’s turn our attention to a particular CUB function—the store function—and examine how it is equivalently expressed in Parallax, and how Parallax’s type system reifies the assumptions implicit in CUB.

In CUB, the store function is implemented as a class, as shown below:

```
1 template<typename T, int BlockDimX,
2         int ItemsPerThread, BlockLoadAlgorithm Algorithm = BLOCK_LOAD_DIRECT,
3         int BlockDimY = 1, int BlockDimZ = 1>
4 class BlockStore {
```

To use it, users must first instantiate an object of this class, and then call it with shared memory.

```
1 using BlockLoad = cub::BlockLoad<int, 128, 4, BLOCK_LOAD_DIRECT>;
2 // Allocate shared memory for BlockLoad
3 __shared__ typename BlockLoad::TempStorage temp_storage;
4 int thread_data[4]; // Thread local data
5 BlockLoad(temp_storage).Load(d_data, thread_data);
```

CUB exposes a leaky abstraction, where information about the number of threads, block sizes, and other details seeps through:

- (1) The CUB documentation needs to specify the number of threads that the function can assume to be available, because within the function, each thread must locate itself in the computation and use its `threadIdx.x` accordingly. If the starting `threadIdx.x` is not 0, the function must compute its relative ID internally.
- (2) The “item per thread” design is interesting and serves two purposes. The first is related to performance: if loops have constant bounds, they can be unrolled, enabling downstream optimizations. The second is related to correctness: the function relies on the assumption that all threads call the function with an equal number of values to load.
- (3) The CUB documentation also needs to specify that `thread_data` can be data local to each thread.
- (4) Finally, the CUB documentation specifies that if shared memory is being overwritten, a `__syncthreads()` call must be made to ensure that all reads have completed.

On the other hand, the same function has a completely different interface in Parallax:

```
1 @cuda("device")
2 @require(block[1], thread[1])
3 def load(src: ptr(const(shared(int))) @ block[1],
4         item_per_thread: int @ block[1],
5         total_size: int @ block[1],
6         thread_data: ptr(int) @ thread[1]):
```

In Parallax, we are able to encapsulate code effectively, reducing the need to communicate numerous implementation details through documentation:

- (1) We do not need to pass in the number of threads at all. Whenever Parallax calls a function, it threads the relative ID through, so each function can be written locally as if it were running alone, rather than having to determine where the thread resides in the global array.
- (2) We do not need to make `item_per_thread` a template argument for *correctness*. Its frequency is set at the function signature, so Parallax will never allow a function to be called with a value living at a narrower perspective than its signature requires.
- (3) In our interface, `thread_data` is explicitly set to a thread-local value. Since it is not marked as `const`, Parallax conservatively assumes it may be written to, and enforces at compile time that only `thread[1]` values are passed in.
- (4) Finally, using our synchronization pass outlined in Section 5, a `__syncthreads()` call will be inserted automatically if `src` is going to be used for writing.

7 Related Work

Parallax builds on a rich tradition for systems languages for GPU programming and theoretical foundations for reasoning about parallel programs. However, Parallax differs from these systems in a crucial way: it treats collective operations, and therefore composability, as first-class concerns. Parallax occupies a distinct niche by pairing low-level control with safety, rather than prioritizing one at the expense of the other.

Imperative Languages for GPUs. CUDA [22], ROCm [2], and OpenCL [?] are imperative, C-style programming languages that expose low-level access to GPU hardware. These languages do not model perspectives for code or data, statically or otherwise; instead, programmers must orchestrate computation on the machine by describing how individual threads execute code. The subtle and often silent failures permitted by this model are the motivation for Parallax.

Descend [20], a new, Rust-inspired language, uses a type system to track aspects of the compute and memory hierarchy, and, is thus the closest to our work. However, Descend’s main concern is preventing data races, while ours is ensuring that collective operations execute in valid contexts. Moreover, Descend lacks support for many collective operations like tensor cores; in fact, because Descend allows threads to read their own IDs and thus induce data-dependent control flow, adding such support would require major changes to the language. It is also worth noting that while Descend does formalize a type system, it does not attempt to prove any properties about it.

Functional Languages for GPUs. Futhark [19], Accelerate [5] and Vertigo [14] are functional array languages with compilers targeting CUDA. These languages expose a high level interface, abstracting away details of the hardware entirely in exchange for stronger safety guarantees. Parallax, however, makes an intentional decision to expose low-level details of the hardware, and thus does not trade performance for safety.

Tile-Based Kernel DSLs. More recently, tile-based GPU languages—Triton [?], Pallas [?], Tilus [13], Helion [?]—offer a middle ground between high-level abstraction and low-level control. However, these languages sidestep the question of composability entirely because they restrict programmers to a single vantage point in the GPU compute hierarchy: a block (Triton, Tilus, Helion), or a warpgroup (Pallas). Gluon [?] and is a new low-level, tile-based language. However, it does not check and cannot enforce that collective operations are executed at the correct layer of the hierarchy.

Task-Based and Scheduling Languages. Languages like Cypress [?], Halide [?], Fireiron [16] and RISE/ELEVATE [17] expose a mapping specification or scheduling language that lets users modify an existing reference program through external commands. Because they operate by transforming a fixed source-of-truth rather than generating the program directly, they expose a fundamentally different programming model than Parallax.

Libraries Built on CUDA. Many GPU libraries, such as CUTLASS [?], CUB [21], and ThunderKittens [?], expose device functions that operate at various levels of the compute hierarchy. These functions are typically organized via C++ namespaces to reflect their intended usage (e.g., warp-level, block-level). However, this organization is purely conventional: it encodes hierarchy through naming rather than enforcing it with a compiler. As a result, correct use requires careful discipline from both the library implementer (to uphold naming and usage invariants) and the user (to correctly interpret them). Any mismatch or subtle misunderstanding between the intended and actual use of these functions goes unchecked by the compiler, leaving room for fragile errors.

Theoretical Foundations. The design of Prism is heavily inspired by existing work on dependency tracking [1]. Dependency tracking calculi allow type systems to track how data and code depend on each other, and have commonly been used to implement secure information flow analyses [11]. In Prism, data that lives at a narrow perspective unable flow into data living at a broader perspective, and we use dependency tracking to capture this restriction in Prism’s type system.

8 Conclusion, Limitations, and Future Work

We have presented Parallax, a new, low-level GPU language that statically guarantees safe usage of compute resources by construction, without sacrificing low-level control. Parallax introduces a new mental model for writing GPU code, which we are excited to make more expressive and ergonomic.

Currently, Parallax does not natively support explicit references to pointers; programmers must rely on built-in constructs to interact with memory. Prior work such as Descend [20], which builds on Rust’s ownership model [?], provides a promising direction for extending Parallax with more flexible pointer semantics.

We also aim to improve the ergonomics of programming with Parallax, particularly by enhancing the pipelining experience. As discussed in Section 6.2, Parallax currently requires pipeline slots to be explicitly name; we can improve this by adding language support for generating pipeline-style code.

Parallax can be extended to more architectures; in particular, it can accommodate newer GPUs—such as Blackwell—that support coarser-grained tensor-core operations. More broadly, we hope to generalize the model presented here to hierarchical compute environments, including distributed systems.

On the theoretical side, we plan to explore a terminating fragment of Prism and prove the liveness property discussed in Section 4.3: that all threads sharing perspective π eventually reach the parts of a program viewed at that π . This amounts to showing that threads sharing the same perspective execute the same code and observe the same data, which we hope to prove using logical relations, following [?] and [?].

We believe that Parallax is a promising low-level substrate that enables confident composition and can serve as a foundation for building higher-level libraries and abstractions.

References

- [1] Martín Abadi, Anindya Banerjee, Nevin Heintze, and Jon G. Riecke. 1999. A core calculus of dependency. In *Proceedings of the 26th ACM SIGPLAN-SIGACT symposium on Principles of programming languages* (San Antonio Texas USA, 1999-01). ACM, 147–160. doi:10.1145/292540.292555
- [2] Advanced Micro Devices, Inc. 2024. AMD ROCm™ Software. <https://www.amd.com/en/products/software/rocm.html>. Accessed: 2025-11-10.
- [3] Danel Ahman and Matija Pretnar. 2021. Asynchronous effects. *Proceedings of the ACM on Programming Languages* 5, POPL (Jan. 2021), 1–28. doi:10.1145/3434305
- [4] Lars Birkedal, Filip Sieczkowski, and Jacob Thamsborg. 2012. A Concurrent Logical Relation. In *Computer Science Logic (CSL ’12) - 26th International Workshop/21st Annual Conference of the EACSL (2012)*. Schloss Dagstuhl – Leibniz-Zentrum für Informatik, 107–121. doi:10.4230/LIPIcs.CSL.2012.107
- [5] Manuel M.T. Chakravarty, Gabriele Keller, Sean Lee, Trevor L. McDonell, and Vinod Grover. 2011. Accelerating Haskell array codes with multicore GPUs. In *Proceedings of the Sixth Workshop on Declarative Aspects of Multicore*

- Programming* (Austin, Texas, USA) (*DAMP '11*). Association for Computing Machinery, New York, NY, USA, 3–14. doi:10.1145/1926354.1926358
- [6] NVIDIA Corporation. 2025. cub::BlockReduce. https://nvidia.github.io/cccl/cub/api/classcub_1_1BlockReduce.html
 - [7] NVIDIA Corporation. 2025. cuBLAS. <https://docs.nvidia.com/cuda/cublas/index.html>
 - [8] NVIDIA Corporation. 2025. cuDNN. <https://docs.nvidia.com/deeplearning/cudnn/latest/>
 - [9] NVIDIA Corporation. 2025. cuSPARSE. <https://docs.nvidia.com/cuda/cusparse/index.html>
 - [10] NVIDIA Corporation. 2025. Warp Shuffle Functions. <https://docs.nvidia.com/cuda/cuda-c-programming-guide/#warp-shuffle-functions>
 - [11] Dorothy E. Denning and Peter J. Denning. 1977. Certification of programs for secure information flow. *Commun. ACM* 20, 7 (July 1977), 504–513. doi:10.1145/359636.359712
 - [12] Edsger Wybe Dijkstra. 1963. Over de sequentialiteit van procesbeschrijvingen. (1963). <https://training-ir7.tdl.org/handle/123456789/594>
 - [13] Yaoyao Ding, Bohan Hou, Xiao Zhang, Allan Lin, Tianqi Chen, Cody Yu Hao, Yida Wang, and Gennady Pekhimenko. 2025. Tilus: A Tile-Level GPGPU Programming Language for Low-Precision Computation. arXiv:2504.12984 [cs.LG] <https://arxiv.org/abs/2504.12984>
 - [14] C. Elliott. 2005. Vertigo — GPU Compiler & Embedded Language for Graphics Processors. <http://conal.net/Vertigo/>. Accessed: 2025-11-10.
 - [15] Azadeh Farzan, Zachary Kincaid, and Andreas Podelski. 2016. Proving Liveness of Parameterized Programs. In *Proceedings of the 31st Annual ACM/IEEE Symposium on Logic in Computer Science*. ACM, New York NY USA, 185–196. doi:10.1145/2933575.2935310
 - [16] Bastian Hagedorn, Archibald Samuel Elliott, Henrik Barthels, Rastislav Bodik, and Vinod Grover. 2020. Fireiron: A Data-Movement-Aware Scheduling Language for GPUs. In *Proceedings of the ACM International Conference on Parallel Architectures and Compilation Techniques* (Virtual Event, GA, USA) (*PACT '20*). Association for Computing Machinery, New York, NY, USA, 71–82. doi:10.1145/3410463.3414632
 - [17] Bastian Hagedorn, Johannes Lenfers, Thomas Kundendhler, Xueying Qin, Sergei Gorlatch, and Michel Steuwer. 2020. Achieving high-performance the functional way: a functional pearl on expressing high-performance optimizations as rewrite strategies. *Proc. ACM Program. Lang.* 4, ICFP, Article 92 (Aug. 2020), 29 pages. doi:10.1145/3408974
 - [18] Troels Henriksen, Niels G. W. Serup, Martin Elsman, Fritz Henglein, and Cosmin E. Oancea. 2017. Futhark: purely functional GPU-programming with nested parallelism and in-place array updates. In *Proceedings of the 38th ACM SIGPLAN Conference on Programming Language Design and Implementation*. ACM, Barcelona Spain, 556–571. doi:10.1145/3062341.3062354
 - [19] Troels Henriksen, Niels G. W. Serup, Martin Elsman, Fritz Henglein, and Cosmin E. Oancea. 2017. Futhark: purely functional GPU-programming with nested parallelism and in-place array updates. *SIGPLAN Not.* 52, 6 (June 2017), 556–571. doi:10.1145/3140587.3062354
 - [20] Bastian Köpcke, Sergei Gorlatch, and Michel Steuwer. 2024. Descend: A Safe GPU Systems Programming Language. 8 (2024), 841–864. Issue PLDI. doi:10.1145/3656411
 - [21] NVIDIA Corporation. 2025. *CUB: CUDA Unbound*. <https://docs.nvidia.com/cuda/cub/index.html> CUDA Toolkit Documentation.
 - [22] NVIDIA Corporation. 2025. *CUDA C++ Programming Guide*. NVIDIA. <https://docs.nvidia.com/cuda/cuda-c-programming-guide/> Accessed: 2025-11-06.
 - [23] NVIDIA Corporation. 2025. Parallel Thread Execution (PTX) ISA — Asynchronous Warpgroup Level Matrix Multiply-Accumulate Instructions.

A Complete Prism Type System, Semantics, and Syntactic Soundness Proofs

A.1 Basic Definitions

Hierarchy Levels $h ::= \mathbf{Grid} \mid \mathbf{Block} \mid \mathbf{Thread}$

Memory Kinds $l ::= \mathbf{Local} \mid \mathbf{Shared} \mid \mathbf{Global}$

Perspectives $\pi : h \times \mathbb{N}$

Base Types $b ::= \mathbf{bool} \mid \mathbf{int} \mid \mathbf{float}$

Types $\tau ::= b \mid b[]^l \mid \mathbf{Fun}(\Gamma, \pi, m) \mid \mathbf{async} \ \tau$

Contexts $\Gamma ::= \cdot \mid \Gamma, x : \pi \ \tau$

Shared Memory Remaining $m : \mathbb{N}$

Perspectives π are part of an algebra parameterized over some constant values T (the number of threads per block) and B (the number of blocks per grid). With these values, we have two isomorphisms:

$$(\mathbf{Block}, 1) \cong (\mathbf{Thread}, T)$$

and

$$(\mathbf{Grid}, 1) \cong (\mathbf{Block}, B)$$

The **group** and **split** operations of Parallax allows us to move along this isomorphism, from left to right. For clarity, in Prism we split these operations into three: a **split** operation that can split perspectives into multiple narrower ones, and a **destruct** operation that directly moves us along the isomorphism, and a **group** that divides our current perspective into equal sized parts.

Perspectives π are also lexicographically ordered in the obvious way. hs are ordered such that $\mathbf{Thread} \leq \mathbf{Block} \leq \mathbf{Grid}$, and $(h_1, n_1) \leq (h_2, n_2)$ iff $n_1 \mid n_2$ and $h_1 \leq h_2$.

We define scalar multiplication $i \times h$ of natural numbers with hs : $i \times (h, n) = (h, in)$.

We also define division of perspectives and hierarchy levels of type $\pi \times \pi \rightarrow \mathbb{N}$. $\mathbf{Grid}/\mathbf{Block} = B$ and $\mathbf{Block}/\mathbf{Thread} = T$.

We lift this to perspectives like so: $(h_1, n_1) / (h_2, n_2) = ((h_1/h_2) \cdot n_1) / n_2$.

Lastly we define a partial \downarrow operator on hs such that $\downarrow \mathbf{Block} = \mathbf{Thread}$ and $\downarrow \mathbf{Grid} = \mathbf{Block}$. Note that $\downarrow \mathbf{Thread}$ is undefined. This operator lifts to π s whose second component is 1 and encodes the leftward component of the isomorphism above: $\downarrow (\mathbf{Block}, 1) = (\mathbf{Thread}, T)$ and $\downarrow (\mathbf{Grid}, 1) = (\mathbf{Block}, B)$.

Note also that the presentation of these rules in the main body of the paper elide the m portion, which tracks the maximum amount of memory a given computation is allowed to use. In the full system presented here, both the typing rules and the operational semantics carry an additional piece of information tracking allocated memory.

A.2 Complete Prism Typing Rules

A.2.1 Expressions.

$$\frac{x : \pi \ \tau \in \Gamma}{\Gamma \vdash^\pi x : \tau} \text{ T-Var} \quad \frac{}{\Gamma \vdash^\pi n : \mathbf{int}} \text{ T-Int} \quad \frac{}{\Gamma \vdash^\pi f : \mathbf{float}} \text{ T-Float}$$

$$\frac{}{\Gamma \vdash^\pi b : \mathbf{bool}} \text{ T-Bool} \quad \frac{\pi < (\mathbf{Grid}, 1)}{\Gamma \vdash^\pi \mathbf{partition_id} : \mathbf{int}} \text{ T-Partition-Id}$$

$$\frac{\Gamma \vdash^{\pi'} e_1 : \tau[]^l \quad \Gamma \vdash^\pi e_2 : \mathbf{int} \quad l = \mathbf{Global} \text{ or } l = \mathbf{Local} \quad \pi \leq \pi'}{\Gamma \vdash^\pi e_1[e_2] : \tau} \text{ T-Arr-Access}$$

$$\frac{\Gamma \vdash^{\pi'} e_1 : \tau[]^{\mathbf{Shared}} \quad \Gamma \vdash^\pi e_2 : \mathbf{int} \quad \pi \leq (\mathbf{Block}, 1) \quad \pi \leq \pi'}{\Gamma \vdash^\pi e_1[e_2] : \tau} \text{ T-Arr-Access-Shared}$$

$$\frac{\Gamma \vdash^\pi e_1 : \mathbf{int} \quad \Gamma \vdash^\pi e_2 : \mathbf{int}}{\Gamma \vdash^\pi e_1 \text{ bop } e_2 : \mathbf{int}} \text{ T-Bop} \quad \frac{\Gamma \vdash^\pi e_1 : \mathbf{int} \quad \Gamma \vdash^\pi e_2 : \mathbf{int}}{\Gamma \vdash^\pi e_1 \text{ cmp } e_2 : \mathbf{bool}} \text{ T-Cmp}$$

A.2.2 Statements.

$$\begin{array}{c}
\frac{f : ^\pi \mathbf{Fun}(x_i : \tau_i, \pi, m') \in \Gamma \quad \Gamma \vdash^\pi e_i : \tau_i \quad m' \leq m}{\Gamma \vdash_m^\pi f(e_1, \dots, e_n)} \text{ T-Function-Call} \\
\\
\frac{\Gamma \vdash_m^{(h, n_1)} s_1 \quad \Gamma \vdash_m^{(h, n_2)} s_2 \quad n_1, n_2 \text{ align to } n}{\Gamma \vdash_m^{(h, n)} \mathbf{split}(n_1, n_2) \{s_1\} \{s_2\}} \text{ T-Split} \quad \frac{\Gamma \vdash_m^{\downarrow \pi} s}{\Gamma \vdash_m^\pi \mathbf{destruct in } s} \text{ T-Destruct} \\
\\
\frac{\Gamma \vdash_m^{(h, n)} s}{\Gamma \vdash_m^{(h, q \cdot n)} \mathbf{group } q \text{ } s} \text{ T-Group} \quad \frac{x : ^{\pi'} \tau \in \Gamma \quad \Gamma \vdash^\pi e : \tau \quad \pi' \leq \pi}{\Gamma \vdash_m^\pi x = e} \text{ T-Assn} \\
\\
\frac{}{\Gamma \vdash_m^\pi \mathbf{init}_\psi} \text{ T-Sync-Init} \quad \frac{}{\Gamma \vdash_m^\pi \mathbf{dec}_\psi} \text{ T-Sync-Dec} \\
\\
\frac{}{\Gamma \vdash_m^\pi \mathbf{wait}_\psi} \text{ T-Sync-Wait} \quad \frac{}{\Gamma \vdash_m^\pi \mathbf{skip}} \text{ T-Skip} \quad \frac{n \leq m}{\Gamma \vdash_m^\pi \mathbf{free } n} \text{ T-Free} \\
\\
\frac{\Gamma \vdash^\pi e : \tau \quad \Gamma, x : ^{\pi'} \tau \vdash_m^\pi s \quad \pi' \leq \pi \quad \tau \text{ not an array type}}{\Gamma \vdash_m^\pi x : \tau @ \pi' = e \text{ in } s} \text{ T-Decl} \\
\\
\frac{\Gamma \vdash^{\pi'} e_1 : \tau []^l \quad \Gamma \vdash^\pi e_2 : \mathbf{int} \quad \Gamma \vdash^\pi e_3 : \tau \quad l = \mathbf{Global} \text{ or } l = \mathbf{Local} \quad \pi' \leq \pi}{\Gamma \vdash_m^\pi e_1 [e_2] = e_3} \text{ T-Arr-Assn} \\
\\
\frac{\Gamma \vdash^{\pi'} e_1 : \tau []^{\mathbf{Shared}} \quad \Gamma \vdash^\pi e_2 : \mathbf{int} \quad \Gamma \vdash^\pi e_3 : \tau \quad \pi \leq (\mathbf{Block}, 1) \quad \pi' \leq \pi}{\Gamma \vdash_m^\pi e_1 [e_2] = e_3} \text{ T-Arr-Assn-Shared} \\
\\
\frac{\Gamma \vdash^\pi e : \mathbf{bool} \quad \Gamma \vdash_{m_1}^\pi s_1 \quad \Gamma \vdash_{m_2}^\pi s_2}{\Gamma \vdash_{\max(m_1, m_2)}^\pi \mathbf{if } e \text{ then } s_1 \text{ else } s_2} \text{ T-If} \\
\\
\frac{\Gamma \vdash^\pi e : \mathbf{bool} \quad \Gamma \vdash_m^\pi s}{\Gamma \vdash_m^\pi \mathbf{while } e \text{ do } s} \text{ T-While} \quad \frac{\Gamma \vdash_{m_1}^\pi s_1 \quad \Gamma \vdash_{m_2}^\pi s_2}{\Gamma \vdash_{\max(m_1, m_2)}^\pi s_1 ; s_2} \text{ T-Seq} \\
\\
\frac{\Gamma, x : ^\pi \tau []^l \vdash_m^\pi s \quad l = \mathbf{Global} \text{ or } l = \mathbf{Local}}{\Gamma \vdash_{m+n \cdot \mathbf{size}(\tau)}^\pi x := \mathbf{alloc } l \text{ } \tau \text{ } n \text{ in } s} \text{ T-Alloc} \\
\\
\frac{\Gamma, x : ^{(\mathbf{Block}, 1)} \tau []^{\mathbf{Shared}} \vdash_m^\pi s}{\Gamma \vdash_{m+n \cdot \mathbf{size}(\tau)}^{(\mathbf{Block}, 1)} x := \mathbf{alloc Shared } \tau \text{ } n \text{ in } s} \text{ T-Alloc-Shared} \\
\\
\frac{\Gamma, y : ^{(h, n/c)} \tau []^l \vdash_m^{(h, n)} s \quad c | n \quad l \neq \mathbf{Local}}{\Gamma, x : ^{(h, n)} \tau []^l \vdash_m^{(h, n)} \mathbf{partition}_\psi x \text{ into } y \text{ by } c \text{ in } s} \text{ T-Partition} \\
\\
\frac{\Gamma, y : ^{(h, n')} \tau []^l \vdash_m^{(h, n')} s \quad n' \leq n \quad l \neq \mathbf{Local}}{\Gamma, x : ^{(h, n)} \tau []^l \vdash_m^{(h, n)} \mathbf{claim}_\psi x \text{ into } y \text{ at } n' \text{ in } s} \text{ T-Claim} \\
\\
\frac{\Gamma, y : ^{\downarrow \pi} \tau []^l \vdash_m^\pi s \quad l \neq \mathbf{Local}}{\Gamma, x : ^\pi \tau []^l \vdash_m^\pi \mathbf{lower}_\psi x \text{ into } y \text{ in } s} \text{ T-Lower} \\
\\
\frac{\Gamma, y : ^{(\mathbf{Thread}, 1)} \mathbf{async } \tau []^l \vdash_m^{(\mathbf{Thread}, 1)} s}{\Gamma, x : ^{(\mathbf{Thread}, 1)} \tau []^l \vdash_m^{(\mathbf{Thread}, 1)} \mathbf{async_partition}_\phi x \text{ into } y \text{ in } s} \text{ T-Async-Partition}
\end{array}$$

$$\begin{array}{c}
\frac{}{\Gamma, x : (\text{Thread}, 1) \text{ async } \tau[]^l, y : (\text{Thread}, 1) \tau[]'^l \vdash_m^{(\text{Thread}, 1)} \text{ async_memcpy}(x, y)} \text{ T-Async-Memcpy} \\
\frac{}{\Gamma, x : \pi \tau[]^l, y : \pi \tau[]'^l \vdash_m^\pi \text{ memcpy}(x, y)} \text{ T-Memcpy}
\end{array}$$

A.3 Complete Prism Semantics

A.3.1 Definitions.

Global Memory $\Sigma ::= \cdot \mid \Sigma, n \mapsto^\pi v$

Shared Memory $\sigma ::= \cdot \mid \sigma, n \mapsto^\pi v$

Local Memory $\eta ::= \cdot \mid \eta, n \mapsto^\pi v$

Block Memory Map $S ::= \forall n \in B, n \mapsto \sigma$

Thread Memory Map $L ::= \forall n \in T, n \mapsto \eta$

Synchronization Map $\Psi : \psi \rightarrow \mathbb{N} \rightarrow \mathbb{N}$

Deferred Computations Map $\Phi : \phi \rightarrow \{s\}$

In real GPUs, thread IDs are only unique within their block. However, in this calculus for simplicity we assume thread IDs are global. One can convert back and forth between this abstracted notion of a thread ID and a block-unique ID via addition modulo T . That is, $t_{\text{real}} = t_{\text{simplified}} \bmod T$ and $t_{\text{simplified}} = t_{\text{real}} + b \cdot T$.

By convention the names for local and shared and global memory do not conflict, as on the GPU they will be separate pointer spaces. Additionally, we freely interchange between using names for variables and integer locations.

In the main body of the paper, for simplicity we elide the synchronization map and the deferred computation map from the operational semantics, as our theorems do not make any guarantees about non-interference. However, as they are part of the full semantics, we include them here for completeness. By convention the synchronization and deferred computation maps are a total functions, initialized to map to $\lambda_{\cdot}.0$ for ψ s and $\lambda_{\cdot}.\{\}$ for ϕ s not explicitly initialize.

The shape of the judgment for a single thread is $\eta, \sigma, \Sigma, t, b, p, \Psi \vdash_m^\pi s \leadsto s' \dashv_{m'} \eta', \sigma', \Sigma', \Psi'$. The t here is the thread ID, the b is the block ID, and the p is the perspective ID. The last of these three is modified and managed by the rules for **split**, **group** and **destruct**, and tracks the relative position of the thread within a group. This semantics is in a small step style.

The shape of the judgment for expressions is $\eta, \sigma, \Sigma \vdash_{\pi}^\pi e \Downarrow v$. The two π s represent the ambient compute context (i.e., the context in which resources are being read), while π' , represents the target compute context (i.e., the compute context of the variable into which the result of the expression is going to be written. This is relevant for computing the value of **partition_id**, which divides the two contexts. As a shorthand, we can divide a perspective by a scalar value like so: $(h, n)/c = (h, n)/(h, c)$.

The overall evaluation of a program is expressed as

$$L, S, \Sigma, P, \Psi, \Phi \leadsto L', S', \Sigma', P', \Psi', \Phi'.$$

In this judgment P serves as a *thread pool*, mapping pairs of thread and block IDs (which don't change) to statements and memory (which can be updated by stepping). One can think of P as tracking which program is running on each thread. This steps according to the following rule:

$$\frac{L(t), S(b), \Sigma, t, b, 0, \Psi, \Phi \vdash_m^{(\text{Grid}, 1)} s \rightsquigarrow s' \dashv_{m'} \eta', \sigma', \Sigma', \Psi', \Phi' \quad P(t, b) = (s, m)}{L, S, \Sigma, P, \Psi, \Phi \rightsquigarrow L[t \mapsto \eta'], S[b \mapsto \sigma'], \Sigma', P[(t, b) \mapsto (s', m')], \Psi', \Phi'} \text{ S-Program}$$

For simplicity of notation, we define an **update** operation that searches the three environments for the one that contains the variable being used (by convention, there is no conflict between the environments, as in reality they exist in three separate address spaces). We also define a similar **get** operation that retrieves a variable from memory, and a **rename** operation that remaps a variable with the same value but under a different name.

$$\mathbf{update}(\eta, \sigma, \Sigma, x, v) = (\eta[x \mapsto^\pi v], \sigma, \Sigma) \text{ when } x \in^\pi \eta$$

$$\mathbf{update}(\eta, \sigma, \Sigma, x, v) = (\eta, \sigma[x \mapsto^\pi v], \Sigma) \text{ when } x \in^\pi \sigma$$

$$\mathbf{update}(\eta, \sigma, \Sigma, x, v) = (\eta, \sigma, \Sigma[x \mapsto^\pi v]) \text{ when } x \in^\pi \Sigma$$

$$\mathbf{get}(\eta, \sigma, \Sigma, x) = \eta(x) \text{ when } x \in^\pi \eta$$

$$\mathbf{get}(\eta, \sigma, \Sigma, x) = \sigma(x) \text{ when } x \in^\pi \sigma$$

$$\mathbf{get}(\eta, \sigma, \Sigma, x) = \Sigma(x) \text{ when } x \in^\pi \Sigma$$

$$\mathbf{rename}(\eta, \sigma, \Sigma, x, y, \pi') = (\eta[y \mapsto^{\pi'} \eta(x)], \sigma, \Sigma) \text{ when } x \in^\pi \eta$$

$$\mathbf{rename}(\eta, \sigma, \Sigma, x, y, \pi') = (\eta, \sigma[y \mapsto^{\pi'} \sigma(x)], \Sigma) \text{ when } x \in^\pi \sigma$$

$$\mathbf{rename}(\eta, \sigma, \Sigma, x, y, \pi') = (\eta, \sigma, \Sigma[y \mapsto^{\pi'} \Sigma(x)]) \text{ when } x \in^\pi \Sigma$$

A.3.2 Perspective Management Rules.

$$\frac{p < n_1 \quad n_1, n_2 \text{ align to } n \quad \eta, \sigma, \Sigma, t, b, p, \Psi, \Phi \vdash_m^{(h, n_1)} s_1 \rightsquigarrow s'_1 \dashv_{m'} \eta', \sigma', \Sigma', \Psi', \Phi'}{\eta, \sigma, \Sigma, t, b, p, \Psi, \Phi \vdash_m^{(h, n)} \mathbf{split}(n_1, n_2) \{s_1\} \{s_2\} \rightsquigarrow \mathbf{split}(n_1, n_2) \{s'_1\} \{s'_2\} \dashv_{m'} \eta', \sigma', \Sigma', \Psi', \Phi'} \text{ S-Split-Left}$$

$$\frac{p < n_1 \quad n_1, n_2 \text{ align to } n}{\eta, \sigma, \Sigma, t, b, p, \Psi, \Phi \vdash_m^{(h, n)} \mathbf{split}(n_1, n_2) \{\mathbf{skip}\} \{s_2\} \rightsquigarrow \mathbf{skip} \dashv_m \eta, \sigma, \Sigma, \Psi, \Phi} \text{ S-Split-Left-Done}$$

$$\frac{p \geq n_1 \quad p < n_1 + n_2 \quad n_1, n_2 \text{ align to } n \quad \eta, \sigma, \Sigma, t, b, p - n_1, \Psi, \Phi \vdash_m^{(h, n_2)} s_2 \rightsquigarrow s'_2 \dashv_{m'} \eta', \sigma', \Sigma', \Psi', \Phi'}{\eta, \sigma, \Sigma, t, b, p, \Psi, \Phi \vdash_m^{(h, n)} \mathbf{split}(n_1, n_2) \{s_1\} \{s_2\} \rightsquigarrow \mathbf{split}(n_1, n_2) \{s_1\} \{s'_2\} \dashv_{m'} \eta', \sigma', \Sigma', \Psi', \Phi'} \text{ S-Split-Right}$$

$$\frac{p \geq n_1 \quad p < n_1 + n_2 \quad n_1, n_2 \text{ align to } n}{\eta, \sigma, \Sigma, t, b, p, \Psi, \Phi \vdash_m^{(h, n)} \mathbf{split}(n_1, n_2) \{s_1\} \{\mathbf{skip}\} \rightsquigarrow \mathbf{skip} \dashv_m \eta, \sigma, \Sigma, \Psi, \Phi} \text{ S-Split-Right-Done}$$

$$\frac{p \geq n_1 + n_2 \quad n_1, n_2 \text{ align to } n}{\eta, \sigma, \Sigma, t, b, p, \Psi, \Phi \vdash_m^{(h, n)} \mathbf{split}(n_1, n_2) \{s_1\} \{s_2\} \rightsquigarrow \mathbf{skip} \dashv_m \eta, \sigma, \Sigma, \Psi, \Phi} \text{ S-Split-None}$$

$$\frac{\eta, \sigma, \Sigma, t, b, t \text{ mod } T, \Psi, \Phi \vdash_m^{(\text{Thread}, T)} s \rightsquigarrow s' \dashv_{m'} \eta', \sigma', \Sigma', \Psi', \Phi'}{\eta, \sigma, \Sigma, t, b, 0, \Psi, \Phi \vdash_m^{(\text{Block}, 1)} \mathbf{destruct in } s \rightsquigarrow \mathbf{destruct in } s' \dashv_{m'} \eta', \sigma', \Sigma', \Psi', \Phi'} \text{ S-Destruct-Block}$$

$$\frac{\eta, \sigma, \Sigma, t, b, b \text{ mod } B, \Psi, \Phi \vdash_m^{(\text{Block}, B)} s \leadsto s' \dashv_{m'} \eta', \sigma', \Sigma', \Psi', \Phi'}{\eta, \sigma, \Sigma, t, b, 0, \Psi, \Phi \vdash_m^{(\text{Grid}, 1)} \text{destruct in } s \leadsto \text{destruct in } s' \dashv_{m'} \eta', \sigma', \Sigma', \Psi', \Phi'} \quad \text{S-Destruct-Grid}$$

$$\frac{}{\eta, \sigma, \Sigma, t, b, 0, \Psi, \Phi \vdash_m^\pi \text{destruct in skip} \leadsto \text{skip} \dashv_m \eta, \sigma, \Sigma, \Psi, \Phi} \quad \text{S-Destruct-Done}$$

$$\frac{\eta, \sigma, \Sigma, t, b, p \text{ mod } n, \Psi, \Phi \vdash_m^{(h, n)} s \leadsto s' \dashv_{m'} \eta', \sigma', \Sigma', \Psi', \Phi'}{\eta, \sigma, \Sigma, t, b, p, \Psi, \Phi \vdash_m^{(h, q \cdot n)} \text{group } q \leadsto \text{group } q \text{ } s'; \dashv_{m'} \eta', \sigma', \Sigma', \Psi', \Phi'} \quad \text{S-Group}$$

$$\frac{}{\eta, \sigma, \Sigma, t, b, p, \Psi, \Phi \vdash_m^\pi \text{group } q \text{ skip} \leadsto \text{skip}; \dashv_m \eta, \sigma, \Sigma, \Psi, \Phi} \quad \text{S-Group-Done}$$

A.3.3 Thread Synchronization. We define a **size** operation on perspectives to compute the size of a perspective (the number of individual compute resources sharing it). The operation is defined thusly:

$$\begin{aligned} \text{size}(\text{Thread}, n) &= n \\ \text{size}(\text{Block}, n) &= n \cdot T \\ \text{size}(\text{Grid}, n) &= n \cdot B \cdot T \end{aligned}$$

$$\frac{\Psi(\psi)(p) = 0}{\eta, \sigma, \Sigma, t, b, p, \Psi, \Phi \vdash_m^\pi \text{wait}_\psi \leadsto \text{skip} \dashv_m \eta, \sigma, \Sigma, \Psi, \Phi} \quad \text{S-Sync-Wait-Done}$$

$$\frac{\Psi(\psi)(p) \neq 0}{\eta, \sigma, \Sigma, t, b, p, \Psi, \Phi \vdash_m^\pi \text{wait}_\psi \leadsto \text{wait}_\psi \dashv_m \eta, \sigma, \Sigma, \Psi, \Phi} \quad \text{S-Sync-Wait-Spin}$$

$$\frac{\Psi' = \Psi(\psi)[p \mapsto \Psi(\psi)(p) - 1]}{\eta, \sigma, \Sigma, t, b, p, \Psi, \Phi \vdash_m^\pi \text{dec}_\psi \leadsto \text{skip} \dashv_m \eta, \sigma, \Sigma, \Psi', \Phi} \quad \text{S-Sync-Dec}$$

$$\frac{\Psi(\psi)(p) = 0 \quad \Psi' = \Psi(\psi)[p \mapsto \text{size}(\pi)]}{\eta, \sigma, \Sigma, t, b, p, \Psi, \Phi \vdash_m^\pi \text{init}_\psi \leadsto \text{skip} \dashv_m \eta, \sigma, \Sigma, \Psi', \Phi} \quad \text{S-Sync-Init-Zero}$$

$$\frac{\Psi(\psi)(p) \neq 0}{\eta, \sigma, \Sigma, t, b, p, \Psi, \Phi \vdash_m^\pi \text{init}_\psi \leadsto \text{skip} \dashv_m \eta, \sigma, \Sigma, \Psi, \Phi} \quad \text{S-Sync-Init-Nonzero}$$

A.3.4 Asynchrony.

$$\frac{\text{rename}(\eta, \sigma, \Sigma, x, y, (\text{Thread}, 1), t, b, p, \Psi, \Phi \vdash_m^{(\text{Thread}, 1)} s \leadsto s' \dashv_{m'} \eta', \sigma', \Sigma', \Psi', \Phi')}{\eta, \sigma, \Sigma, t, b, p, \Psi, \Phi \vdash_{m'}^{(\text{Thread}, 1)} \text{async_partition}_\phi x \text{ into } y \text{ in } s \leadsto \text{async_partition}_\phi x \text{ into } y \text{ in } s' \dashv_{m'} \eta', \sigma', \Sigma', \Psi', \Phi'} \quad \text{S-Async-Partition-Congr}$$

$$\frac{\Phi = \Phi'[\phi \mapsto \Phi'(\phi) \cup \{s\}]}{\eta, \sigma, \Sigma, t, b, p, \Psi, \Phi \vdash_m^{(\text{Thread}, 1)} \text{async_partition}_\phi x \text{ into } y \text{ in skip} \leadsto (\text{async_partition}_\phi x \text{ into } y \text{ in } s) \dashv_m \eta, \sigma, \Sigma, \Psi, \Phi'} \quad \text{S-Async-Partition-Unwind}$$

$$\frac{\Phi(\phi) = \emptyset}{\eta, \sigma, \Sigma, t, b, p, \Psi, \Phi \vdash_m^{(\text{Thread}, 1)} \text{async_partition}_{\phi} x \text{ into } y \text{ in skip} \leadsto \text{skip} \dashv_m \eta, \sigma, \Sigma, \Psi, \Phi} \text{ S-Async-Partition-Done}$$

$$\frac{}{\eta, \sigma, \Sigma, t, b, p, \Psi, \Phi \vdash_m^{(\text{Thread}, 1)} \text{async_memcpy}(x, y) \leadsto \text{skip} \dashv_m \eta, \sigma, \Sigma, \Psi, \Phi[\phi \mapsto \Phi(\phi) \cup \{\text{memcpy}(x, y)\}]} \text{ S-Async-Memcpy}$$

$$\frac{(\eta', \sigma', \Sigma') = \text{update}(\eta, \sigma, \Sigma, x, \text{get}(\eta, \sigma, \Sigma, y))}{\eta, \sigma, \Sigma, t, b, p, \Psi, \Phi \vdash_m^{\pi} \text{memcpy}(x, y) \leadsto \text{skip} \dashv_m \eta', \sigma', \Sigma', \Psi, \Phi} \text{ S-Memcpy}$$

A.3.5 Variables and Memory.

$$\frac{\eta, \sigma, \Sigma \vdash_{\pi'} e \Downarrow v \quad \pi' \leq \pi}{\eta, \sigma, \Sigma, t, b, p, \Psi, \Phi \vdash_m^{\pi} x : \tau @ \pi' := e \text{ in } s \leadsto s \dashv_m \eta[x \mapsto^{\pi'} v], \sigma, \Sigma, \Psi, \Phi} \text{ S-Decl}$$

$$\frac{}{\eta, \sigma, \Sigma, t, b, p, \Psi, \Phi \vdash_m^{\pi} \text{free } n \dashv_{m-n} \eta, \sigma, \Sigma, \Psi, \Phi} \text{ S-Free}$$

$$\frac{\eta, \sigma, \Sigma, t, b, p, \Psi, \Phi \vdash_m^{\pi} x := \text{alloc Local } \tau \ n \text{ in } s \leadsto s; \text{free } (n \cdot \text{size}(\tau)) \dashv_{m+n \cdot \text{size}(\tau)} \eta[x \mapsto^{\pi} \langle x, n \rangle], \sigma, \Sigma, \Psi, \Phi}{\eta, \sigma, \Sigma, t, b, p, \Psi, \Phi \vdash_m^{\pi} x := \text{alloc Local } \tau \ n \text{ in } s \leadsto s; \text{free } (n \cdot \text{size}(\tau)) \dashv_{m+n \cdot \text{size}(\tau)} \eta[x \mapsto^{\pi} \langle x, n \rangle], \sigma, \Sigma, \Psi, \Phi} \text{ S-Alloc-Local}$$

$$\frac{\pi = (\text{Block}, 1)}{\eta, \sigma, \Sigma, t, b, p, \Psi, \Phi \vdash_m^{\pi} x := \text{alloc Shared } \tau \ n \text{ in } s \leadsto s; \text{free } (n \cdot \text{size}(\tau)) \dashv_{m+n \cdot \text{size}(\tau)} \eta, \sigma[x \mapsto^{\pi} \langle x, n \rangle], \Sigma, \Psi, \Phi} \text{ S-Alloc-Shared}$$

$$\frac{}{\eta, \sigma, \Sigma, t, b, p, \Psi, \Phi \vdash_m^{\pi} x := \text{alloc Global } \tau \ n \text{ in } s \leadsto s; \text{free } (n \cdot \text{size}(\tau)) \dashv_{m+n \cdot \text{size}(\tau)} \eta, \sigma, \Sigma[x \mapsto^{\pi} \langle x, n \rangle], \Psi, \Phi} \text{ S-Alloc-Global}$$

$$\frac{x \in \pi' \quad \eta, \sigma, \Sigma \quad \eta, \sigma, \Sigma \vdash_{\pi'} e \Downarrow v \quad (\eta', \sigma', \Sigma') = \text{update}(\eta, \sigma, \Sigma, x, v) \quad \pi' \leq \pi}{\eta, \sigma, \Sigma, t, b, p, \Psi, \Phi \vdash_m^{\pi} x = e \leadsto \text{skip} \dashv_m \eta', \sigma', \Sigma', \Psi, \Phi} \text{ S-Assn}$$

$$\frac{\eta, \sigma, \Sigma, t, b, p \vdash_{\pi} e_1 \Downarrow \langle l, n \rangle \quad i < n \quad \pi' \leq \pi \quad \eta, \sigma, \Sigma, t, b, p \vdash_{\pi} e_2 \Downarrow i \quad x \in \pi' \quad \eta, \sigma, \Sigma \quad \eta, \sigma, \Sigma \vdash_{\pi} e_3 \Downarrow v \quad (\eta', \sigma', \Sigma') = \text{update}(\eta, \sigma, \Sigma, l + i, v)}{\eta, \sigma, \Sigma, t, b, p, \Psi, \Phi \vdash_m^{\pi} e_1[e_2] = e_3 \leadsto \text{skip} \dashv_m \eta', \sigma', \Sigma', \Psi, \Phi} \text{ S-Arr-Assn}$$

$$\frac{s' = s[(y + c \cdot p) / y] \quad (\eta', \sigma', \Sigma') = \text{rename}(\eta, \sigma, \Sigma, x, y, \pi / c)}{\eta, \sigma, \Sigma, t, b, p, \Psi, \Phi \vdash_m^{\pi} \text{partition} \psi x y c s \leadsto \text{init}_{\psi}; s'; \text{dec}_{\psi}; \text{wait}_{\psi} \dashv_m \eta', \sigma', \Sigma', \Psi, \Phi} \text{ S-Partition}$$

$$\frac{(\eta', \sigma', \Sigma') = \text{rename}(\eta, \sigma, \Sigma, x, y, (h, n_1))}{\eta, \sigma, \Sigma, t, b, p, \Psi, \Phi \vdash_m^{(h, n_1 + n_2)} \text{claim}_{\psi} x \text{ into } y \text{ at } n_1 \text{ in } s \leadsto \text{init}_{\psi}; \text{split}(n_1, n_2) \{s'\} \{\text{skip}\}; \text{dec}_{\psi}; \text{wait}_{\psi} \dashv_m \eta', \sigma', \Sigma', \Psi, \Phi} \text{ S-Claim}$$

$$\frac{(\eta', \sigma', \Sigma') = \text{rename}(\eta, \sigma, \Sigma, x, y, \downarrow \pi)}{\eta, \sigma, \Sigma, t, b, p, \Psi, \Phi \vdash_m^{\pi} \text{lower}_{\psi} x \text{ into } y \text{ in } s \leadsto \text{init}_{\psi}; s; \text{dec}_{\psi}; \text{wait}_{\psi} \dashv_m \eta', \sigma', \Sigma', \Psi, \Phi} \text{ S-Lower}$$

A.3.6 Control Flow.

$$\begin{array}{c}
\frac{\eta, \sigma, \Sigma \vdash_{\pi}^{\pi} e \Downarrow \text{true}}{\eta, \sigma, \Sigma, t, b, p, \Psi, \Phi \vdash_m^{\pi} \text{if } e \text{ then } s_1 \text{ else } s_2 \rightsquigarrow s_1 \dashv_m \eta, \sigma, \Sigma, \Psi, \Phi} \text{S-If-True} \\
\frac{\eta, \sigma, \Sigma \vdash_{\pi}^{\pi} e \Downarrow \text{false}}{\eta, \sigma, \Sigma, t, b, p, \Psi, \Phi \vdash_m^{\pi} \text{if } e \text{ then } s_1 \text{ else } s_2 \rightsquigarrow s_2 \dashv_m \eta, \sigma, \Sigma, \Psi, \Phi} \text{S-If-False} \\
\frac{}{\eta, \sigma, \Sigma, t, b, p, \Psi, \Phi \vdash_m^{\pi} \text{while } e \text{ do } s \rightsquigarrow \text{if } e \text{ then } (s; \text{while } e \text{ do } s) \text{ else skip } \dashv_m \eta, \sigma, \Sigma, \Psi, \Phi} \text{S-While} \\
\frac{\eta, \sigma, \Sigma, t, b, p, \Psi, \Phi \vdash_m^{\pi} s_1 \rightsquigarrow s'_1 \dashv_m^{\pi'} \eta', \sigma', \Sigma', \Psi', \Phi'}{\eta, \sigma, \Sigma, t, b, p, \Psi, \Phi \vdash_m^{\pi} s_1; s_2 \rightsquigarrow s'_1; s_2 \dashv_m^{\pi'} \eta', \sigma', \Sigma', \Psi', \Phi'} \text{S-Seq-First} \\
\frac{}{\eta, \sigma, \Sigma, t, b, p, \Psi, \Phi \vdash_m^{\pi} \text{skip}; s_2 \rightsquigarrow s_2 \dashv_m \eta, \sigma, \Sigma, \Psi, \Phi} \text{S-Seq-Done} \\
\frac{\Sigma(f) = \{[x_1 : \tau_1, \dots, x_n : \tau_n], s\} \quad \sigma, \Sigma \vdash_{\pi}^{\pi} e_i \Downarrow v_i \quad m' \leq m}{\eta, \sigma, \Sigma, t, b, p, \Psi, \Phi \vdash_m^{\pi} f(e_1, \dots, e_n) \rightsquigarrow s \dashv_m \eta[x_i \mapsto v_i], \sigma, \Sigma, \Psi, \Phi} \text{S-Function-Call}
\end{array}$$

A.3.7 Expressions.

$$\begin{array}{c}
\frac{\pi < (\text{Grid}, 1) \quad \pi' \leq \pi}{\eta, \sigma, \Sigma \vdash_{\pi'}^{\pi} \text{partition_id} \Downarrow \pi / \pi' - 1} \text{E-Partition-Id} \\
\frac{}{\eta, \sigma, \Sigma \vdash_{\pi'}^{\pi} x \Downarrow \text{get}(\eta, \sigma, \Sigma, x)} \text{E-Var} \\
\frac{\eta, \sigma, \Sigma \vdash_{\pi'}^{\pi} e_1 \Downarrow \langle l, n \rangle \quad \eta, \sigma, \Sigma \vdash_{\pi'}^{\pi} e_2 \Downarrow i \quad i < n \quad \pi' \leq \pi}{\eta, \sigma, \Sigma \vdash_{\pi'}^{\pi} e_1[e_2] \Downarrow \text{get}(\eta, \sigma, \Sigma, l+i)} \text{E-Arr-Access} \\
\frac{}{\eta, \sigma, \Sigma \vdash_{\pi'}^{\pi} \vdash n \Downarrow n} \text{E-Int} \quad \frac{}{\eta, \sigma, \Sigma \vdash_{\pi'}^{\pi} \vdash b \Downarrow b} \text{E-Bool} \\
\frac{\eta, \sigma, \Sigma \vdash_{\pi'}^{\pi} e_1 \Downarrow v_1 \quad \eta, \sigma, \Sigma \vdash_{\pi'}^{\pi} e_2 \Downarrow v_2 \quad v = v_1 \text{ bop } v_2}{\eta, \sigma, \Sigma \vdash_{\pi'}^{\pi} e_1 \text{ bop } e_2 \Downarrow v} \text{E-Bop} \\
\frac{\eta, \sigma, \Sigma \vdash_{\pi'}^{\pi} e_1 \Downarrow v_1 \quad \eta, \sigma, \Sigma \vdash_{\pi'}^{\pi} e_2 \Downarrow v_2 \quad v = v_1 \text{ cmp } v_2}{\eta, \sigma, \Sigma \vdash_{\pi'}^{\pi} e_1 \text{ cmp } e_2 \Downarrow v} \text{E-Cmp}
\end{array}$$

A.4 Theorems and Proofs

Note that in this section we assume no out of bounds array accesses. In general Parallax (and by extension Prism) makes no guarantees about array out of bounds.

A.4.1 More definitions. As a premise to our type safety theorems, we need to assume we have a well-typed environment, written $\Gamma \vdash \eta, \sigma, \Sigma$. We define what this means inductively

$$\begin{array}{c}
\frac{}{\eta, \sigma, \Sigma \vdash n : \text{int}} \text{V-Int} \quad \frac{}{\eta, \sigma, \Sigma \vdash b : \text{bool}} \text{V-Bool} \quad \frac{}{\eta, \sigma, \Sigma \vdash f : \text{float}} \text{V-Float} \\
\frac{\forall i < n, \eta, \sigma, \Sigma \vdash \text{get}(\eta, \sigma, \Sigma, x+i) : \tau}{\eta, \sigma, \Sigma \vdash \langle x, n \rangle : \tau[]^l} \text{V-Array} \quad \frac{\Gamma, x_i :^{\pi} \tau_i \vdash_m^{\pi} s \quad \Gamma \vdash \cdot, \text{fns } \Sigma}{\eta, \sigma, \Sigma \vdash \{x_i : \tau_i, s\} : \text{Fun}(x_i : \tau_i, \pi, m)} \text{V-Function} \\
\frac{}{\cdot \vdash \eta, \sigma, \Sigma} \text{G-Empty} \quad \frac{\eta(x) =^{\pi} v \quad \eta, \sigma, \Sigma \vdash v : \text{int}}{\Gamma, x :^{\pi} \text{int} \vdash \eta, \sigma, \Sigma} \text{G-Int}
\end{array}$$

$$\begin{array}{c}
\frac{\eta(x) =^\pi v \quad \eta, \sigma, \Sigma \vdash v : \mathbf{bool}}{\Gamma, x :^\pi \mathbf{bool} \vdash \eta, \sigma, \Sigma} \quad \mathbf{G-Bool} \qquad \frac{\eta(x) =^\pi v \quad \eta, \sigma, \Sigma \vdash v : \mathbf{float}}{\Gamma, x :^\pi \mathbf{float} \vdash \eta, \sigma, \Sigma} \quad \mathbf{G-Float} \\
\\
\frac{\eta(x) =^\pi v \quad \eta, \sigma, \Sigma \vdash v : \tau[]}{\Gamma, x :^\pi \tau[]^{\mathbf{Local}} \vdash \eta, \sigma, \Sigma} \quad \mathbf{G-Local} \qquad \frac{\sigma(x) =^\pi v \quad \eta, \sigma, \Sigma \vdash v : \tau[]}{\Gamma, x :^\pi \tau[]^{\mathbf{Shared}} \vdash \eta, \sigma, \Sigma} \quad \mathbf{G-Shared} \\
\\
\frac{\Sigma(x) =^\pi v \quad \eta, \sigma, \Sigma \vdash v : \tau[]}{\Gamma, x :^\pi \tau[]^{\mathbf{Global}} \vdash \eta, \sigma, \Sigma} \quad \mathbf{G-Global} \qquad \frac{\Sigma(x) =^\pi v \quad \eta, \sigma, \Sigma \vdash v : \mathbf{Fun}(\Gamma', \pi, m)}{\Gamma, x :^\pi \mathbf{Fun}(\Gamma', \pi, m) \vdash \eta, \sigma, \Sigma} \quad \mathbf{G-Function}
\end{array}$$

We can prove a couple simple lemmas about well-typed environments under operations like **rename**, **update**, and **get**.

LEMMA A.1. (*Well-typed get*) If $\Gamma \vdash \eta, \sigma, \Sigma$ and $x :^\pi \tau \in \Gamma$ then $\eta, \sigma, \Sigma \vdash \mathbf{get}(\eta, \sigma, \Sigma, x) : \tau$.

LEMMA A.2. (*Well-typed rename*) If $\Gamma, x :^\pi \tau \vdash \eta, \sigma, \Sigma$ then $\Gamma, x :^\pi \tau, y :^{\pi'} \tau' \vdash \mathbf{rename}(\eta, \sigma, \Sigma, x, y, \pi')$.

LEMMA A.3. (*Well-typed update*) If $\Gamma \vdash \eta, \sigma, \Sigma$ and $\eta, \sigma, \Sigma \vdash v : \tau$ then $\Gamma, x :^\pi \tau \vdash \mathbf{update}(\eta, \sigma, \Sigma, x, v)$.

We also define a well-formedness precondition on p with respect to π :

$$(h, n) \vdash p ::= p < n$$

We also define well-formedness for the async stack:

$$\Gamma \vdash \Phi ::= \forall \phi, s \in \Phi(\phi), \Gamma \vdash_m^{(\mathbf{Thread}, 1)} s$$

A.5 Proofs

LEMMA A.4. (*Expression Safety*) If $\Gamma \vdash^\pi e : \tau$ and $\Gamma \vdash \eta, \sigma, \Sigma$ and $\pi \vdash p$ and $\pi' \leq \pi$, then there is some v such that $\eta, \sigma, \Sigma \vdash_{\pi'} e \Downarrow v$ and $v : \tau$.

PROOF. This proof proceeds by induction on the typing relation for expressions. Despite the fact that this property implies termination, we do not need a logical relation to prove it because the expression language is very simple.

The **T-Int**, **T-Float**, and **T-Bool** cases are trivial, using the rules **E-Int**, **E-Bool** and **E-Float** to compute values. In the case for **T-Partition-Id**, the π premises of the typing rules and our assumption that $\pi' \leq \pi$ match the premises of the evaluation rule, so this rule is simple as well. The cases for **T-Bop** and **T-Cmp** follow directly from the inductive hypotheses, assuming a valid and correctly implemented set of binary operators and comparators.

The only interesting cases are **T-Arr-Access** and **T-Arr-Access-Shared**.

In both cases our inductive hypotheses and inversion give us that $\pi' \leq \pi$, and e_1 evaluates to a $\langle x, n \rangle$, and that all the values between x and $x+n$ in the appropriate environment are typed at τ . We also know that e_2 evaluates to an integer i . We assume that all array accesses are in bounds, so $i < l$, which is sufficient to use the **E-Arr-Access** rule to complete this case, and the proof. \square

LEMMA A.5. (*Expression Determinism*) If $\eta, \sigma, \Sigma, t, b, p \vdash_{\pi'}^\pi e \Downarrow v_1$ and $\eta, \sigma, \Sigma, t, b, p \vdash_{\pi'}^\pi e \Downarrow v_2$ then $v_1 = v_2$.

PROOF. Straightforward by induction on the semantic derivation. \square

LEMMA A.6. (*Expression Well-Typedness*) If $\Gamma \vdash_{\pi'}^\pi e : \tau$ and $\Gamma \vdash \eta, \sigma, \Sigma$ and $\pi \vdash p$, and $\eta, \sigma, \Sigma, t, b, p \vdash_{\pi'}^\pi e \Downarrow v$, then $v : \tau$.

PROOF. By our expression safety lemma our well-typed expression must evaluate to a well-typed value v' . By our determinism lemma v' must be the same as v , so v is well-typed. \square

LEMMA A.7. (*Statement Progress*) If $\Gamma \vdash \eta, \sigma, \Sigma$ and $\Gamma \vdash_m^\pi s$ and $\pi \vdash p$, then either s is **skip** or there is some s' such that $\eta, \sigma, \Sigma, t, b, p, \Psi, \Phi \vdash_m^\pi s \rightsquigarrow s' \dashv_{m''} \eta', \sigma', \Sigma', \Psi', \Phi'$

PROOF. This proceeds by induction on the typing derivation.

Perspective Management Rules.

- Case **T-Split**:

In this case we have by our assumption that $\pi \vdash p$ that $p < n$. We also have that n_1, n_2 **align to** n , so $n_1 + n_2 \leq n$. There are three cases to consider, then: when $p < n_1$, when $p \geq n_1$ and $p < n_1 + n_2$, and when $p \geq n_1 + n_2$.

In the first case, we have by our inductive hypothesis that s is either **skip** or that it can step in an (h, n_1) context. In the former case we can use the **S-Split-Left-Done** rule and in the latter we can use the **S-Split-Left rule**.

The second case is almost symmetric. The only additional work we have to do is to argue that $(h, n_2) \vdash p - n_1$, or equivalently that $p - n_1 < n_2$. This, however, is immediate from our assumption that $p < n_1 + n_2$.

In the last case, we just use the **S-Split-None** rule to step to **skip**.

- Case **T-Destruct**

By our inductive hypothesis, we know that s can step at $\downarrow \pi$ if $\downarrow \pi \vdash p$. The \downarrow operation is only defined at **(Block, 1)** or **(Grid, 1)**, so we only need to consider the cases where π is one of those.

In the former case p becomes $t \bmod T$ while $\downarrow \pi$ is **(Thread, T)**. $t \bmod T < T$ for any t so this satisfies the requirement that $\pi \vdash p$, which lets us use our inductive hypothesis: s is either **skip** or can step. If it can step, we can use this to satisfy the premise of **S-Destruct-Block** to step in this case. If it is **skip**, then we use the rule **S-Destruct-Done** to step instead.

The latter case is the same, except using the fact that $b \bmod B < B$ and the **S-Destruct-Grid** rule.

- Case **T-Group**

In this case we have by assumption that $(h, q \cdot n) \vdash p$, i.e., that $p < q \cdot n$.

In this case we have our IH that if $(h, n) \vdash p'$ for some p' , then s is either **skip** or steps with (h, n) perspective with p' as our perspective ID.

We choose p' to be $p \bmod n$. This is always $< n$, so $(h, n) \vdash p'$. This lets us use our IH to get that s is either **skip** (in which case we can use the **S-Group-Done** rule to step) or itself steps, which lets us use the **S-Group** rule to step.

Thread Synchronization Rules.

- Case **T-Sync-Wait**

$\Psi(\psi)(p)$ is either zero or it is not. In the former case we use the **S-Sync-Wait-Done** rule and in the latter we use **S-Sync-Wait-Spin**.

- Case **T-Sync-Dec**

We use the **S-Sync-Dec** rule to step.

- Case **T-Sync-Init**

We use the **S-Sync-Init-Zero** or **S-Sync-Init-Nonzero** rules depending on whether $\Psi(\psi)(p)$ is zero or not.

Asynchrony Rules.

- Case **T-Async-Partition**

In this case, we have via our IH that if $\Gamma, y : (\text{Thread}, 1) \text{ async } \tau []^l \vdash \eta', \sigma', \Sigma'$, then we can either step s with **(Thread, 1)** perspective under environments η', σ' , and Σ' , or s is **skip**.

We have by assumption that $\Gamma, x : (\text{Thread}, 1) \text{ async } \tau[]^l \vdash \eta, \sigma, \Sigma$. By our environment renaming lemma, this gives us what we need to use our IH with $(\eta', \sigma', \Sigma')$ as $\text{rename}(\eta, \sigma, \Sigma, x, y, (\text{Thread}, 1))$.

Thus s either steps or is **skip**. In the former case we can step with **S-Async-Partition-Congr**, and in the latter we can use either **S-Async-Partition-Unwind** or **S-Async-Partition-Done** depending on whether $\Phi(\phi)$ is empty or not.

- Case **T-Async-Memcpy**
Immediate via use of the **S-Async-Memcpy** rule.
- Case **T-Memcpy**
Immediate via use of the **S-Memcpy** rule.

Memory Rules.

- Case **T-Decl**
Via our lemma about expression type safety and our hypothesis that e is well-typed, we obtain the premises necessary to use the **S-Decl** rule to step.
- Case **T-Arr-Assn**
Each of e_1 , e_2 and e_3 must evaluate to a well-typed value by the expression type safety lemma. In particular, both e_1 evaluates to some $\langle l, n \rangle$ and e_2 evaluates to some i . We assume all array accesses are in bounds, so this is sufficient to use the **S-Arr-Assn** rule to step.
- Case **T-Arr-Assn-Shared**
Same as previous case.
- Case **T-Free**
Trivial via the **S-Free** rule.
- Case **T-Partition**
Trivial via the **S-Partition** rule.
- Case **T-Claim**
Trivial via the **S-Claim** rule.
- Case **T-Lower**
Trivial via the **S-Lower** rule.
- Case **T-Alloc**
We assume that l is not **Shared**, so we can use the **S-Alloc-Local** or **S-Alloc-Global** rule, depending on whether l is **Local** or **Global**.
- Case **T-Alloc-Shared**
Trivial via the **S-Alloc-Shared** rule.

Control Rules.

- Case **T-Skip**
Trivial
- Case **T-While**
Trivial, all while loops step via the **S-While** rule
- Case **T-If**
By our proof of expression type safety, the expression e steps to either the boolean value true or false. We can thus use either the **S-If-True** or **S-If-False** rules to step.
- Case **T-Seq**
In this case we know by our IH that s_1 is either **skip** or can step. In the former case we use the **S-Seq-Done** rule and in the latter we use the **S-Seq-First** rule.
- Case **T-Function-Call**
In this case we know by our expression safety lemma that each of the arguments will evaluate to a well-typed value. We also have by assumption that f has a function type,

which by inversion on the **V-Function** rule tells us that it is a closure type. Additionally our assumption that $\Gamma \vdash \eta, \sigma, \Sigma$ tell us that Σ contains f at the same type that Γ does. These premises are sufficient to use the **S-Function-Call** rule. \square

LEMMA A.8. (*Statement Preservation*) If $\Gamma \vdash \eta, \sigma, \Sigma$ and $\Gamma \vdash_m^\pi s$ and $\eta, \sigma, \Sigma, t, b, p, \Psi, \Phi \vdash_{m'}^\pi s \rightsquigarrow s' \dashv_{m''} \eta', \sigma', \Sigma, \Psi', \Phi'$ and $\pi \vdash p$ and $m \geq m'$ and $\Gamma \vdash \Phi$, then there is some Γ' such that $\Gamma \subseteq \Gamma'$ and $\Gamma' \vdash_m^\pi s'$ and $\Gamma' \vdash \eta', \sigma', \Sigma'$ and $m \geq m''$ and $\Gamma' \vdash \Phi'$.

PROOF. We proceed by induction on the derivation of $\Gamma \vdash_m^\pi s$.

Perspective Management Rules.

- Case **T-Split**

In this case we have by assumption that $(h, n) \vdash p$, $\Gamma \vdash \eta, \sigma, \Sigma$, and $\eta, \sigma, \Sigma, t, b, p, \Psi, \Phi \vdash_{m'}^\pi \mathbf{split}(n_1, n_2) \{s_1\} \{s_2\} \rightsquigarrow s' \dashv_{m''} \eta', \sigma', \Sigma, \Psi', \Phi'$. Our inductive hypotheses give us that if for any p , if $(h, n_1) \vdash p$ and s_1 steps with perspective ID p , or if $(h, n_2) \vdash p$ and s_2 steps with perspective ID p then their results are well typed.

By inversion on our semantic derivation, we are in one of 5 cases.

In the **S-Split-Left** case $p < n_1$ and s_1 steps to s'_1 . This is sufficient to tell us that s'_1 is well-typed and the output environments of that relation η', σ' , and Σ' are all well-typed by $\Gamma' \supseteq \Gamma$, and that the memory is properly bounded by the typing rules.

We can thus use the **T-Split-Left** rule to conclude that the result of this case is well-typed.

The **S-Split-Left-Done** rule is trivial via the **T-Skip** rule.

The **Right** cases are symmetric, with the observation that when $p \geq n_1$ and $p < n_1 + n_2$ then $p - n_1 < n_2$.

The last **S-Split-None** rule is trivial via the **T-Skip** rule.

- Case **T-Destruct**

In this case we have by assumption that $\downarrow \pi$ is defined, so π is either (**Block**, 1) or (**Grid**, 1). We also assume that $\eta, \sigma, \Sigma, t, b, p, \Psi, \Phi \vdash_{m'}^\pi \mathbf{destruct\ in} s \rightsquigarrow s' \dashv_{m''} \eta', \sigma', \Sigma, \Psi', \Phi'$. We also have by our inductive hypothesis that for any p such that $\downarrow \pi \vdash p$ and s'' such that s steps to s'' at p , then that step preserved well-typedness.

By inversion on the step relation, we are in one of three cases.

If the rule used **S-Destruct-Block**, then we know that s steps to s'' and p is $t \bmod T$. (**Thread**, T) $\vdash t \bmod T$ for any t , so we can use our inductive hypothesis to conclude that the s'' stepped to by s is well-typed, as are its environments and memory usage. The **T-Destruct** rule then gives us our desired goal.

The **S-Destruct-Grid** case proceeds similarly, while the **S-Destruct-Done** case is trivial.

- Case **T-Group**

In this case we have by assumption that

$\eta, \sigma, \Sigma, t, b, p, \Psi, \Phi \vdash_{m'}^{(h, q \cdot n)} \mathbf{group} q s \rightsquigarrow s' \dashv_{m''} \eta', \sigma', \Sigma, \Psi', \Phi'$. We have by our inductive hypothesis that s steps to s'' at some perspective ID p' and $(h, n) \vdash p'$, then s'' is well typed, as are the other outputs of that step.

By inversion on the step relation, we are in one of two cases. The **S-Group-Done** case is trivial, so we shall focus on the **S-Group** case. In this case we have that s steps to s'' at perspective ID $p \bmod n$. It is always the case that $(h, n) \vdash p \bmod n$ for any p , so we can use our inductive hypothesis to conclude that s'' is well-typed, as are its output environments and memory usage. From there, it is a simple application of the **T-Group** rule to conclude that $\mathbf{group} q s'$ is well-typed, and to finish the case.

Thread Synchronization Rules. These rules are all trivial: with one exception all thread synchronization primitives step to skip without changing environment or memory, and are thus obviously well-typed. **S-Sync-Wait-Spin** does not produce skip, but it steps to the same statement as we already assumed typechecks in the premise of the lemma, so is straightforward nonetheless.

If we wanted to say something about deadlock freedom we'd have more work here, but we aren't doing that, so these rules are easy.

Asynchrony Rules.

- **Case T-Async-Partition** In this case we have that s is well-typed in a context where x has been renamed into y , with $(\mathbf{Thread}, 1)$ perspective. We also have that $\Gamma, x :^{(\mathbf{Thread}, 1)} \tau[]^l \vdash \eta, \sigma, \Sigma$ and $\Gamma, x :^{(\mathbf{Thread}, 1)} \tau[]^l \vdash \Phi$.
By inversion, we are in one of three cases.
In the **S-Async-Partition-Done** case, we are done.
In the **S-Async-Partition-Congr** case, our inductive hypothesis gives us that there is some Γ' such that $\Gamma, y :^{(\mathbf{Thread}, 1)} \mathbf{async} \tau[]^l \subseteq \Gamma'$ and $\Gamma' \vdash \Phi'$ and $\Gamma' \vdash \mathbf{rename}(\eta, \sigma, \Sigma, x, y, (\mathbf{Thread}, 1))$ via our well-typed renaming lemma. This lets us use the **T-Async-Partition** rule to check this case, with a choice of Γ' as $\Gamma', x :^{(\mathbf{Thread}, 1)} \tau[]^l$.
In the **S-Async-Partition-Unwind** case, our assumption that Φ is well-typed tells us that $\Gamma, y :^{(\mathbf{Thread}, 1)} \vdash (\mathbf{Thread}, 1)_m s$. Thus, we can use the **T-Async-Partition** rule to type this case.
- **Case T-Async-Memcpy**
In this case the statement and environment typing are trivial, we need only to show that the async stack remains well typed.
In this case we have that x and y have the same type at $(\mathbf{Thread}, 1)$. This is sufficient for us to check $x = y$ at $(\mathbf{Thread}, 1)$, meaning that adding that instruction to the stack maintains its well-typedness.
- **Case T-Memcpy**
Immediate via use of the **S-Memcpy** rule. We just need to show that the environment remains well typed, which we know via our lemmas about **update** and **get**.

Memory Rules.

- **Case T-Decl**
By inversion, we are using **S-Decl** rule for evaluation. Our well-typed expression lemma gives us that v is well-typed, so it follows from our assumptions and our lemmas about extending environments that the extended η and s are well-typed by $\Gamma, x :^\pi \tau$.
- **Case T-Free**
Trivial.
- **Case T-Alloc**
By inversion we are either in the **S-Alloc-Local** or **S-Alloc-Global** rules. In either case, we assume that $\Gamma, x :^\pi \tau[]^l$ checks s , meaning we can use our extended environment lemmas and the **T-Seq** and **T-Free** rules to check these cases.
- **Case T-Alloc-Shared**
Same as previous case.
- **Case T-Partition**
In this case we have by assumption that $\Gamma, y :^{(h, n/c)} \tau[]^l s$ and l is not local and c divides n . By inversion on our step relation we must be in the **S-Partition** case, so we have $\eta, \sigma, \Sigma, t, b, p, \Psi, \Phi \vdash_m^{\pi} \mathbf{partition} \psi x y c s \rightsquigarrow \mathbf{init}_{\psi}; s'; \mathbf{dec}_{\psi}; \mathbf{wait}_{\psi} \dashv_m \eta', \sigma', \Sigma', \Psi, \Phi$ where

$s' = s[(y + c \cdot p)/y]$ and $(\eta', \sigma', \Sigma') = \mathbf{rename}(\eta, \sigma, \Sigma, x, y, \pi/c)$. Via our well-typed renaming lemma contexts we know that we can check the renamed environments in the extended environment $\Gamma, y : \pi/c \tau[]^l, x : \pi \tau[]^l$, and this context is also sufficient to check s' (via a substitution-preserves-typing lemma that is obvious). Using the **T-Skip** rule this is exactly what we need to show to complete this case, as the thread sync primitives check trivially via their typing rules.

- **Case T-Claim**
Essentially the same as **T-Partition**.
- **Case T-Lower**
Essentially the same as **T-Partition**.

Control Rules.

- **Case T-If**
We have by assumption that **if** e **then** s_1 **else** s_2 steps to some s' , and by inversion we know that either e evaluates to **true** and s' is s_1 , or e evaluates to **false** and s' is s_2 . In either case, our inductive hypotheses is sufficient to tell us that these are well-typed. In particular, in both cases our IHs tell us that the amount of memory used by stepping each branch of the **if** is less than the amount of memory computed by the type system for each branch. Because the whole **if** expression checks using the greater of the memory usage of m_1 or m_2 (i.e., the memory usage on each branch), the resulting usage for the whole conditional is also bounded by the type system.
- **Case T-Skip**
Trivial: **skip** does not step
- **Case T-Seq**
In this case we have that s_1 and s_2 are both well-typed (with m_1 memory and m_2 memory respectively), and $m = \mathbf{max}(m_1, m_2)$. We also have by inversion that s_1 either steps to **skip** or s'_1 , and our inductive hypothesis tells us that s'_1 is well-typed. In the former case we can use the **S-Seq-Done** rule to trivially finish the case. In the latter, our IH allows us to finish the case, since m_1 is always $\leq \mathbf{max}(m_1, m_2)$
- **Case T-While**
By inversion, we have that

$$\begin{aligned} & \eta, \sigma, \Sigma, t, b, p, \Psi, \Phi \vdash_m^\pi \mathbf{while} \ e \ \mathbf{do} \ s \\ & \leadsto \mathbf{if} \ e \ \mathbf{then} \ (s; \mathbf{while} \ e \ \mathbf{do} \ s) \ \mathbf{else} \ \mathbf{skip} \dashv_m \eta, \sigma, \Sigma, \Psi, \Phi. \end{aligned}$$

We also have by assumption that e and s are well-typed. With this information, through use of the **T-If**, **T-Seq**, **T-While**, and **T-Skip** rules, we can conclude that the result of this rule is also well-typed.

- **Case T-Function-Call**
By inversion we have that

$$\begin{aligned} & \eta, \sigma, \Sigma, t, b, p, \Psi, \Phi \vdash_m^\pi f(e_1, \dots, e_n) \\ & \leadsto \mathbf{call} \ s \ \mathbf{with} \ (x_i : \tau_i \mapsto v_i) @ \{\eta, \sigma, \Sigma, m'\} \dashv_m \eta, \sigma, \Sigma, \Psi, \Phi, \end{aligned}$$

and also that $\Sigma(f) = \{[x_1 : \tau_1, \dots, x_n : \tau_n], s\}$, and $\eta, \sigma, \Sigma, t, b, p \vdash^\pi e_i \Downarrow v_i$, and $m' \leq m$.

We also have from the premises of our case that $f : \pi \mathbf{Fun}(x_1 : \tau_1, \dots, x_n : \tau_n, \pi, m') \in \Gamma$, and $\Gamma \vdash^\pi e_i : \tau_i$, and $m' \leq m$.

Our lemma for expression well-typedness tells us that each v_i is a well-typed value, and our assumption that $\Gamma \vdash \eta, \sigma, \Sigma$ tells us that $\Sigma(f)$ is a well-typed function and thus that

fns $\Gamma, x_i :^{\pi_i} \tau_i \vdash_m^{\pi} s$. We can take the union of this with Γ to produce $\Gamma, x_i :^{\pi_i} \tau_i \vdash_m^{\pi} s$, which clearly checks s and the output environments. \square

1765
1766
1767
1768
1769
1770
1771
1772
1773
1774
1775
1776
1777
1778
1779
1780
1781
1782
1783
1784
1785
1786
1787
1788
1789
1790
1791
1792
1793
1794
1795
1796
1797
1798
1799
1800
1801
1802
1803
1804
1805
1806
1807
1808
1809
1810
1811
1812
1813