

FIREIRON: A SCHEDULING LANGUAGE FOR HIGH-PERFORMANCE LINEAR ALGEBRA ON GPUS

BASTIAN HAGEDORN | **SAM ELLIOTT** | **HENRIK BARTHELS** | **RAS BODIK** | **VINOD GROVER**
 University of Münster | lowRISC | RWTH Aachen University | Uni of Washington | NVIDIA

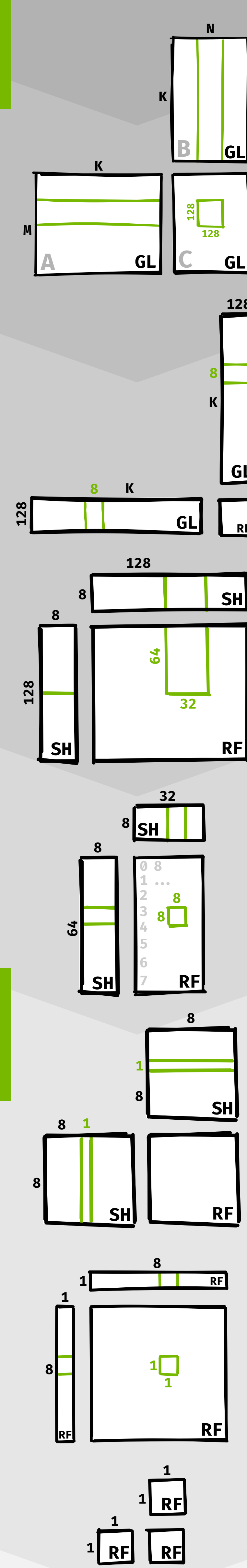
EXPLOITING HIERARCHICAL STRUCTURE IN HIGH-PERFORMANCE LINEAR ALGEBRA GPU IMPLEMENTATIONS

01

Most high-performance kernels for GPUs are written in a hierarchical style. The original problem to be computed is decomposable into smaller sub-problems of the same kind. These sub-problems are then assigned to and computed by the different levels of the compute hierarchy. The figure below visualizes this observation using a simple matrix multiplication kernel. Step by step, the matrix multiplication is decomposed into a hierarchy of tiles and data is transferred to lower levels of the memory hierarchy until eventually every thread computes a single FMA instruction. Here, FMA can be viewed as a matrix multiplication of matrices which only contain a single element.

```

1 __global__ void MatMul(const float A[M * K], const float B[K * N], float C[M * N]) {
2
3     shared float ASH[128][8], BSH[8][128];
4     float ARF[8][1], BRF[1][8], CRF[8][8];
5
6     iBlock ← 128 * blockIdx.x;
7     jBlock ← 128 * blockIdx.y;
8
9     CRF ← 0;
10
11     for (k ← 0; k < K / 8; k++) {
12
13         GLbToSh(A → ASH (8x128), start at (iBlock, jBlock))
14         GLbToSh(B → BSH (128x8), start at (iBlock, jBlock))
15         __syncthreads();
16
17         iWarp ← iBlock + warpIdx.x * 64;
18         jWarp ← jBlock + warpIdx.y * 32;
19
20         iThread ← iWarp + threadIdx.x * 8;
21         jThread ← jWarp + threadIdx.y * 8
22
23         for (kk ← 0; kk < 8; kk++)
24
25             ShToPvt(ASH → ARF (8x1), start at (iThread, jThread))
26             ShToPvt(BSH → BRF (1x8), start at (iThread, jThread))
27
28             for (i ← 0; i < 8; i++)
29                 for (j ← 0; j < 8; j++)
30
31                     CRF[i][j] += ARF[i][0] * BRF[0][j];
32
33             endfor
34         endfor
35
36     endfor
37
38     PvtToGlb(CRF → C (128x128), start at iBlock, jBlock)
39
40 } // end kernel
    
```



FIREIRON: A SCHEDULING LANGUAGE EXPRESS COMPLEX IMPLEMENTATIONS AS COMPOSITIONS OF SIMPLE PRIMITIVES

02

Fireiron introduces two main concepts: Specifications and Decompositions. A Specification (spec) is a data-structure describing the computation to implement. A spec contains enough information such that a programmer would be able to manually write an implementation. This especially entails that a spec keeps track of the shapes, locations and storage layouts of its input and output tensors, as well as which level of the compute hierarchy (i.e., Kernel, Block, Warp or Thread) is responsible for computing this operation. A Decomposition describes how to (partially) implement a given spec. More specifically, a decomposition is a function Spec → Spec which, given a spec, returns a new spec that represents the smaller decomposed sub-problem. Fireiron provides two main decompositions, tile and load, which allow implementations to use the compute and memory hierarchy of a GPU.



03 DECOMPOSING MATRIX MULTIPLICATION GRADUALLY DESCENDING THE COMPUTE AND MEMORY HIERARCHY

KERNEL-LEVEL In order to generate a high-performance Matrix Multiplication kernel with Fireiron, we start with a kernel-level MatMul specification...

Storage Layout

```

MatMul(ComputeHierarchy: Kernel,
      A: Matrix((M x K), FP16, GL, ColMajor),
      B: Matrix((K x N), FP16, GL, ColMajor),
      C: Matrix((M x N), FP16, GL, ColMajor))
.tile(128, 128).to(Block)
    
```

Location in Memory Hierarchy

Shape **Scalar Data Type**

BLOCK-LEVEL

MatMul-specific decompositions: In order to express advanced decompositions for storing computed results, and accumulating intermediate results in registers, we introduce a new decomposition .epilog(l,i,d). Similar to load, i and d are decompositions. Here, i describes the initialization of a new buffer in location l (usually in registers) used for accumulating the results. The .split(k) decomposition is used to create tiles in the K-dimension. This allows efficient use of shared memory.

```

.epilog(RF, init, store)
.split(8)
.load(A, SH, contiguousPrefetch) // = contiguousPrefetch = Move
.load(B, SH, crosswisePrefetch)
.tile(64, 32).to(Warp)
    
```

Move: Another Spec representing data movement between levels of the memory hierarchy. **.unroll:** A Refinement, more in a second...

WARP-LEVEL

```

.tile(8, 8).to(Thread)
.layout(ColMajor)
    
```

We also support **Tensor Cores!** In order to use them, computations must be decomposed to **warp-level WMMA** or **quad-pair mma** specifications

Refinements allow fine-grained control in decompositions. For example, **.unroll** is used to indicate loops shall be unrolled during code generation, and **.layout** is used to specify the mapping between tiles and the elements of the current compute hierarchy.

THREAD-LEVEL

```

.split(1)
.load(A, RF, aToRF)
.load(B, RF, bToRF)
.tile(1, 1)
.done
    
```

04 PERFORMANCE GENERATING FAST KERNELS COMPETITIVE TO cuBLAS

