

# Descend

## A Safe GPU Systems Programming Language

---

Bastian Köpcke, Sergei Gorlatch, Michel Steuwer

University of Münster, Technische Universität Berlin

June 27, 2024

**Computing on GPUs is great!**

**GPU Programming is too hard!**

```
1 __global__ void transpose(const float *input, float *output) {  
2     __shared__ float tmp[1024];  
3     for (int j = 0; j < 32; j += 8) {  
4         tmp[threadIdx.y+j*32+threadIdx.x] =  
5             input[(blockIdx.y*32+threadIdx.y+j)*2048 + blockIdx.x*32+threadIdx.x];  
6     }  
7     __syncthreads();  
8     for (int j = 0; j < 32; j += 8) {  
9         output[(blockIdx.x*32+threadIdx.y+j)*2048 + blockIdx.y*32+threadIdx.x] =  
10            tmp[threadIdx.x*32+threadIdx.y+j];  
11     }  
12 }
```

*Kernel functions* are executed by thousands of threads

```
1 __global__ void transpose(const float *input, float *output) {  
2     __shared__ float tmp[1024];  
3     for (int j = 0; j < 32; j += 8) {  
4         tmp[threadIdx.y+j*32+threadIdx.x] =  
5             input[(blockIdx.y*32+threadIdx.y+j)*2048 + (blockIdx.x*32+threadIdx.x)];  
6     }  
7     __syncthreads();  
8     for (int j = 0; j < 32; j += 8) {  
9         output[(blockIdx.x*32+threadIdx.y+j)*2048 + blockIdx.y*32+threadIdx.x] =  
10            tmp[threadIdx.x*32+threadIdx.y+j];  
11     }  
12 }
```

*Kernel functions* are executed by thousands of threads

Threads are hierarchically organized into multi-dimensional *blocks*  
Blocks are hierarchically are organized into a multi-dimensional *grid*

```
1 __global__ void transpose(const float *input, float *output) {  
2     __shared__ float tmp[1024];  
3     for (int j = 0; j < 32; j += 8) {  
4         tmp[threadIdx.y+j*32+threadIdx.x] ←  
5             input[(blockIdx.y*32+threadIdx.y+j)*2048 + blockIdx.x*32+threadIdx.x];  
6     }  
7     __syncthreads();  
8     for (int j = 0; j < 32; j += 8) {  
9         output[(blockIdx.x*32+threadIdx.y+j)*2048 + blockIdx.y*32+threadIdx.x] =  
10            tmp[threadIdx.x*32+threadIdx.y+j];  
11     }  
12 }
```

explicit indexing to determine where a thread accesses memory

*Kernel functions* are executed by thousands of threads

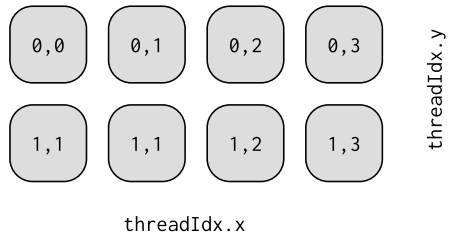
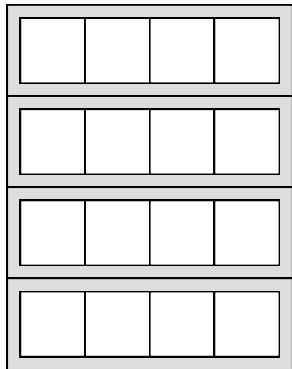
Threads are hierarchically organized into multi-dimensional *blocks*  
Blocks are hierarchically are organized into a multi-dimensional *grid*

```
1 __global__ void transpose(const float *input, float *output) {  
2     __shared__ float tmp[1024];  
3     for (int j = 0; j < 32; j += 8) {  
4         tmp[threadIdx.y+j*32+threadIdx.x]  
5             input[(blockIdx.y*32+threadIdx.y+j)*2048 + blockIdx.x*32+threadIdx.x];  
6     }  
7     __syncthreads();  
8     for (int j = 0; j < 32; j += 8) {  
9         output[(blockIdx.x*32+threadIdx.y+j)*2048 + blockIdx.y*32+threadIdx.x] =  
10            tmp[threadIdx.x*32+threadIdx.y+j];  
11     }  
12 }
```

explicit indexing to determine where a thread accesses memory

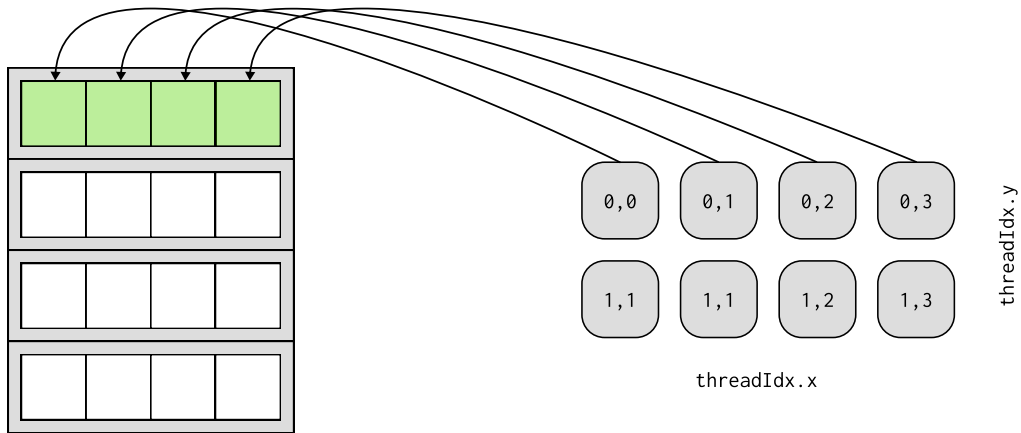
Is the implementation correct?

`tmp[threadIdx.y+j*4+threadIdx.x]`      `j=0`

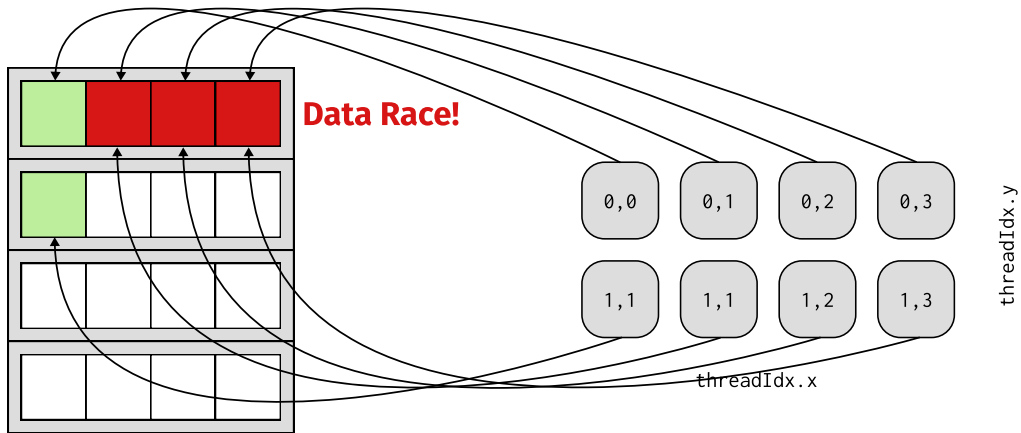




`tmp[threadIdx.y+j*4+threadIdx.x]`      `j=0`

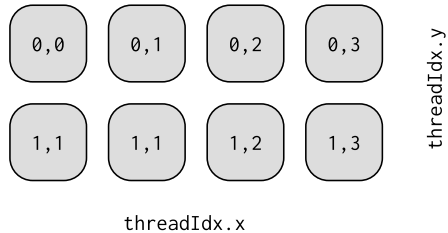
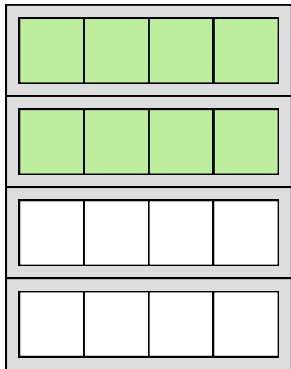


`tmp[threadIdx.y+j*4+threadIdx.x]`      `j=0`



```
tmp[threadIdx.y+j*4+threadIdx.x]    j=0
```

```
tmp[(threadIdx.y+j)*4+threadIdx.x]
```



**How can we create a  
safe GPU programming language?**

# Descend

## A Safe GPU Systems Programming Language

*"if the program successfully compiles,  
then it is free of data races and other memory problems"*

```
1 fn transpose(  
2   input: &gpu.global [[f32;2048];2048],  
3   output: &uniq gpu.global [[f32;2048];2048],  
4   [grid.blocks.forall(X).forall(Y)] tmp: &uniq gpu.shared [[f32;32];32]  
5 ) -[grid: gpu.grid<XY<64,64>,XY<32,8>>]-> () {  
6   sched(Y,X) block in grid {  
7     sched(Y,X) thread in block {  
8       for i in 0..4 {  
9         tmp.transpose.group_by_row::<32,4>[[thread]][i] =  
10          input.group_by_tile::<32,32>.transpose[[block]].group_by_row::<32,4>[[thread]][i]  
11       };  
12  
13       for i in 0..4 {  
14         output.group_by_tile::<32,32>[[block]].group_by_row::<32,4>[[thread]][i]=  
15         tmp.group_by_row::<32,4>[[thread]][i]  
16       }  
17     }  
18   }  
19 }
```

```
1 fn transpose(  
2   input: &gpu.global [[f32;2048];2048],  
3   output: &uniq gpu.global [[f32;2048];2048],  
4   [grid.blocks.forall(X).forall(Y)] tmp: &uniq gpu.shared [[f32;32];32]  
5 ) - [grid: gpu.grid<XY<64,64>,XY<32,8>>-> () {  
6   sched(Y,X) block in grid {  
7     sched(Y,X) thread in block {  
8       for i in 0..4 {  
9         tmp.transpose.group_by_row::<32,4>[[thread]][i] =  
10          input.group_by_tile::<32,32>.transpose[[block]].group_by_row::<32,4>[[thread]][i]  
11        };  
12  
13        for i in 0..4 {  
14          output.group_by_tile::<32,32>[[block]].group_by_row::<32,4>[[thread]][i]=  
15          tmp.group_by_row::<32,4>[[thread]][i]  
16        }  
17      }  
18    }  
19  }
```

function is explicitly annotated with a 2D **grid**  
of  $64 \times 64$  **blocks** with  $32 \times 8$  **threads** each

```
1 fn transpose(  
2   input: &gpu.global [[f32;2048];2048],  
3   output: &uniq gpu.global [[f32;2048];2048],  
4   [grid.blocks.forall(X).forall(Y)] tmp: &uniq gpu.shared [[f32;32];32]  
5 ) -[grid: gpu.grid<XY<64,64>,XY<32,8>>-> () {  
6   sched(Y,X) block in grid {  
7     sched(Y,X) thread in block {  
8       for i in 0..4 {  
9         tmp.transpose.group_by_row::<32,4>[[thread]][i] =  
10          input.group_by_tile::<32,32>.transpose[[block]].group_by_row::<32,4>[[thread]][i]  
11       };  
12  
13       for i in 0..4 {  
14         output.group_by_tile::<32,32>[[block]].group_by_row::<32,4>[[thread]][i]=  
15         tmp.group_by_row::<32,4>[[thread]][i]  
16       }  
17     }  
18   }  
19 }
```

collectively executed by the entire grid



```
1 fn transpose(  
2   input: &gpu.global [[f32;2048];2048],  
3   output: &uniq gpu.global [[f32;2048];2048],  
4   [grid.blocks.forall(X).forall(Y)] tmp: &uniq gpu.shared [[f32;32];32]  
5 ) -[grid: gpu.grid<XY<64,64>,XY<32,8>>-> () {  
6   sched(Y,X) block in grid {  
7     sched(Y,X) thread in block {  
8       for i in 0..4 {  
9         tmp.transpose.group_by_row::<32,4>[[thread]][i] =  
10          input.group_by_tile::<32,32>.transpose[[block]].group_by_row::<32,4>[[thread]][i]  
11       };  
12  
13       for i in 0..4 {  
14         output.group_by_tile::<32,32>[[block]].group_by_row::<32,4>[[thread]][i]=  
15         tmp.group_by_row::<32,4>[[thread]][i]  
16       }  
17     }  
18   }  
19 }
```

collectively executed by each block in the grid

```
1 fn transpose(  
2   input: &gpu.global [[f32;2048];2048],  
3   output: &uniq gpu.global [[f32;2048];2048],  
4   [grid.blocks.forall(X).forall(Y)] tmp: &uniq gpu.shared [[f32;32];32]  
5 ) -[grid: gpu.grid<XY<64,64>,XY<32,8>>-> () {  
6   sched(Y,X) block in grid {  
7     sched(Y,X) thread in block {  
8       for i in 0..4 {  
9         tmp.transpose.group_by_row::<32,4>[[thread]][i] =  
10        input.group_by_tile::<32,32>.transpose[[block]].group_by_row::<32,4>[[thread]][i]  
11      };  
12  
13      for i in 0..4 {  
14        output.group_by_tile::<32,32>[[block]].group_by_row::<32,4>[[thread]][i]=  
15        tmp.group_by_row::<32,4>[[thread]][i]  
16      }  
17    }  
18  }  
19 }
```

collectively executed by each thread in the block

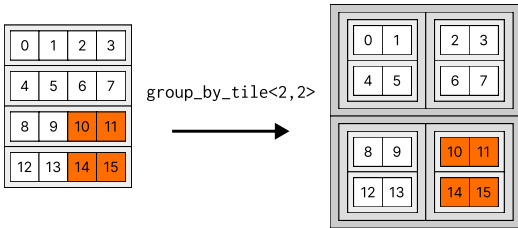
```
1 fn transpose(  
2   input: &gpu.global [[f32;2048];2048],  
3   output: &uniq gpu.global [[f32;2048];2048],  
4   [grid.blocks.forall(X).forall(Y)] tmp: &uniq gpu.shared [[f32;32];32]  
5 ) -[grid: gpu.grid<XY<64,64>,XY<32,8>>-> () {  
6   sched(Y,X) block in grid {  
7     sched(Y,X) thread in block {  
8       for i in 0..4 {  
9         tmp.transpose.group_by_row::<32,4>[[thread]][i] =  
10          input.group_by_tile::<32,32>.transpose[[block]].group_by_row::<32,4>[[thread]][i]  
11       };  
12  
13       for i in 0..4 {  
14         output.group_by_tile::<32,32>[[block]].group_by_row::<32,4>[[thread]][i] =  
15          tmp.group_by_row::<32,4>[[thread]][i]  
16       }  
17     }  
18   }  
19 }
```

safely write in parallel over blocks and threads using **views**

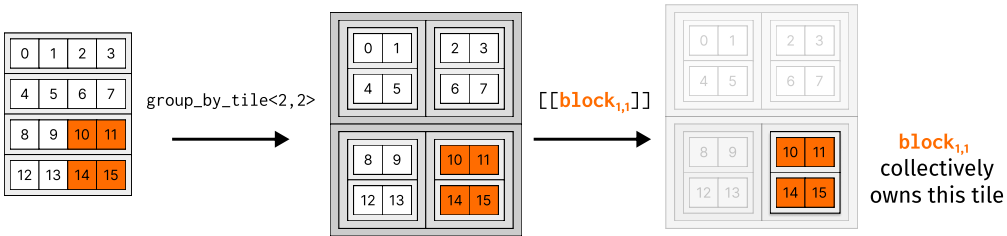
```
output.group_by_tile<2,2>[[block1,1]].group_by_row<2,2>[[threadq,1]]
```

0	1	2	3
4	5	6	7
8	9	10	11
12	13	14	15

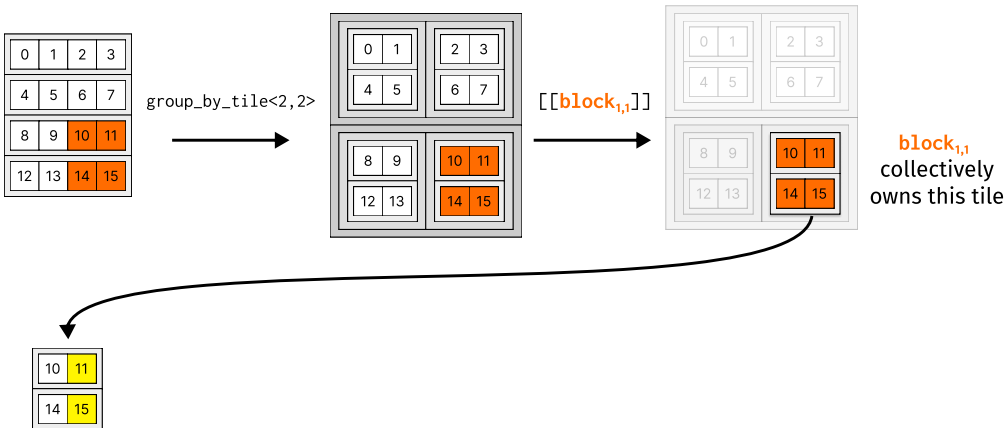
```
output.group_by_tile<2,2>[[block1,1]].group_by_row<2,2>[[thread0,1]]
```



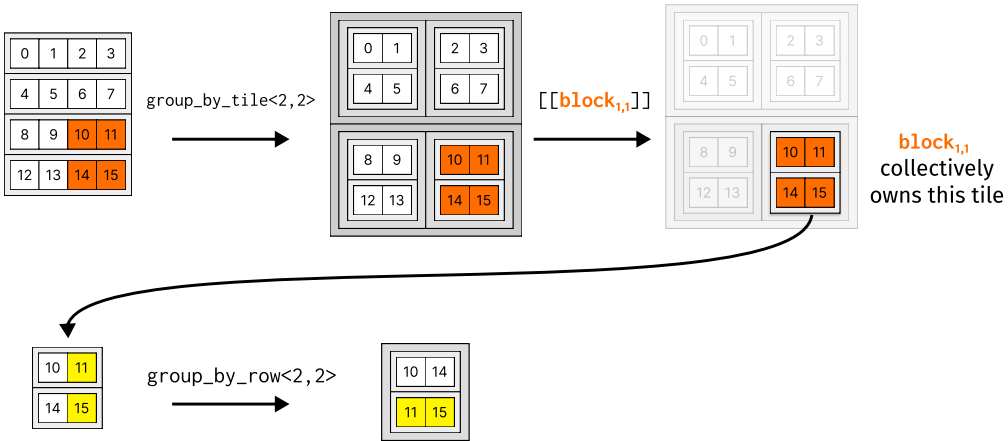
```
output.group_by_tile<2,2>[[block1,1]].group_by_row<2,2>[[thread0,1]]
```



```
output.group_by_tile<2,2>[[block1,1]].group_by_row<2,2>[[thread0,1]]
```

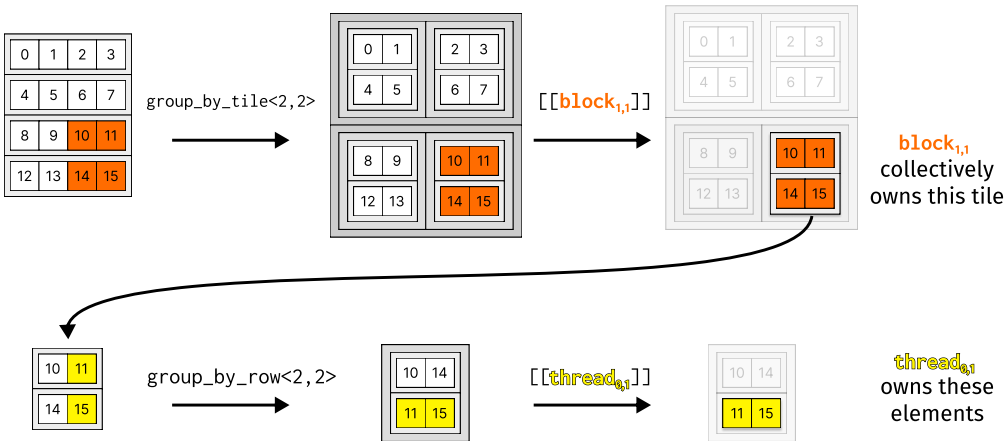


`output.group_by_tile<2,2>[[block1,1]].group_by_row<2,2>[[threadq,1]]`

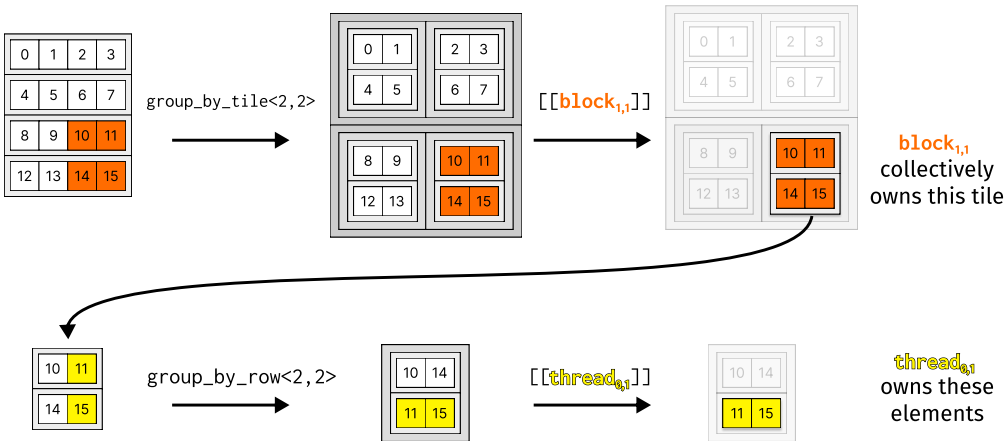




`output.group_by_tile<2,2>[[block1,1]].group_by_row<2,2>[[thread0,1]]`

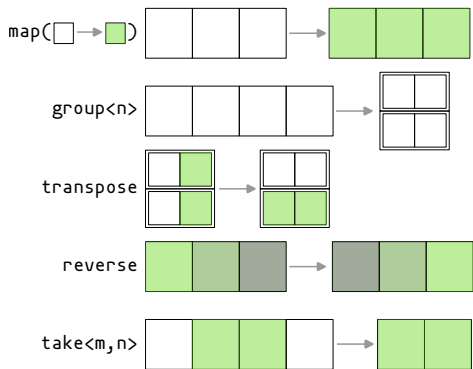


`output.group_by_tile<2,2>[[block1,1]].group_by_row<2,2>[[thread0,1]]`

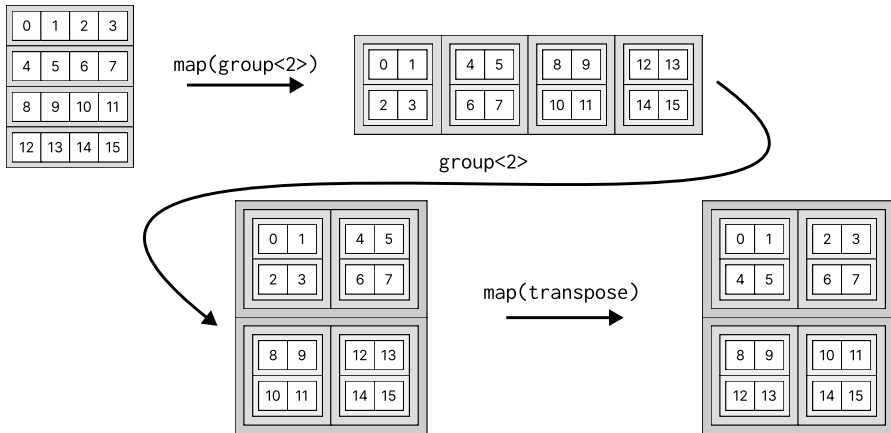


Descend's extended ownership model ensures that memory is accessed safely via views

# Descend's five basic view primitives



```
group_by_tile<2,2> = map(group<2>).group<2>.map(transpose)
```



**Views are composed to express complex memory access patterns**

```

1  fn transpose(
2    input: &gpu.global [[f32;2048];2048],
3    output: &uniq gpu.global [[f32;2048];2048],
4    [grid.blocks.forall(X).forall(Y)] tmp: &uniq gpu.shared [[f32;32];32]
5  ) -[grid: gpu.grid<XY<64,64>,XY(32,0)>] {
6    sched(Y,X) block in grid {
7      sched(Y,X) thread in block {
8        for i in 0..4 {
9          tmp.group_by_row::<32,4>[[thread]][i] =
10             input.group_by_tile::<32,32>.transpose[[block]].group_by_row::<32,4>[[thread]][i]
11        };
12
13        for i in 0..4 {
14          output.group_by_tile::<32,32>[[block]].group_by_row::<32,4>[[thread]][i] =
15             tmp.transpose.group_by_row::<32,4>[[thread]][i]
16        }
17      }
18    }
19  }

```

threads write safely into memory tmp

```
1 fn transpose(  
2   input: &gpu.global [[f32;2048];2048],  
3   output: &uniq gpu.global [[f32;2048];2048],  
4   [grid.blocks.forall(X).forall(Y)] tmp: &uniq gpu.shared [[f32;32];32]  
5 ) -[grid: gpu.grid<XY<64,64>,XY<32,8>>-> () {  
6   sched(Y,X) block in grid {  
7     sched(Y,X) thread in block {  
8       for i in 0..4 {  
9         tmp.group_by_row::<32,4>[[thread]][i] =  
10        input.group_by_tile::<32,32>.transpose[[block]].group_by_row::<32,4>[[thread]][i]  
11      };  
12  
13      for i in 0..4 {  
14        output.group_by_tile::<32,32>[[block]].group_by_row::<32,4>[[thread]][i] =  
15        tmp.transpose.group_by_row::<32,4>[[thread]][i]  
16      }  
17    }  
18  }  
19 }
```

threads read safely from memory tmp

```
1 fn transpose(  
2   input: &gpu.global [[f32;2048];2048],  
3   output: &uniq gpu.global [[f32;2048];2048],  
4   [grid.blocks.forall(X).forall(Y)] tmp: &uniq gpu.shared [[f32;32];32]  
5 ) -[grid: gpu.grid<XY<64,64>,XY<32,8>>-> () {  
6   sched(Y,X) block in grid {  
7     sched(Y,X) thread in block {  
8       for i in 0..4 {  
9         tmp.group_by_row::<32,4>[[thread]][i] =  
10        input.group_by_tile::<32,32>.transpose[[block]].group_by_row::<32,4>[[thread]][i]  
11      };  
12  
13      for i in 0..4 {  
14        output.group_by_tile::<32,32>[[block]].group_by_row::<32,4>[[thread]][i] =  
15        tmp.transpose.group_by_row::<32,4>[[thread]][i]  
16      }  
17    }  
18  }  
19 }
```

threads read safely from memory tmp

```

1  fn transpose(
2      input: &gpu.global [[f32;2048];2048],
3      output: &uniq gpu.global [[f32;2048];2048],
4      [grid.blocks.forall(X).forall(Y)] tmp: &uniq gpu.shared [[f32;32];32]
5  ) -[grid: gpu.grid<XY<64,64>,XY<32,8>>-> () {
6      sched(Y,X) block in grid {
7          sched
8          for
9              t
10
11      };
12
13      for
14          o
15
16      }
17  }
18  }
19  }

```

missing synchronization leads to an error

Descend

error: conflicting memory access

| line 15: tmp.transpose.group\_by\_row::<32,4>[[thread]][i]


| .....  
|

| line 9: tmp.group\_by\_row::<32,4>[[thread]][i]

| .....  
|



```
1 fn transpose(  
2   input: &gpu.global [[f32;2048];2048],  
3   output: &uniq gpu.global [[f32;2048];2048],  
4   [grid.blocks.forall(X).forall(Y)] tmp: &uniq gpu.shared [[f32;32];32]  
5 ) -[grid: gpu.grid<XY<64,64>,XY<32,8>>-> () {  
6   sched(Y,X) block in grid {  
7     sched(Y,X) thread in block {  
8       for i in 0..4 {  
9         tmp.group_by_row::<32,4>[[thread]][i] =  
10        input.group_by_tile::<32,32>.transpose[[block]].group_by_row::<32,4>[[thread]][i]  
11      };  
12      sync(block);  
13      for i in 0..4 {  
14        output.group_by_tile::<32,32>[[block]].group_by_row::<32,4>[[thread]][i] =  
15        tmp.transpose.group_by_row::<32,4>[[thread]][i]  
16      }  
17    }  
18  }  
19 }
```



explicit barrier synchronization on the entire block

```
1 fn transpose(  
2   input: &gpu.global [[f32;2048];2048],  
3   output: &uniq gpu.global [[f32;2048];2048],  
4   [grid.blocks.forall(X).forall(Y)] tmp: &uniq gpu.shared [[f32;32];32]  
5 ) -[grid: gpu.grid<XY<64,64>,XY<32,8>>-> () {  
6   sched(Y,X) block in grid {  
7     sched(Y,X) thread in block {  
8       for i in 0..4 {  
9         tmp.group_by_row:<32,4>[[thread]][i] =  
10        input.group_by_tile:<32,32>.transpose[[block]].group_by_row:<32,4>[[thread]][i]  
11      };  
12      sync(block);  
13      for i in 0..4 {  
14        output.group_by_tile:<32,32>[[block]].group_by_row:<32,4>[[thread]][i] =  
15        tmp.transpose.group_by_row:<32,4>[[thread]][i]  
16      }  
17    }  
18  }  
19 }
```

**Descend ensures that there are not data races**

# Current State of Descend

Further guarantees provided by Descend:

- Respect separation between CPU & different GPU memories
- Assumptions between CPU and GPU code match
- Safe memory accesses in branching code
- Absence of barrier divergence

What we did:

- Formalized the type system
- Implemented compiler to CUDA

## Descend: A Safe GPU Systems Programming Language

BASTIAN KÖPCKE, University of Münster, Germany  
SERGEI GORLATCH, University of Münster, Germany  
MICHEL STEUWER, Technische Universität Berlin, Germany

Graphics Processing Units (GPU) offer tremendous computational power by following a throughput oriented paradigm where many thousand computational units operate in parallel. Programming such massively parallel hardware is challenging. Programmers must correctly and efficiently coordinate thousands of threads and their accesses to various shared memory spaces. Existing mainstream GPU programming languages, such as CUDA and OpenCL, are based on C/C++ inheriting their fundamentally unsafe ways to access memory via raw pointers. This facilitates easy to make, but hard to detect bugs, such as data races and deadlocks.

In this paper, we present Descend: a safe GPU programming language. In contrast to prior safe high-level GPU programming approaches, Descend is an imperative GPU systems programming language in the spirit of Rust, enforcing safe CPU and GPU memory management in the type system by tracking Ownership and Lifetimes. Descend introduces a new holistic GPU programming model where computations are hierarchically scheduled over the GPU's execution resources: grid, blocks, warps, and threads. Descend's extended borrow checking ensures that execution resources safely access memory regions without data races. For this, we introduced views describing safe parallel access patterns of memory regions, as well as atomic variables. For memory accesses that can't be checked by our type system, users can annotate limited code sections as unsafe.

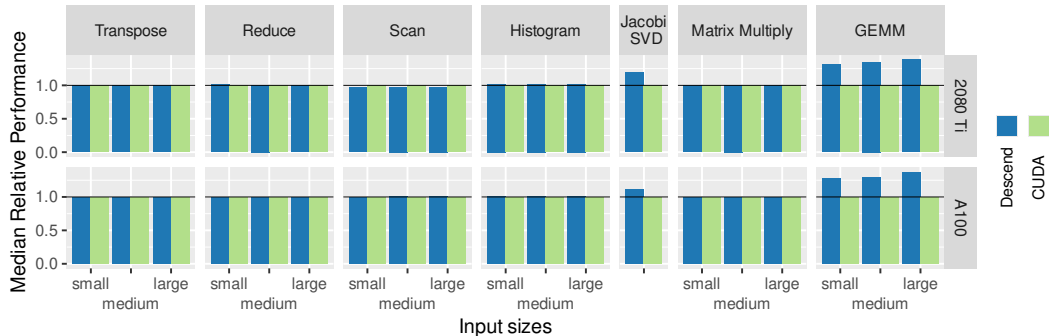
We discuss the memory safety guarantees offered by Descend and evaluate our implementation using multiple benchmarks, demonstrating that Descend is capable of expressing real-world GPU programs showing competitive performance compared to manually written CUDA programs lacking Descend's safety guarantees.

CCS Concepts • Software and its engineering → Parallel programming languages, Imperative languages.

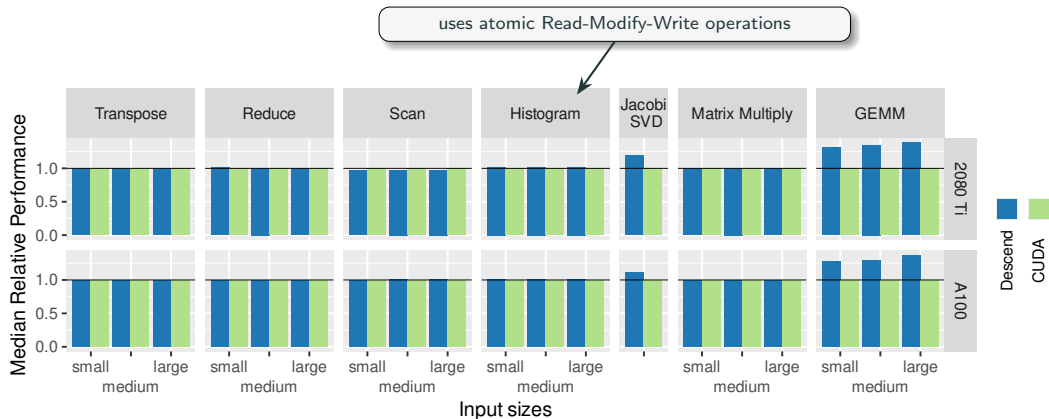
## 3 Typing Rules

$\Delta; \Gamma \mid e; A \mid t \vdash \boxed{\delta} : \Gamma' \dashv A'$	
<p><b>T-READ-LOCAL</b></p> $\begin{array}{c} \text{isPlace}(p) \\ \Delta; \Gamma \mid y; e; e \mid A \vdash [t; \delta]_p \dashv \Gamma' \mid A' \\ \Gamma'(p) = (\delta_p, \tau_p) \quad e = \tau_p \\ \Delta; \Gamma' \vdash^{wz} \delta_p \rightarrow \delta_p \dashv \Gamma' \\ \Gamma' \vdash^{wz} \text{borrow } p \rightarrow \{ \text{wz} \} \end{array}$ $\Delta; \Gamma \mid y; e; e \mid A \vdash [p = t; \text{unit}] \dashv \Gamma' \mid [p \mapsto \delta] \dashv A'$	<p><b>T-WRITE-SHARED-MEM</b></p> $\begin{array}{c} \neg \text{isPlace}(p) \quad y; e; t; e; \text{GpuThread} \\ \Delta; \Gamma \mid y; e; e \mid A \vdash [t; \delta]_p \dashv \Gamma' \mid A' \\ \Delta; \Gamma' \mid y; e; e \vdash \tau_p \vdash^{wz} p : \delta_p \\ \Delta; \Gamma' \vdash^{wz} \delta_p \rightarrow \delta_p \dashv \Gamma' \\ \Delta; \Gamma' \vdash^{wz} \text{borrow}^+ p \rightarrow \{ \text{wz} \} \end{array}$ $\Delta; \Gamma \mid y; e; e \mid A \vdash [p = t; \text{unit}] \dashv \Gamma' \mid A'$
<p><b>T-READ-BY-COPY</b></p> $\begin{array}{c} \text{isPlace}(p) \quad \text{isCopyable}(\delta) \\ \Delta; \Gamma \mid y; e; e \vdash \tau_p \vdash^{wz} p : \delta \\ \Delta; \Gamma \mid y; e; e \mid A \vdash^{wz} \text{borrow}^+ p \rightarrow \{ \text{wz} \} \\ \Delta; \Gamma \mid y; e; e \mid A \vdash [p; \delta] \dashv \Gamma \mid A \end{array}$	<p><b>T-READ-BY-MOVE</b></p> $\begin{array}{c} \text{isPlace}(p) \quad \neg \text{isCopyable}(\delta) \\ \Delta; \Gamma \mid y; e; e \vdash \tau_p \vdash^{wz} p : \delta \\ \Delta; \Gamma \mid y; e; e \mid A \vdash^{wz} \text{borrow}^+ p \rightarrow \{ \text{wz} \} \\ \Delta; \Gamma \mid y; e; e \mid A \vdash [p; \delta] \dashv \Gamma \mid [p \mapsto \delta'] \mid A \end{array}$
<p><b>T-READ-SHARED-MEM</b></p> $\begin{array}{c} \neg \text{isPlace}(p) \quad \text{isCopyable}(\delta) \\ \Delta; \Gamma \mid y; e; e \vdash \tau_p \vdash^{wz} p : \delta \\ \Delta; \Gamma \mid y; e; e \mid A \vdash^{wz} \text{borrow}^+ p \rightarrow \{ \text{wz} \} \\ \Delta; \Gamma \mid y; e; e \mid A \vdash [p; \delta] \dashv \Gamma \mid A, \delta \end{array}$	<p><b>T-BORROW</b></p> $\begin{array}{c} \neg \text{isPlace}(p) \quad \Gamma(r) = \emptyset \\ \Delta; \Gamma \mid y; e; e \vdash \tau_p \vdash^{wz} p : \delta, \mu \\ \Delta; \Gamma \mid y; e; e \mid A \vdash^{wz} \text{borrow}^+ p \rightarrow \{ \text{wz} \} \\ \Delta; \Gamma \mid y; e; e \mid A \vdash [kr \omega \mu; \delta] \dashv \Gamma \mid A \end{array}$
<p><b>T-LET</b></p> $\begin{array}{c} \Delta; \Gamma' \vdash^{wz} \delta' \rightarrow \delta \dashv \Gamma', (f) \quad \forall r \in \text{free-regions}(\delta), \Gamma', (f) \vdash r \text{ rnz} \\ \Delta; \Gamma \vdash \delta : \text{dty} \quad \Delta; \Gamma \mid y; e; e \mid A \vdash [t; \delta'] \dashv \Gamma' \mid A' \\ \Delta; \Gamma \mid y; e; e \mid A \vdash [\text{let } x; \delta = t; \text{unit}] \dashv \Gamma', (f, x; e; e) \mid A' \end{array}$	

# Evaluation

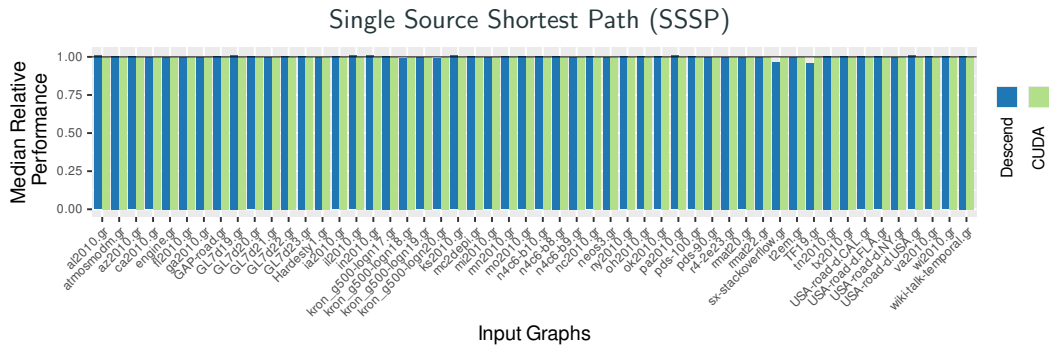


Programs generated from Descend achieve competitive performance to handwritten CUDA.



Programs generated from Descend achieve competitive performance to handwritten CUDA.

# unsafe: Interfacing with CUDA



unsafe code provides a quick method of interfacing with existing CUDA code

# Conclusion

- *Descend* assists programmers in managing CPU and GPU memory and enforcing previously implicit assumptions
- *Descend* extends the concept of ownership to the thread hierarchy and uses views to ensure safe memory accesses on the GPU
- The evaluation shows that *Descend* is expressive enough to write programs that achieve performance on-par with handwritten CUDA

Website for Descend: <https://descend-lang.org>

bastian.koepcke@uni-muenster.de