

Scan Stack: A Search-based Concurrent Stack for GPU

Implementation Report

Aditya Azad (241110008)

Sparsh Sehgal (24111071)

Vraj Patel (241110080)

April 21, 2025

Abstract

In this report, we present our implementation of the Scan Stack, a concurrent search-based GPU stack based on the work by Noah South. The Scan Stack is designed to take advantage of GPU memory access patterns, memory coalescence, and thread structures (warps) to increase throughput. Our implementation focuses on the core structure and operations of the stack, including push and pop operations, as well as optimization techniques like elimination and lowest-area scanning. We evaluate the implementation through various tests and analyze its performance characteristics.

Contents

1	Introduction	3
2	Base Structure	3
2.1	Structure Definition	3
2.2	Memory Access Pattern	4
2.3	Lowest-Area Optimization	4
3	Stack Operations	4
3.1	Push Operation	4
3.2	Pop Operation	4
4	Contention Handling	6
4.1	Basic Contention Handling	6
4.2	Elimination Techniques	6
4.2.1	Warp- and Block-Level Elimination	6
4.2.2	Grid-Wide Elimination-Backoff	8
4.2.3	Benefits and Trade-offs	8
5	Performance Analysis	9
5.1	Throughput Analysis	9
5.2	Factors Affecting Performance	9
6	Implementation Challenges	9
6.1	The Step Over Problem	9
6.2	Avoiding Race Conditions	10
6.3	Balancing Elimination and Direct Access	10
6.4	Memory Coalescing Optimization	10
7	Conclusion	10

1 Introduction

Concurrent data structures play a crucial role in the overall performance of GPGPU applications. The stack, one of the most fundamental data structures, finds numerous applications where data is processed in a Last In First Out (LIFO) fashion. Although concurrent stacks are well researched for multi-core CPUs, there is less research on GPU-friendly implementations.

In this report, we detail our implementation of the Scan Stack, a search-based concurrent stack for GPUs as proposed by Noah South. The key features of our implementation include:

- A search-based approach that eliminates the need for a global top pointer
- Optimization of the search area to improve efficiency
- GPU-friendly push and pop operations
- Contention handling mechanisms including elimination

The implementation is designed to take advantage of GPU memory access patterns and the SIMT (Single Instruction Multiple Thread) execution model, providing high throughput even under heavy contention.

2 Base Structure

The Scan Stack is based on a one-dimensional array structure, making it well-suited for GPU implementation. The array structure offers benefits of GPU-friendly memory access patterns when threads access memory in a coalesced manner.

2.1 Structure Definition

The core structure of our implementation is defined as follows:

```
1 struct ScanStack
2 {
3     int *array;
4     int capacity;
5
6     __device__ void init()
7     {
8         for (int i = 0; i < capacity; ++i)
9             {
10                array[i] = EMPTY;
11            }
12    }
13    // Push and Pop operations are implemented here
14};
```

Key constants used in the implementation:

```
1 __device__ const int EMPTY = -1;
2 __device__ const int INVALID = -2;
3
4 #define LA_GRANULARITY 1
```

The array represents the stack with special values to indicate:

- EMPTY (-1): Indicates an empty cell
- INVALID (-2): Used during pop operations to mark cells that are being processed

2.2 Memory Access Pattern

The implementation takes advantage of GPU memory coalescing by structuring operations to access memory in patterns that benefit from this hardware feature. When multiple threads within a warp access sequential memory locations, the GPU can combine these into fewer memory transactions, significantly improving performance.

2.3 Lowest-Area Optimization

To improve the efficiency of stack operations, we implemented the "lowest necessary area" optimization. This technique reduces the search space for operations by scanning at intervals to find an approximate location of the stack's top.

```
1 // Example from push operation
2 int startIndex = 0;
3 for (int idx = LA_GRANULARITY; idx < capacity; idx += LA_GRANULARITY)
4 {
5     if (array[idx] == EMPTY)
6     {
7         startIndex = idx - LA_GRANULARITY;
8         if (startIndex < 0)
9             startIndex = 0;
10        break;
11    }
12 }
13 if (startIndex == 0)
14 {
15     startIndex = capacity - 1;
16 }
```

This optimization significantly reduces the number of memory accesses needed to find the top of the stack, especially when the stack is mostly full (for push operations) or mostly empty (for pop operations).

3 Stack Operations

3.1 Push Operation

The push operation inserts an integer value into the stack. The implementation follows these steps:

The key aspects of the push operation include:

- Scanning from the bottom of the stack upwards until an empty cell is found
- Using atomicCAS to safely insert the value
- Handling contention by backing off and rescanning if insertion fails

3.2 Pop Operation

The pop operation removes and returns the top value from the stack:

Key aspects of the pop operation:

- Scanning from the top of the stack downwards until a non-empty cell is found
- Using atomicCAS to safely extract the value and mark the cell as invalidated
- Handling invalid cells by re-validating them when appropriate

Algorithm 1 Scan Stack Push Operation

```
1: Find starting index using lowest-area optimization
2: for  $i \leftarrow startIndex$  to  $capacity - 1$  do
3:   if  $array[i] = EMPTY$  then
4:      $observed \leftarrow atomicCAS(\&array[i], EMPTY, value)$ 
5:     if  $observed = EMPTY$  then
6:       return  $value$  ▷ Push successful
7:     else
8:        $i \leftarrow i - 1$ 
9:       while  $array[i] < 0$  and  $i > 0$  do
10:         $i \leftarrow i - 1$ 
11:      end while
12:      Set flag to retry scan
13:    end if
14:  end if
15: end for
16: return  $INVALID$  ▷ Stack is full or other failure
```

Algorithm 2 Scan Stack Pop Operation

```
1: Find starting index using lowest-area optimization
2: for  $i \leftarrow startIndex$  to 0 do
3:   if  $array[i] = INVALID$  and  $array[i + 1] = EMPTY$  then
4:      $atomicCAS(\&array[i], INVALID, EMPTY)$  ▷ Re-validate empty cell
5:     Adjust scan position
6:   end if
7:   if  $array[i] > -1$  then ▷ Non-empty cell found
8:      $value \leftarrow array[i]$ 
9:      $observed \leftarrow atomicCAS(\&array[i], value, INVALID)$ 
10:    if  $observed = value$  then
11:      return  $value$  ▷ Pop successful
12:    else
13:      Adjust scan position and retry
14:    end if
15:  end if
16: end for
17: return  $INVALID$  ▷ Stack is empty or other failure
```

- Handling contention when multiple threads attempt to pop the same value

4 Contention Handling

4.1 Basic Contention Handling

When multiple threads attempt to access the top of the stack simultaneously, contention can significantly impact performance. Our implementation handles this in several ways:

- The scanning approach itself distributes contention by allowing threads to independently locate the top
- Failed CAS operations are handled by backing off and rescanning
- Invalid cells are used to prevent race conditions between push and pop operations

```

1 // Example of contention handling in push
2 if (observed == EMPTY)
3 {
4     int old = atomicCAS(&array[i], EMPTY, val);
5     if (old == EMPTY)
6     {
7         return val;
8     }
9     else
10    {
11        int j = i - 1;
12        while (j >= 0 && array[j] < 0)
13        {
14            j--;
15        }
16        if (j < 0)
17        {
18            j = 0;
19        }
20        i = j;
21        retrace = true;
22    }
23 }

```

4.2 Elimination Techniques

To minimize contention on the central stack, we combine two complementary elimination strategies: a lightweight “local” scheme that works entirely within a CUDA block, and a more powerful “global” elimination-backoff that spans the entire grid.

4.2.1 Warp- and Block-Level Elimination

This GPU-native optimization pairs complementary operations (push and pop) within the same warp or block using only shared memory and simple scans—no expensive atomics or backoff. In practice, the first lane of each warp scans its 32 entries for a matching push/pop, exchanges values, and clears both requests. Any unmatched pushes then let a single leader scan the entire block (via one `atomicCAS`) to match remaining pops.

Kernel pseudocode

```
1 // within performOperations kernel, after loading opType, opValue...
2 int warpId = tid / warpSize;
3 int laneId = tid % warpSize;
4 __syncthreads();
5
6 // 1) Warp-local elimination
7 if (laneId == 0) {
8     for (int k = warpId*warpSize; k < min((warpId+1)*warpSize,
9         blockDim.x); ++k) {
10         if (opType[k] == OP_PUSH) {
11             for (int m = warpId*warpSize; m < min((warpId+1)*warpSize,
12                 blockDim.x); ++m) {
13                 if (opType[m] == OP_POP) {
14                     ops[m].result = opValue[k];
15                     ops[k].result = opValue[k];
16                     opType[k] = OP_NONE;
17                     opType[m] = OP_NONE;
18                     break;
19                 }
20             }
21         }
22     }
23     __syncthreads();
24
25 // 2) Block-wide elimination
26 if (laneId == 0) {
27     int pushIdx = -1;
28     for (int k = warpId*warpSize; k < min((warpId+1)*warpSize,
29         blockDim.x); ++k)
30         if (opType[k] == OP_PUSH) { pushIdx = k; break; }
31
32     if (pushIdx >= 0) {
33         int pv = opValue[pushIdx];
34         opType[pushIdx] = 2; // in-progress
35         for (int j = 0; j < blockDim.x; ++j) {
36             if (opType[j] == OP_POP) {
37                 if (atomicCAS(&opType[j], OP_POP, OP_NONE) == OP_POP) {
38                     ops[j].result = pv;
39                     ops[pushIdx].result = pv;
40                     opType[pushIdx] = OP_NONE;
41                     break;
42                 }
43             }
44         }
45         if (opType[pushIdx] == 2)
46             opType[pushIdx] = OP_PUSH; // no match found
47     }
48 }
49 __syncthreads();
```

This local scheme is essentially free when matches occur, but it can only eliminate within a single block of threads.

4.2.2 Grid-Wide Elimination-Backoff

To handle contention across the entire GPU grid (and adapt dynamically to varying loads), we implemented the lock-free, linearizable *elimination-backoff* algorithm. Here's how it works:

- **Global arrays** (`d_location[]` and `d_collision[]`) live in device memory, sized for all threads.
- After a failed direct push/pop CAS on stack, a thread publishes its intent in `d_location[tid]`, picks a random slot in `d_collision`, and does an `atomicExch` to meet a partner.
- If it finds an opposite operation, they exchange values via a small CAS on `d_location`, return immediately, and *never* touch the stack.
- If no partner arrives, the thread “back-off” spins for an exponentially growing delay, then retries the stack CAS. This both spreads out retries and increases the chance of elimination as contention grows.

Device wrapper pseudocode

```
1  __device__ int eliminateOrStackOp(ScanStack *stack,
2                                     int op, int val, int tid) {
3      // 1) Try direct on central stack
4      int r = (op==OP_PUSH) ? stack->push(val)
5                          : stack->pop();
6      if (r != INVALID_STACK) return r;
7
8      // 2) Elimination-backoff loop
9      ThreadInfo ti = {tid, op, val, 1u};
10     while (true) {
11         d_location[tid] = ti;
12         int slot = randomSlot(tid);
13         int him = atomicExch(&d_collision[slot], tid);
14
15         if (him != EMPTY_ELIM &&
16             d_location[him].op == -op) {
17             // complementary op found
18             if (op==OP_PUSH) return val;
19             else return d_location[him].value;
20         }
21
22         // backoff spin
23         __nanosleep(ti.spin);
24         ti.spin = min(ti.spin*2, MAX_BACKOFF);
25
26         // retry central stack
27         r = (op==OP_PUSH) ? stack->push(val)
28                          : stack->pop();
29         if (r != INVALID_STACK) return r;
30     }
31 }
```

In the final stage of the kernel, any operation that survived block elimination calls `eliminateOrStackOp` instead of directly popping or pushing.

4.2.3 Benefits and Trade-offs

- Warp/block elimination is ultra-low-overhead but limited to a block.

- Grid-wide elimination-backoff scales across all threads, adapts to contention, and is formally proven linearizable/lock-free, but at the cost of additional atomics and global memory.

Together, they deliver both low latency under light load and high throughput under heavy contention.”

5 Performance Analysis

5.1 Throughput Analysis

Based on the Scan Stack research paper, we can expect the following performance characteristics from our implementation:

- The elimination technique significantly improves throughput, especially when there is a mix of push and pop operations
- The lowest-area optimization reduces scanning overhead, especially for large stacks
- The implementation scales well with increasing operation counts

5.2 Factors Affecting Performance

Several factors can influence the performance of the Scan Stack:

- **LA_GRANULARITY:** The interval size for the lowest-area scan affects the trade-off between scan overhead and accuracy
- **Thread block size:** The size of thread blocks affects the effectiveness of elimination
- **Operation mix:** A balanced mix of push and pop operations allows for more elimination opportunities
- **Stack size:** The capacity of the stack impacts the scanning overhead

6 Implementation Challenges

During our implementation, we encountered several challenges:

6.1 The Step Over Problem

One of the most significant challenges in concurrent stack implementations is the ”step over” problem. This occurs when threads scanning for the top of the stack pass each other due to their different scanning directions and conditions.

In the Scan Stack, this problem can manifest when:

- A push operation is scanning upward looking for an empty cell
- Simultaneously, a pop operation is scanning downward looking for a non-empty cell
- The operations ”pass” each other, potentially causing the stack to enter an invalid state

Figure 1 illustrates this problem:

Our implementation addresses this problem through:

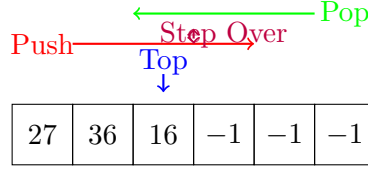


Figure 1: Illustration of the step over problem in Scan Stack

- Invalidation cells: When a pop operation removes a value, it temporarily marks the cell as INVALID (-2) before it can be marked as EMPTY (-1)
- Retracing and scanning logic: When operations detect they may have passed the true top, they adjust their scanning position
- Atomic operations: Using atomicCAS ensures that only one thread succeeds when multiple threads contend for the same cell

The invalidation mechanism is particularly important as it prevents push operations from inserting values into cells that might create "gaps" in the stack structure, thus maintaining stack integrity.

6.2 Avoiding Race Conditions

Beyond the step over problem, ensuring the stack remains in a valid state despite concurrent access presents additional challenges. The use of atomic operations and invalidation mechanisms was crucial to prevent various race conditions that could occur during concurrent operations.

6.3 Balancing Elimination and Direct Access

Finding the right balance between elimination and direct stack access was challenging. While elimination reduces contention, it also introduces overhead for synchronization and shared memory operations. Our implementation had to carefully consider when elimination is beneficial versus direct stack access.

6.4 Memory Coalescing Optimization

Ensuring memory access patterns that benefit from coalescing required careful structuring of the code. The lowest-area optimization had to be implemented in a way that maintains coalesced access patterns while reducing scan overhead.

7 Conclusion

Our implementation of the Scan Stack provides an efficient concurrent stack for GPU applications. The design takes advantage of GPU-specific features such as memory coalescing and the warp execution model, while handling the challenges of concurrent access through non-blocking techniques.

Key contributions of our implementation include:

- A fully functional implementation of the search-based Scan Stack
- Incorporation of the lowest-area optimization for efficient top location
- Implementation of elimination techniques to reduce contention

- A design that scales well with increasing operation counts

The Scan Stack demonstrates that array-based concurrent data structures can be effectively implemented on GPUs by leveraging their unique architectural features. The combination of search-based operations and elimination techniques provides a scalable solution for high-throughput stack operations in GPU-parallel applications.

The source code for this implementation is available on GitHub at:

<https://github.com/Vrajb24/ScanStack-Implementation.git>

To compile and run the Scan Stack code, use the following command:

```
1 nvcc -O2 --allow-unsupported-compiler -ccbin=g++-9 -std=c++17
  -arch=sm_86 -lineinfo -res-usage -src-in-ptx Scan_stack.cu -o
  output && ./output
```

To compile and run the CAS Stack code, use the following command:

```
1 nvcc -O2 --allow-unsupported-compiler -ccbin=g++-9 -std=c++17
  -arch=sm_86 -lineinfo -res-usage -src-in-ptx CAS_stack.cu -o output
  && ./output
```

8 Future Work

Potential areas for further improvement include:

- Dynamic adjustment of LA_GRANULARITY based on stack state
- Exploration of alternative elimination strategies
- Extension to support multiple data types
- Integration with other GPU data structures for complex applications

References

- [1] South, N. (2022). *Scan Stack: A Search-based Concurrent Stack for GPU*. University of Mississippi.
- [2] Hendler, D., Shavit, N., & Yerushalmi, L. (2004). *A scalable lock-free stack algorithm*. In Proceedings of the sixteenth annual ACM symposium on Parallelism in algorithms and architectures (pp. 206-215).
- [3] Misra, P., & Chaudhuri, M. (2012). *Performance evaluation of concurrent lock-free data structures on GPUs*. In 2012 IEEE 18th International Conference on Parallel and Distributed Systems (pp. 53-60). IEEE.