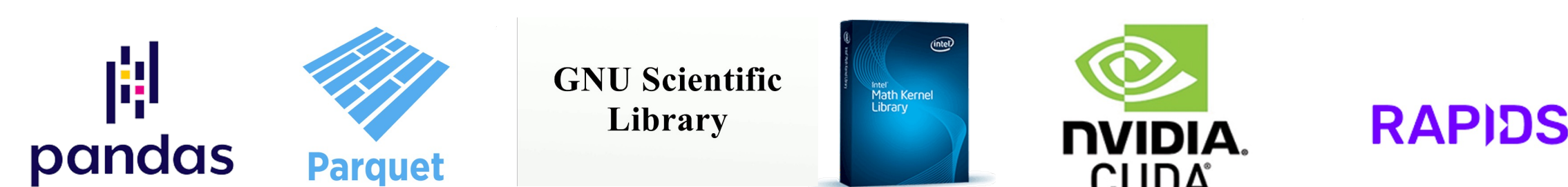


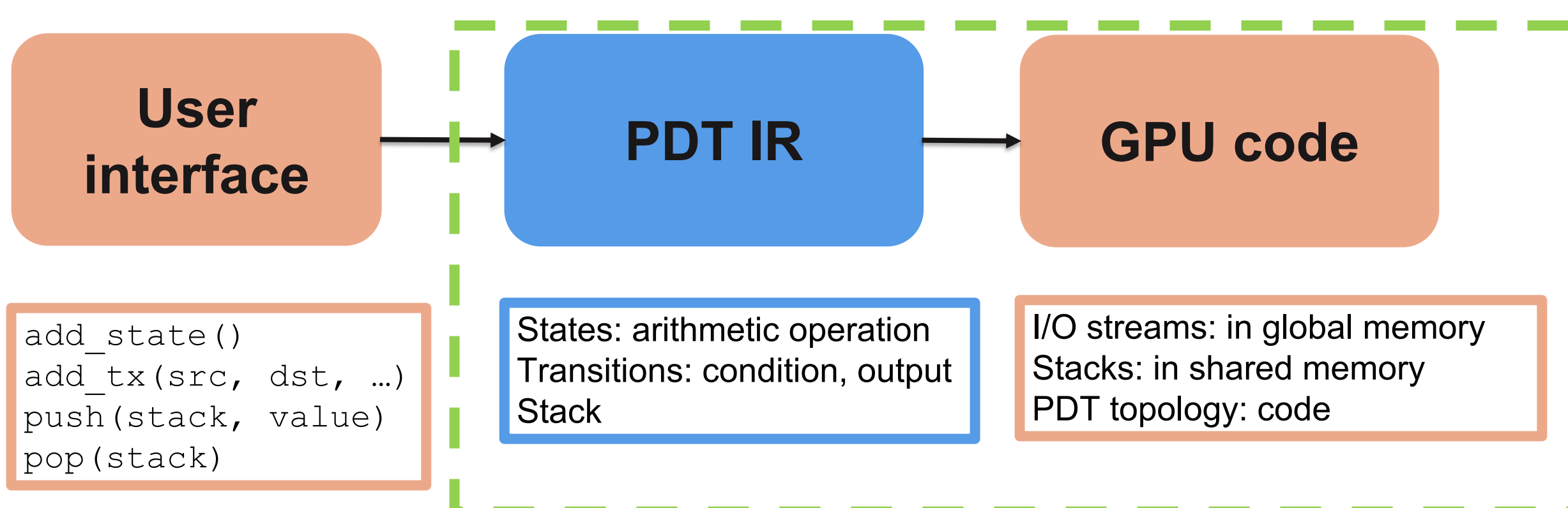
Motivation

- Data transformation kernels are at the core of many data analytics, data processing and scientific applications
- Custom CPU and GPU libraries and hardware accelerators for specific data transformation tasks provide efficiency but lack generality and extensibility
- Accelerating the computational abstraction at the core of data transformation can benefit a broad spectrum of applications



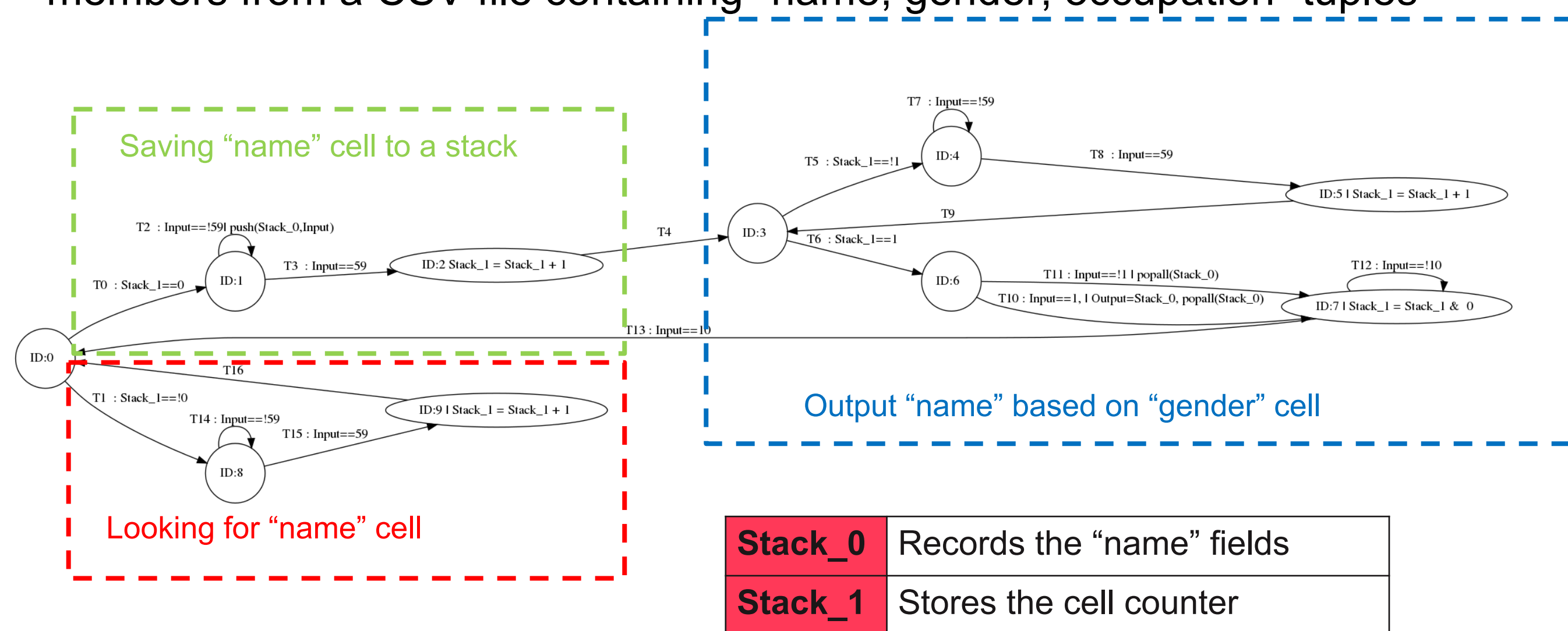
Objective

- Compilation framework that generates accelerated GPU kernels for data transformation tasks expressed using Pushdown Transducers (PDTs)



Pushdown Transducers (PDT)

Example (CSV parsing): PDT that extracts the name of all female members from a CSV file containing "name, gender, occupation" tuples



Conversion Steps

- Convert the list of states into an *if-else* block in which each *if*-condition guards the content of a state
- For each state:
 - convert the associated arithmetic operation (if any) to the appropriate arithmetic operation on stack values
 - convert the list of outgoing transitions into an *if-else* block, with an *if*-statement per transition. Generate *if*-conditions and statements according to the appropriate input, output and stack read/write operation
- Allocate stacks in shared memory
- Store the remaining context information (pointers in I/O streams, active state indicator) into local variables to be stored in the register file

Basic Code

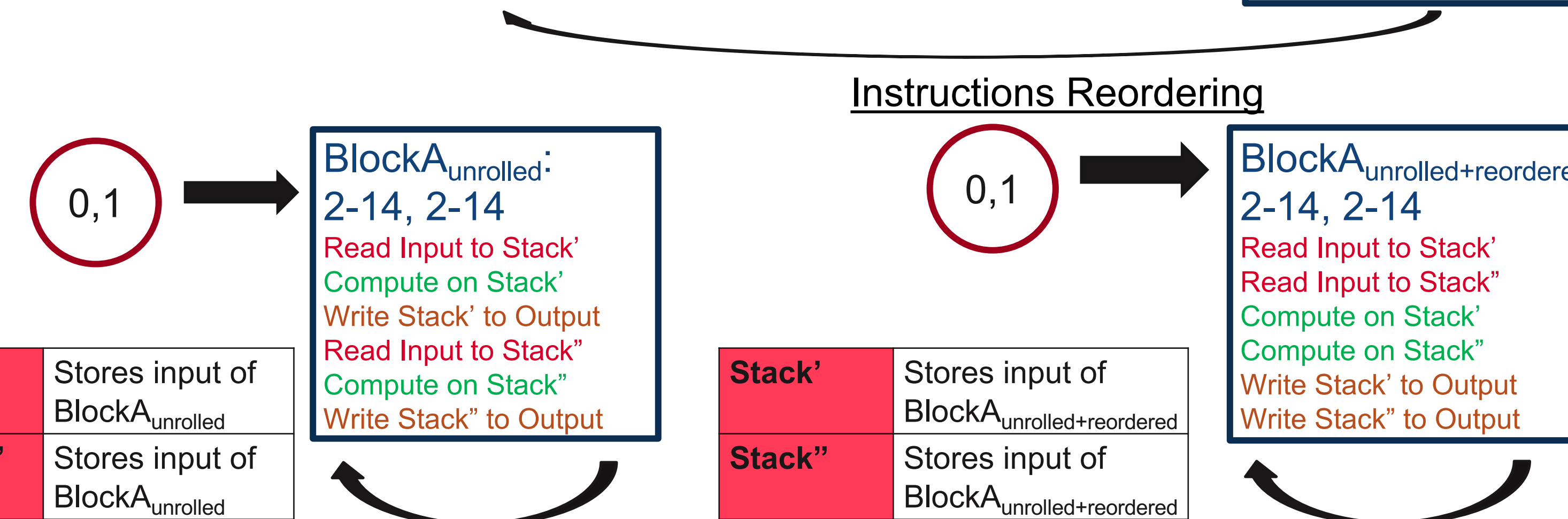
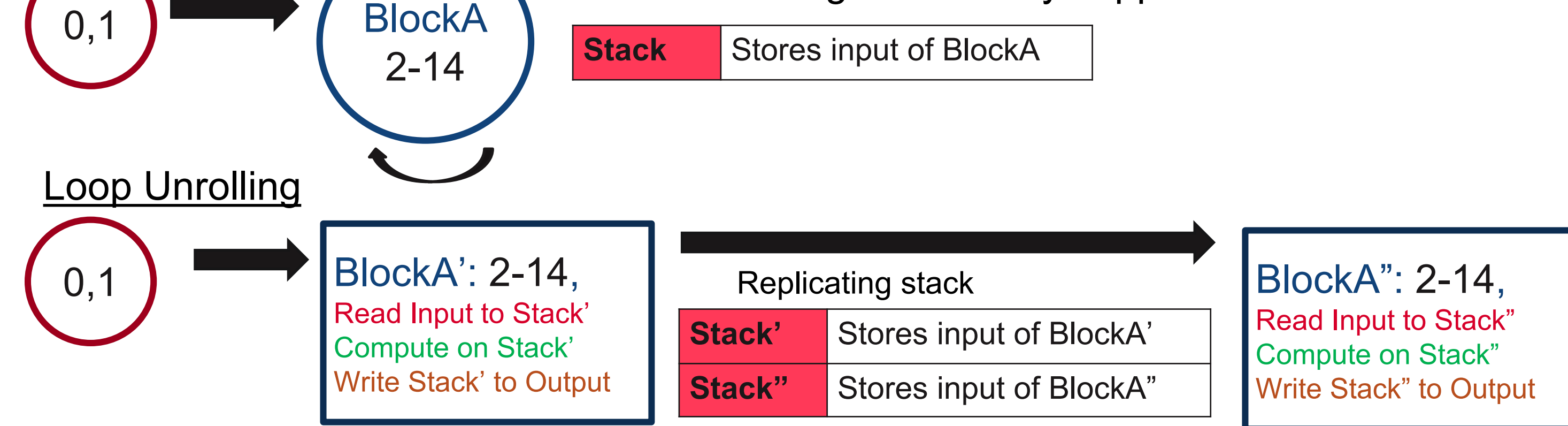
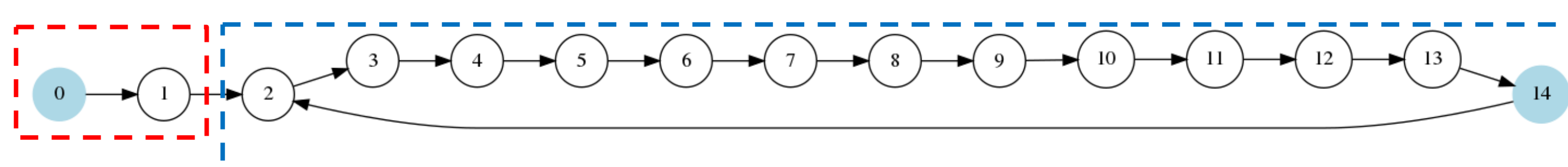
Snapshot of generated PDT traversal code

```

29 while ( !done) && ( currentState != 10){
30   if(currentState == 0){
31     if((var[(threadIdx.x * 2 + 1) * STACK_DEPTH + top_1-1] == 0)){ Tx 0
32       currentState = 1;
33     }
34     else if((var[(threadIdx.x * 2 + 1) * STACK_DEPTH + top_1-1] != 0)){ Tx 1
35       currentState = 0;
36     }
37   }
38   else if ( currentState == 1){
39     if((input[0][base_0 + processed_0] != 59)){ Tx 2
40       var[(threadIdx.x * 2 + 0) * STACK_DEPTH + top_0-1] = input[0][base_0 + processed_0];
41       top_0++;
42       currentState = 1;
43     }
44     else if((input[0][base_0 + processed_0] == 59)){ Tx 3
45       currentState = 2;
46     }
47     processed_0++;
48   }
49   else if ( currentState == 2){
50     var[(threadIdx.x * 2 + 1) * STACK_DEPTH + top_1-1] = var[(threadIdx.x * 2 + 0) * STACK_DEPTH + top_0-1];
51     currentState = 3;
52   }
53   else if ( currentState == 3){
54     if((var[(threadIdx.x * 2 + 1) * STACK_DEPTH + top_1-1] != 6)){ Tx 5
55       currentState = 4;
56     }
57     else if((var[(threadIdx.x * 2 + 1) * STACK_DEPTH + top_1-1] == 6)){ Tx 7
58       currentState = 0;
59     }
60   }
61   else if ( currentState == 4){
62     if((input[0][base_0 + processed_0] != 59)){ Tx 8
63       currentState = 4;
64     }
65     else if((input[0][base_0 + processed_0] == 59)){ Tx 8
66       currentState = 5;
67     }
68     processed_0++;
69   }
    
```

Optimizations

- Goals:**
- Reduce code size
 - Reduce control flow operations
 - Improve memory access patterns
- No optimizations**



- Advantages: Loop unrolling enables more aggressive optimizations such as reordering of reads (from input) and writes (to output)
- Disadvantages: Loop unrolling increases register pressure

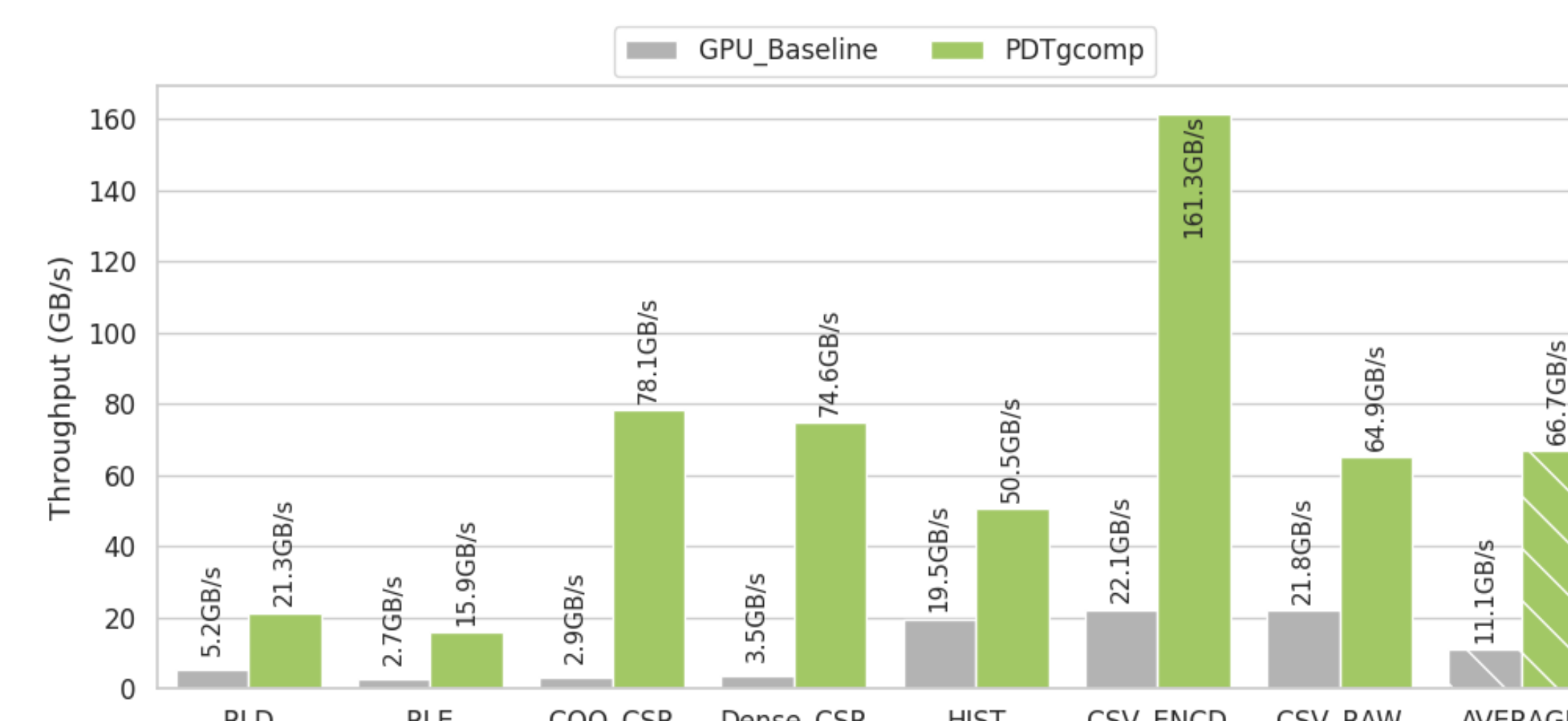
Experimental Setup

Transformation class	Input dataset	CPU Baseline	GPU Baseline
Data Enc/Dec	Cantebery Corpus	Parquet [1]	Nvidia Thrust [4]
Matrix Transformation	Texas A&M Sparse Matrix	Intel MKL [3]	Nvidia cuSparse [4]
Histogram	RDU Accident and Crime Report	GSL Histogram [2]	Nvidia Cub [4]
CSV Query	RDU Accident and Crime Report	Pandas [5]	Rapids AI [6]
System			
CPU	2x Intel Xeon E5-2630 2.2GHz		
GPU	NVIDIA TITAN XP 12GB, 30SMs		
OS/CUDA	Ubuntu 18.04, CUDA toolkit 11.7		

Performance



PDTgcomp speedup over custom CPU libraries



Throughput of PDTgcomp and custom GPU libraries in GB/s

- Average 82x and 6x speedup over custom CPU and GPU implementations, respectively
- On CPU, highest speedup on matrix transformation (200x) and CSV querying (120x)
- On GPU, highest throughput on CSV querying (161GB/s)

Conclusion and Future Work

- We demonstrated a method to generate efficient GPU code implementing data transformation tasks expressed using PDT
- Future work encompasses more compiler optimizations to further improve performance.

References and Acknowledgement

[1] Apache Parquet: <https://parquet.apache.org>
 [2] GNU scientific library: <https://www.gnu.org>
 [3] Intel MKL: <https://www.intel.com/content/www/us/en/develop/documentation/get-started-with-mkl-for-dpppp/top.html>
 [4] Cuda toolkit: <https://docs.nvidia.com/cuda>
 [5] Pandas: <https://pandas.pydata.org>
 [6] Open GPU data science: <https://rapids.ai>
 This work was supported by National Science Foundation awards CNS-1812727 and CCF-1907863.