

# Parallel Pk-Pk Join Algorithm Using GPU Servers

-by Abdhesh Dash

-Under the guidance of Prof. Ke Yi (CSE, HKUST)

GitHub- [MPC-PkPkJoin](#)

## **Abstract**

This report details the design and implementation of an efficient parallel Primary key-to-Primary key (Pk-Pk) join algorithm using CUDA-C++. Leveraging the CUB and Thrust libraries, the algorithm incorporates techniques such as BSP (Bulk Synchronous Parallel) sorting and Prefix sum. The performance of the implementation is evaluated against the Crystal library by Anil Shanbhag using identical datasets, demonstrating notable improvements in efficiency and execution time. The results show that the GPU-based approach significantly accelerates the Join operation, making it a viable option for handling large-scale data in real-time applications.

## **Introduction**

The exponential growth of data in various domains necessitates the development of efficient data processing techniques. Database operations, particularly join operations, are critical in extracting meaningful information from vast datasets. Traditional CPU-based methods often fall short in handling large-scale data efficiently due to their limited parallel processing capabilities. In contrast, Graphics Processing Units (GPUs) offer substantial computational power and parallelism using parallel blocks and threads, making them suitable for data-intensive tasks. This project aims to harness the power of GPUs to implement an efficient Pk-Pk join algorithm using CUDA-C++.

## **Background of the Project and Related Work**

The utilization of GPUs for database operations is an active research area, driven by the need to enhance the performance of data processing tasks. The Crystal library by Anil Shanbhag is a significant contribution in this field, offering a collection of data processing primitives optimized by GPU execution. Crystal's modular approach and its extensible set of GPU kernels for SQL query execution provide a solid foundation for developing high-performance database operations. The paper on Output-Optimal Massively Parallel Algorithms for Similarity Joins by Xiao Hu, Prof. Ke Yi and Yufei Tao is like the base for this project. They have explained the ideas of exclusive prefix sum, segmented prefix sum and BSP sorting in detail in that paper. In fact Xiao Hu has guided me on how to implement this in CUDA and the professor himself has also helped. Our work builds on from scratch, focusing on optimizing the Pk-Pk join operation, a fundamental database join where the primary keys of two tables are matched to produce a combined output.

## Description of the Problem and Algorithm

The primary key-to-primary key (Pk-Pk) join operation involves matching the primary keys of two tables and combining the corresponding records. The challenge in implementing this on a GPU lies in efficiently managing memory access patterns and synchronization across threads. Our algorithm addresses these issues through a series of optimized CUDA kernels. The key steps include:

1. **Input Data:** The input data is the collection of the primary keys of both the tables in a single array.
2. **Data Partitioning:** The input array is divided into smaller partitions to fit into the GPU's shared memory of each block.
3. **Radix Sort:** The data in each block is sorted by radix sort using the Block-wide radix sort function of CUB Library which does so in linear time.
4. **BSP Sorting:** Now, to get the whole dataset sorted, the Bulk Synchronous Parallel (BSP) approach is used, ensuring efficient sorting within shared memory.
5. **Evaluate Split values:** To get the global sorted data using the BSP approach, first its needed to get  $(p-1)$  splitters where ' $p$ ' is the number of blocks used. To get this,  $\log(p)$  samples from each block is selected randomly. These  $p \cdot \log(p)$  samples are put in a block and locally sorted. The split values are assigned to the equi-spaced values in this block starting from zero-th index.
6. **Distribution in Blocks:** Now that the splitter vales are got, its needed to distribute data of each block into some other blocks which will store the global sorted array. So the each data less than first split is put into first block, data less than first split is put into first block and so on. The block id for each element is got by simply using a binary search on the splitter values.
7. **Prefix Sum:** But the distribution can have data race because 2 threads may try to put data into the same block together. So, to avoid this, atomic operations can be used but that will increase the time complexity. Hence, prefix sum is used to count the occurrence of all previous key values and segmented prefix sum is computed for each value in the blocks to get the exact block id and thread id for each element.
8. **Assigning to Blocks:** The block id and thread id for each element is used to assign it to its corresponding thread and block. Now, the primary keys are globally sorted.
9. **Join Operation:** The globally sorted blocks are now traversed in parallel, and matching entries are got by only comparing the current element to its previous element. The result is written to the output array. The values corresponding to the matching pairs are got from the hash map of keys to values.

## Implementation Details

The implementation leverages several CUDA libraries, including CUB and Thrust, to optimize performance. The core components of the implementation are detailed below, along with code snippets illustrating key parts of the algorithm.

## CUDA Code Snippets:

### SortDataBlockWise

First the input array is divided into smaller blocks to fit into the GPU's shared memory. The data in each block is sorted by radix sort using the efficient Block-wide radix sort function of CUB Library which does so in linear time.

```
int thread_keys[ITEMS_PER_THREAD];
    int block_offset = blockIdx.x * (BLOCK_THREADS * ITEMS_PER_THREAD);
    int valid_items = num_elements - block_offset > BLOCK_THREADS *
ITEMS_PER_THREAD ? BLOCK_THREADS * ITEMS_PER_THREAD : num_elements -
block_offset;

    // Initialize thread_keys with a known value for safer debugging
    for (int i = 0; i < ITEMS_PER_THREAD; i++) {
        thread_keys[i] = (block_offset + threadIdx.x * ITEMS_PER_THREAD + i)
< num_elements ? d_in[block_offset + threadIdx.x * ITEMS_PER_THREAD + i] :
INT_MAX;
    }

    // Load data
    BlockLoadT(temp_storage.load).Load(d_in + block_offset, thread_keys,
valid_items);

    __syncthreads(); // Barrier for smem reuse

    // Collectively sort the keys
    BlockRadixSortT(temp_storage.sort).Sort(thread_keys);

    __syncthreads(); // Barrier for smem reuse

    // Store the sorted segment
    BlockStoreT(temp_storage.store).Store(d_out + block_offset, thread_keys,
valid_items);
```

### Find Splits

To get the global sorted data using the BSP approach, first its needed to get  $(p-1)$  splitters where ' $p$ ' is the number of blocks used. To get this,  $\log(p)$  samples from each block is selected randomly. These  $p * \log(p)$  samples are put in a block and locally sorted. The split values are assigned to the equi-spaced values in this block starting from the zero index.

```
__global__ void Splitterss(int* d_splitters,int* d_samples,int
sample_size,int p) {
    int tid = blockIdx.x * blockDim.x + threadIdx.x;
    if (tid < p-1) {
        d_splitters[tid] = d_samples[(tid + 1) * sample_size / p];
    }
}
```

## Assigning Blocks

It's now needed to distribute data of each block into some other blocks which will store the global sorted array. Like for example, data less than first split is put into first block, data less than first split is put into second block and so on. The block id for each element is got by simply using a binary search on the splitter values as is implemented below.

```
int tid = blockIdx.x * blockDim.x + threadIdx.x;
if (tid < numData) {
    int item = data[tid];
    int left = 0;
    int right = numSplitters;
    while (left < right) {
        int mid = left + (right - left) / 2;
        if (splitters[mid] > item) {
            right = mid;
        } else {
            left = mid + 1;
        }
    }
    output[tid] = (left < numSplitters) ? left : numSplitters; // 'left' is the partition index
}
```

## Segmented Prefix Sum

Here, segmented prefix sum is used to count the occurrence of each split value in each block and it will be used to get the exact thread id for each element. Here, an efficient prefix sum algorithm is performed within segments which has  $O(\log(n))$  time complexity.

```
__global__ void segmentedPrefixSum(int *input, int *output, int n, int blockSize) {
    extern __shared__ int shared[];

    int tid = threadIdx.x;
    int global_tid = blockIdx.x * blockDim.x + tid;

    // Load input into shared memory
    if (global_tid < n) {
        shared[tid] = input[global_tid];
    } else {
        shared[tid] = -1; // Set to a value that won't match any segment
    }
    __syncthreads();

    // Initialize prefix sum within segments
    if (global_tid < n) {
        if (tid == 0 || shared[tid] != shared[tid - 1]) {
            output[global_tid] = 0;
        } else {
            output[global_tid] = output[global_tid - 1] + 1;
        }
    }
}
```

```

__syncthreads();

// Perform prefix sum within segments in shared memory
for (int stride = 1; stride < blockSize; stride *= 2) {
    int temp = 0;
    if (tid >= stride && shared[tid] == shared[tid - stride]) {
        temp = output[global_tid - stride];
    }
    __syncthreads();
    if (global_tid < n && tid >= stride && shared[tid] == shared[tid - stride]) {
        output[global_tid] += temp;
    }
    __syncthreads();
}

```

## Count Split Values

This counts the total number of values of a split in a block. The split values are counted in shared memory and atomic Add operation is used for concurrency control between any 2 threads. Then the shared counts are written into global memory again using atomic Add operation.

```

for (int i = tid; i < num_splits; i += block_size) {
    shared_counts[i] = 0;
}
__syncthreads();

// Count split indices in shared memory
for (int i = tid + bid * block_size; i < num_elements; i += block_size * gridDim.x) {
    int split_idx = split_indices[i];
    atomicAdd(&shared_counts[split_idx], 1);
}
__syncthreads();

// Write shared memory counts to global memory
for (int i = tid; i < num_splits; i += block_size) {
    atomicAdd(&split_counts[i * gridDim.x + bid], shared_counts[i]);
}

```

## Prefix Sum

This prefix sum is done to get the exact global index of each element. So, the prefix sum is done for each element in till a certain block. The prefix sum is efficiently implemented using the Exclusive Sum in Device Scan of the CUB library.

```

// Determine temporary device storage requirements
void *d_temp_storage = nullptr;
size_t temp_storage_bytes = 0;
(cub::DeviceScan::ExclusiveSum(d_temp_storage, temp_storage_bytes,
d_input, d_output, num_items));

// Allocate temporary storage

```

```

    (cudaMalloc(&d_temp_storage, temp_storage_bytes));

    // Perform exclusive prefix sum
    (cub::DeviceScan::ExclusiveSum(d_temp_storage, temp_storage_bytes,
d_input, d_output, num_items));

```

## Global Array Assignment

The block id and thread id for each element is used to assign it to its corresponding thread and block. Now, the primary keys of both tables are globally sorted.

```

int tid = blockIdx.x * blockDim.x + threadIdx.x;

if (tid < size) {
    int ind_in_cnt = blockIdx.x + p * NewBlockId[tid];
    int final_ind =
segment_sum[tid]+d_split_counts_prefixsum[ind_in_cnt];
    if (d_dst[final_ind]!=0) {
        printf("%d final index, %d segment, %d prefix, %d block, %d
tid\n",
final_ind,segment_sum[tid],d_split_counts_prefixsum[ind_in_cnt],NewBlockId[ti
d],tid);
    }
    d_dst[final_ind] = d_src[tid];
}

```

## Join After Sort

The globally sorted blocks are now traversed in parallel, and matching entries are got by only comparing the current element to its previous element. The result is written to the output array. The values corresponding to the matching pairs are got from the hash map of keys to values.

```

int tid = blockIdx.x * blockDim.x + threadIdx.x;

if (tid < numElements - 1)
{
    // Check if current element is equal to the next element
    if (d_final_array[tid] == d_final_array[tid + 1])
    {
        int k=d_final_array[tid];
        results[3*tid] = k;
        results[3*tid+1] = hmap1[k];
        results[3*tid+2] = hmap2[k];
    }
}

```

## Experimental Setup and Datasets Used

The experimental setup involved comparing the performance of our algorithm with the Crystal library on identical datasets. The datasets consisted of two tables with 1 million rows each, containing randomly generated integer keys and values. The details of the GPU used for the experiments are:

### Hardware:

- **LGPU1:**
  - CPU: 2 x Intel Xeon CPU E5-2650 v4
  - Memory: 256GB RAM
  - GPU:
    - 2 x Nvidia RTX2080 8GB GDDR6
    - 8 x Nvidia Titan Xp 12GB GDDR5X
  - Storage: /home/data/ or project directories
    - \*(local /home/data is not backed up, please backup your critical data regularly.)
- **LGPU2:**
  - CPU: Intel Xeon Gold 6326
  - Memory: 256GB RAM
  - GPU: 7 x Nvidia RTX3090 24GB GDDR6, 1 x Nvidia RTX4090 24GB GDDR6
  - Storage: /ssddata/, /ssddata1 or project directories
    - \*(local /ssddata and /ssddata1 are not backed up, please backup your critical data regularly.)

### Software:

- Operating System: Almalinux 9
- Application Software: (in alphabetical order of application)
  1. CUDA Toolkit (/usr/local/cuda-\*)

CUDA version 10.1.1, with the CUB and Thrust libraries are used for optimized parallel operations. The primary goal was to measure execution time and resource utilization, comparing the results against the Crystal library.

## Experimental Results and Discussion

Our implementation demonstrated a significantly better performance than the CPU Pk-Pk join operation. The BSP sorting and prefix sum optimizations contributed to faster execution times and better memory utilization. The following table summarizes the performance metrics:

```
Data set 1 size : 10^6
Data set 2 size : 10^6
Key mod value : 10^8
NVPROF statistics while running the code on randomly generated data set
==3785409== NVPROF is profiling process 3785409, command: PKPKJOIN
==3785409== Profiling application: PKPKJOIN
==3785409== Profiling result:
Type Time(%) Time Calls Avg Min Max Name
GPU activities: 97.44% 534.81ms 3 178.27ms 5.2992ms 264.77ms [CUDA memcpy HtoD]
0.89% 4.8801ms 1 4.8801ms 4.8801ms countSplits(int*, int*, int, int)
0.73% 4.0189ms 1 4.0189ms 4.0189ms 4.0189ms BlockSortKernel2(int*, int*, int*, int, int)
```

```

0.23% 1.2761ms 1 1.2761ms 1.2761ms initCurand(curandStateXORWOW*,  

unsigned long, int)  

0.22% 1.2142ms 1 1.2142ms 1.2142ms void  

cub::DeviceScanKernel<cub::DeviceScanPolicy<int>::Policy600, int*, int*, cub::ScanTileState<int, bool=1>,  

cub::Sum, cub::detail::InputValue<int, int*>, int, int>(cub::DeviceScanPolicy<int>::Policy600, int*, int*, int,  

int, bool=1, cub::ScanTileState<int, bool=1>)  

0.14% 792.56us 1 792.56us 792.56us Assign(int*, int*, int*, int*,  

int*, int, int)  

0.13% 733.74us 1 733.74us 733.74us 733.74us BlockSortKernel(int*, int*, int)  

0.12% 652.90us 3 217.63us 20.833us 575.26us [CUDA memset]  

0.02% 136.07us 1 136.07us 136.07us segmentedPrefixSum(int*, int*,  

int, int)  

0.02% 110.85us 1 110.85us 110.85us 110.85us findSplitsKernel(int const *,  

int*, int const *, int, int)  

0.01% 50.403us 1 50.403us 50.403us 50.403us sampleElements(curandStateXORWOW*,  

int*, int*, int, int)  

0.01% 42.819us 1 42.819us 42.819us 42.819us JoinKernel(int*, int*, int, int*,  

int*)  

0.00% 24.290us 5 4.8580us 4.7360us 5.1210us void  

cub::RadixSortScanBinsKernel<cub::DeviceRadixSortPolicy<int, cub::NullType, unsigned int>::Policy800, unsigned  

int>(cub::NullType*, int)

0.00% 21.538us 2 10.769us 9.0250us 12.513us void

cub::DeviceScanInitKernel<cub::ScanTileState<int, bool=1>::(int, int)
API calls: 75.57% 2.55572s 28 91.276ms 3.9860us 2.55544s cudaLaunchKernel
15.89% 537.32ms 3 179.11ms 4.4682ms 268.42ms cudaMemcpy
7.84% 265.19ms 17 15.599ms 6.0940us 262.82ms cudaMalloc
0.36% 12.085ms 12 1.0071ms 9.4470us 6.4133ms cudaFree
0.29% 9.8906ms 1140 8.6750us 159ns 2.0199ms cuDeviceGetAttribute
0.05% 1.6618ms 1 1.6618ms 1.6618ms 1.6618ms cudaDeviceSynchronize
0.00% 69.403us 3 23.134us 12.629us 42.077us cudaMemcpy
0.00% 52.933us 10 5.2930us 3.5430us 13.267us cuDeviceGetName
0.00% 29.748us 13 2.2880us 509ns 13.739us
cudaOccupancyMaxActiveBlocksPerMultiprocessorWithFlags
0.00% 25.133us 10 2.5130us 1.3340us 7.7580us cuDeviceGetPCIBusId
0.00% 15.853us 1 15.853us 15.853us 15.853us cudaFuncGetAttributes
0.00% 15.261us 98 155ns 130ns 336ns cudaGetLastError
0.00% 6.0110us 9 667ns 322ns 2.4390us cudaGetDevice
0.00% 5.6820us 34 167ns 140ns 303ns cudaPeekAtLastError
0.00% 4.9500us 20 247ns 148ns 697ns cuDeviceGet
0.00% 3.6010us 3 1.2000us 495ns 2.1630us cudaDeviceGetAttribute
0.00% 3.2800us 10 328ns 277ns 557ns cuDeviceTotalMem
0.00% 2.2290us 10 222ns 179ns 338ns cuDeviceGetUuid
0.00% 1.4420us 3 480ns 247ns 936ns cuDeviceGetCount
0.00% 366ns 1 366ns 366ns 366ns cuModuleGetLoadingMode

0.00% 319ns 1 319ns 319ns 319ns cudaGetDeviceCount

```

## Performance Analysis

The speedup achieved by our implementation can be attributed to several factors:

- Less Memory Latency:** Only 2 reads and 1 write is done into the Global memory using the PCI bus which is the bottleneck for the algorithm.
- Efficient Sorting:** The BSP sorting algorithm optimizes data sorting within GPU shared memory, reducing the overhead associated with global memory access.
- Parallel Prefix Sum:** By utilizing the CUB library's efficient parallel prefix sum implementation; we minimized the time required for intermediate result aggregation.
- Parallel Distribution into Blocks:** The segmented prefix sums and prefix sum values are computed efficiently in parallel which are then used for the efficient distribution of the elements into the Global sorted array. This saves a lot of time than atomic operations.

## Comparison with Crystal Library

While the Crystal library provides a robust set of GPU primitives for SQL query execution, our implementation focuses on optimizing the specific case of Pk-Pk joins. By tailoring the algorithm to this particular operation, we achieved a good performance metric but still slower than Crystal library because of their high level optimization in every intricacy of the CUDA code and advanced algorithms and techniques due to which it has achieved around 10x faster execution time over the same dataset. It is important to note that the Crystal library also offers broader functionality and flexibility for various database operations. Here is the output statistics of the Crystal library on the same dataset.

### Crystal library statistics on same data set

```
==3786656== NVPROF is profiling process 3786656, command: anyfile
{"time_memset":0.031168,"time_build":296.066,"time_probe":0.211168}
{"num_dim":1000000,"num_fact":1000000,"radix":0,"time_partition_build":0,"time_partition_probe":0,"time_partition_total":0,"time_build":296.066,"time_probe":0.211168,"time_extra":0.031168,"time_join_total":296.309}
==3786656== Profiling application: anyfile
==3786656== Profiling result:
Type Time(%) Time Calls Avg Min Max Name
GPU activities: 91.26% 10.717ms 4 2.6793ms 2.6776ms 2.6802ms [CUDA memcpy HtoD]
6.90% 810.25us 1 810.25us 810.25us 810.25us void build_kernel<int=128,
int=4>(int*, int*, int, int*, int)
1.68% 197.80us 1 197.80us 197.80us 197.80us void probe_kernel<int=128,
int=4>(int*, int*, int, int*, int, __int64*)
0.16% 18.593us 2 9.2960us 1.4400us 17.153us [CUDA memset]
API calls: 50.83% 295.28ms 2 147.64ms 11.869us 295.27ms cudaLaunchKernel
45.27% 263.00ms 6 43.834ms 141.28us 261.81ms cudaMalloc
1.74% 10.124ms 4 2.5310ms 1.9767ms 2.8489ms cudaMemcpy
1.68% 9.7874ms 1140 8.5850us 137ns 2.0829ms cuDeviceGetAttribute
0.23% 1.3461ms 4 336.52us 6.5660us 801.21us cudaEventSynchronize
0.18% 1.0747ms 5 214.94us 179.62us 265.58us cudaFree
0.01% 77.355us 2 38.677us 5.1440us 72.211us cudaMemset
0.01% 54.811us 10 5.4810us 3.6430us 13.683us cuDeviceGetName
0.01% 35.246us 9 3.9160us 2.1610us 8.3800us cudaEventRecord
0.00% 26.509us 10 2.6500us 1.2820us 7.0470us cuDeviceGetPCIBusId
0.00% 25.217us 2 12.608us 1.3310us 23.886us cudaEventCreate
0.00% 24.596us 6 4.0990us 1.5690us 11.039us cudaEventCreateWithFlags
0.00% 23.406us 12 1.9500us 489ns 9.1270us cudaGetDevice
0.00% 16.366us 51 320ns 131ns 6.1930us cudaGetLastError
0.00% 6.7070us 4 1.6760us 978ns 2.4250us cudaEventElapsedTime
0.00% 6.5960us 5 1.3190us 958ns 1.8130us cudaEventDestroy
0.00% 3.8810us 20 194ns 143ns 673ns cuDeviceGet
0.00% 3.2820us 10 328ns 269ns 577ns cuDeviceTotalMem
0.00% 2.7340us 3 911ns 231ns 1.8670us cuDeviceGetCount
0.00% 2.1820us 10 218ns 189ns 296ns cuDeviceGetUuid

0.00% 355ns 1 355ns 355ns 355ns cuModuleGetLoadingMode
```

## Conclusion

The parallel Pk-Pk join algorithm implemented using CUDA-C++ demonstrates significant performance benefits over existing CPU-based solutions. By leveraging advanced CUDA techniques and libraries, our implementation efficiently handles large datasets, providing a robust solution for database join operations on GPUs. Future work will explore further optimizations and extend the approach to more complex query operations.

## Future Work

To build on the current implementation, several avenues for future research and development are proposed:

1. **Support for Multi-way Joins:** Extending the algorithm to handle multi-way joins between multiple tables will increase its applicability to more complex queries.
2. **Similarity Joins:** Developing a parallel algorithm for similarity join.
3. **Large Datasets:** Join between two tables having very large number of entries i.e. in 100 millions needs much more optimization in the code given that ours can handle only up to 1 million entries per table.
4. **Dynamic Load Balancing:** Implementing dynamic load balancing mechanisms to handle varying data distributions and improve overall performance.
5. **Integration with Database Systems:** Integrating the algorithm into existing database management systems to provide seamless GPU-accelerated query execution.