



Katholieke
Universiteit
Leuven

Department of
Computer Science

ADVANCED COMPUTER ARCHITECTURE

project

Thibaut Beck
Wannes Paesschesoone
Academic year 2025–2026

Contents

1	Introduction	2
2	Theory	2
2.1	Point Cloud	2
2.2	Voxel	2
2.3	Pointcloud to Voxel Grid Filtering	2
2.4	Voxel Size	2
2.5	Morton Codes	2
2.6	CUDA Thrust Library	2
2.7	LAS Files	2
3	Implementation	3
3.1	Shared Morton Encoding Utilities	3
3.2	Morton and Sort Voxelizer	3
3.3	Dynamic Hash Map Voxelizer	5
4	Results and Analysis	8
4.1	Timing remarks and methodology	8
4.2	General Observations	8
4.3	Performance Analysis per Method	9
4.3.1	Morton-Code Voxelizer	9
4.3.2	Hash-Table Voxelizer	9
4.4	Impact of Voxel Size	9
4.5	Impact of Block Size	9
4.6	Timing Breakdown	10
4.6.1	Morton Method	10
4.6.2	Hash Method	10
4.7	Conclusions	10
4.7.1	Optimal Configurations	10
4.7.2	Explanation	10
4.7.3	Performance Limitations	10
4.8	gpu vs cpu performance	11
4.9	Figures of the overall gpu times	11
5	Visualisation	13
6	Possible Improvements	13

1 Introduction

This report talks about the implementation of a point cloud to voxel grid filter. The reason for implementing this filter is to reduce the number of points in a point cloud, while preserving the overall structure and appearance of the original data. Doing this on a CPU takes a long time for large point clouds, so the filter is implemented on a GPU using CUDA to take advantage of the parallel processing capabilities of modern graphics hardware.

2 Theory

2.1 Point Cloud

A **point cloud** is a collection of data points defined in a three-dimensional coordinate system. Each point represents a position in space, typically described by (x, y, z) coordinates. Point clouds are commonly acquired using 3D scanners, LiDAR sensors, or photogrammetry systems. They are used in applications such as 3D modeling, robotics, mapping, and computer vision.

2.2 Voxel

A **voxel** (volumetric pixel) is the smallest unit of a 3D grid, similar to how a pixel is the smallest unit of a 2D image. Voxels divide 3D space into uniform cubes, allowing volumetric representations of shapes or environments. They are used in 3D reconstruction, simulation, gaming, and medical imaging.

2.3 Pointcloud to Voxel Grid Filtering

Pointcloud to voxel grid filtering converts an unstructured point cloud into a regular 3D grid. Space is divided into equal-sized voxels, and each point is assigned to its corresponding voxel. This reduces data density, removes redundant points, and creates a structured representation that is easier and faster to process for tasks such as downsampling and spatial queries.

2.4 Voxel Size

The **voxel size** is a critical parameter in pointcloud to voxel grid filtering. It determines the dimensions of each cubic cell in the 3D grid. A larger voxel size groups more points together, resulting in a coarser, more downsampled representation of the original point cloud. Conversely, a smaller voxel size produces a finer representation with higher spatial resolution but increased memory and computational requirements. The choice of voxel size depends on the application and the desired balance between accuracy and efficiency.

2.5 Morton Codes

Morton codes (also called Z-order curves) are a method of encoding multi-dimensional coordinates into a single integer value. They work by bit-interleaving the binary coordinates (e.g., x, y, z). This encoding preserves spatial locality, making it useful for data structures such as octrees and for accelerating spatial queries on GPUs.

2.6 CUDA Thrust Library

The **CUDA Thrust library** is a parallel algorithms library for NVIDIA GPUs. It provides high-level abstractions similar to the C++ Standard Template Library (STL), including parallel sorting, scanning, reduction, and vector operations. Thrust simplifies GPU programming by offering ready-to-use, highly optimized parallel primitives.

2.7 LAS Files

LAS files are a standardized file format used for storing LiDAR point cloud data. The format supports 3D coordinates, intensity values, classification labels, GPS time, color information, and other metadata. LAS is widely used in geospatial applications, surveying, and remote sensing due to its efficiency and interoperability.

3 Implementation

In this chapter, the practical implementation details of the point cloud to voxel grid filter are presented. To leverage the massive parallelism of modern hardware, the algorithms were developed using CUDA for GPU acceleration. Two distinct parallel strategies were designed to solve the voxelization problem: a sorting-based approach using Morton encoding, and a scattering approach using a GPU-resident dynamic hash map. This hashmap uses the same Morton encoding for voxel indexing, ensuring consistency between both methods. Both voxelizers share the same Morton encoding utilities so that a 3D voxel index is always mapped to the same 64-bit key, independent of the backend. On top of that, the Morton-and-sort voxelizer introduces its own per-voxel accumulation structure and reduction functor, while the hash-based implementation uses a dedicated hash bucket structure.

3.1 Shared Morton Encoding Utilities

Bit expansion helper. The first shared building block is the mapping from 3D integer voxel coordinates (i_x, i_y, i_z) to a single 64-bit Morton code. This is reused in *both* voxelizers: in the sorting-based pipeline it is used to generate keys for `thrust::sort_by_key`, while in the hash-based pipeline the same 64-bit Morton code serves directly as the hash key for open addressing.

The helper function `splitBy3` expands a 21-bit integer into a 64-bit value in which each original bit is separated by two zero bits:

```
1 __device__ __host__ inline uint64_t splitBy3(uint64_t v) {
2     v = (v | (v << 32)) & 0x1f00000000ffffULL;
3     v = (v | (v << 16)) & 0x1f0000ff0000ffULL;
4     v = (v | (v << 8)) & 0x100f00f00f00f0fULL;
5     v = (v | (v << 4)) & 0x10c30c30c30c30c3ULL;
6     v = (v | (v << 2)) & 0x1249249249249249ULL;
7     return v;
8 }
```

Listing 1: Bit expansion of a 21-bit coordinate into Morton layout (shared).

Morton code generation. Using `splitBy3`, the actual Morton code is formed by interleaving the expanded bits of the three coordinates:

```
1 __device__ __host__ inline uint64_t mortonEncode(uint32_t x,
2                                                    uint32_t y,
3                                                    uint32_t z) {
4     return (splitBy3((uint64_t)x) << 2) |
5            (splitBy3((uint64_t)y) << 1) |
6            splitBy3((uint64_t)z);
7 }
```

Listing 2: Morton encoding of 3D voxel indices (shared).

Both functions are marked `__device__ __host__ inline`, so they can be called from host-side test code as well as from all CUDA kernels in both voxelization backends.

3.2 Morton and Sort Voxelizer

To implement the Morton-and-sort voxelizer efficiently on the GPU, the pipeline is decomposed into: (i) Morton code generation via the shared utilities, (ii) per-voxel accumulation of point attributes, and (iii) CUDA kernels that bridge raw device arrays with Thrust’s sort-and-reduce primitives. The accumulation structures and reduction functor described below are specific to this implementation and are *not* used by the hash-based voxelizer.

Per-voxel accumulation structure. To compute voxel centroids and average colors in a single pass after sorting, we define a compact accumulation structure `PointAccum`. It stores the running sums of coordinates and RGB components, as well as the number of points that contributed to the voxel:

```
1 struct PointAccum {
2     float    sumX, sumY, sumZ;
3     uint32_t sumR, sumG, sumB;
4     uint32_t count;
```

```

5
6  __device__ __host__ PointAccum()
7      : sumX(0), sumY(0), sumZ(0),
8        sumR(0), sumG(0), sumB(0),
9        count(0) {}
10
11  __device__ __host__ PointAccum(float x, float y, float z,
12                                uint8_t r, uint8_t g, uint8_t b)
13      : sumX(x), sumY(y), sumZ(z),
14        sumR(r), sumG(g), sumB(b),
15        count(1) {}
16 };

```

Listing 3: Accumulation structure for voxel-wise reduction (Morton-and-sort only).

This structure is used as the value type in the subsequent reduction-by-key stage and is therefore only required by the sorting-based voxelizer.

Binary reduction functor. To use `PointAccum` with `thrust::reduce_by_key`, we define a binary functor `PointAccumOp`. It describes how two partial voxel accumulators are merged:

```

1 struct PointAccumOp {
2     __device__ __host__
3     PointAccum operator()(const PointAccum& a,
4                           const PointAccum& b) const {
5         PointAccum result;
6         result.sumX = a.sumX + b.sumX;
7         result.sumY = a.sumY + b.sumY;
8         result.sumZ = a.sumZ + b.sumZ;
9         result.sumR = a.sumR + b.sumR;
10        result.sumG = a.sumG + b.sumG;
11        result.sumB = a.sumB + b.sumB;
12        result.count = a.count + b.count;
13        return result;
14    }
15 };

```

Listing 4: Binary operator for reducing `PointAccum` values (Morton-and-sort only).

This functor is specific to the reduction-based pipeline and is not needed by the dynamic hash map implementation.

Morton code computation kernel. The first CUDA kernel in this voxelizer takes the input point positions and converts them into Morton codes using the shared `mortonEncode` helper. Each thread processes exactly one point:

```

1 __global__ void computeMortonCodesKernel(
2     const float* x,
3     const float* y,
4     const float* z,
5     uint64_t* mortonCodes,
6     float minX, float minY, float minZ,
7     float invVoxelSize,
8     size_t numPoints)
9 {
10     int idx = blockIdx.x * blockDim.x + threadIdx.x;
11     if (idx >= numPoints) return;
12
13     // Compute voxel indices (offset by min to ensure non-negative indices)
14     uint32_t ix = (uint32_t)floorf((x[idx] - minX) * invVoxelSize);
15     uint32_t iy = (uint32_t)floorf((y[idx] - minY) * invVoxelSize);
16     uint32_t iz = (uint32_t)floorf((z[idx] - minZ) * invVoxelSize);
17
18     mortonCodes[idx] = mortonEncode(ix, iy, iz);
19 }

```

Listing 5: Kernel computing Morton codes for each point (Morton-and-sort).

The kernel assumes a structure-of-arrays (SoA) layout (`x`, `y`, `z` in separate buffers) for coalesced memory access. The inverse voxel size is precomputed on the host to avoid divisions in device code.

Accumulator creation kernel. After sorting the points by their Morton codes, a `PointAccum` object is created for each point in sorted order. Rather than physically reordering all input arrays, a separate index array encodes the permutation induced by the sort:

```

1 __global__ void createPointAccumKernel(
2     const float* x,
3     const float* y,
4     const float* z,
5     const uint8_t* r,
6     const uint8_t* g,
7     const uint8_t* b,
8     const uint32_t* indices,
9     PointAccum* accums,
10    size_t numPoints)
11 {
12     int idx = blockIdx.x * blockDim.x + threadIdx.x;
13     if (idx >= numPoints) return;
14
15     uint32_t origIdx = indices[idx];
16     accums[idx] = PointAccum(
17         x[origIdx], y[origIdx], z[origIdx],
18         r[origIdx], g[origIdx], b[origIdx]
19     );
20 }

```

Listing 6: Kernel creating `PointAccum` objects after sorting (Morton-and-sort).

Integration in the outer loop. The outer host function for the Morton-and-sort voxelizer (not shown in full) orchestrates the following steps:

1. **Preprocessing:** Compute the global bounding box and `invVoxelSize` on the CPU.
2. **Device setup:** Upload positions and colors to device vectors; initialize an index array with `thrust::sequence`.
3. **Morton computation:** Launch `computeMortonCodesKernel` to fill the Morton code buffer using the shared `mortonEncode`.
4. **Sorting:** Call `thrust::sort_by_key` on the pair `(d.mortonCodes, d.indices)`.
5. **Accumulator creation:** Launch `createPointAccumKernel` to build a `PointAccum` per sorted point.
6. **Reduction by voxel:** Invoke `thrust::reduce_by_key` with `PointAccum` values and `PointAccumOp` to obtain one accumulator per occupied voxel.
7. **Final averaging:** Copy the reduced accumulators back to the host and divide sums by `count` to obtain centroid and color per voxel.

3.3 Dynamic Hash Map Voxelizer

The second approach implements a custom open-addressing hash table directly in GPU global memory. This method is designed for maximum throughput, avoiding the global synchronization and $O(N \log N)$ complexity of the sorting-based pipeline. Instead of ordering points, it scatters them into a large hash table and accumulates voxel statistics in-place using atomics.

This implementation *reuses* only the shared Morton encoding helpers (`splitBy3` and `mortonEncode`). The per-voxel accumulation is performed directly inside each hash bucket, so there is no need for `PointAccum` or `PointAccumOp`.

Hash bucket layout. Each entry in the GPU hash table is represented by a tightly packed `HashBucket` structure:

```

1 struct __align__(16) HashBucket {
2     unsigned long long key; // 64-bit Morton code
3     float sumX, sumY, sumZ;
4     uint32_t sumR, sumG, sumB;
5     uint32_t count;
6 };

```

```

7
8 #define EMPTY_KEY 0xFFFFFFFFFFFFFFFFULL

```

Listing 7: Aligned hash bucket structure and empty key sentinel (hash voxelizer).

The `__align__(16)` qualifier improves memory coalescing, and `EMPTY_KEY` serves as a sentinel to identify unused slots. The 64-bit key field stores the Morton code computed by the shared `mortonEncode` function.

Hash table initialization. Before insertion, the hash table is cleared by setting all keys to `EMPTY_KEY` and zeroing the accumulators:

```

1 __global__ void initHashMapKernel(HashBucket* table, size_t capacity) {
2     size_t idx = blockIdx.x * blockDim.x + threadIdx.x;
3     if (idx < capacity) {
4         table[idx].key = EMPTY_KEY;
5         table[idx].count = 0;
6
7         table[idx].sumX = 0.0f;
8         table[idx].sumY = 0.0f;
9         table[idx].sumZ = 0.0f;
10        table[idx].sumR = 0;
11        table[idx].sumG = 0;
12        table[idx].sumB = 0;
13    }
14 }

```

Listing 8: Kernel initializing the GPU hash table (hash voxelizer).

Hash-Based Voxel Accumulation using Atomic Open Addressing The core function of the hash-based voxelizer is the `populateHashMapKernel` CUDA kernel, which implements a scatter-and-accumulate pattern using a 64-bit Morton code derived from 3D integer coordinates (i_x, i_y, i_z) via `mortonEncode` as the key. Each thread calculates an initial hash slot $H = \text{mortonCode} \pmod{\text{capacity}}$ and uses linear probing to find the correct `HashBucket`. Concurrency is managed by atomically claiming an empty slot using `atomicCAS` (Compare-and-Swap) with `EMPTY_KEY`, or by finding a slot already containing its key. Once the correct slot is secured, the thread uses atomic additions (`atomicAdd`) to safely accumulate the point’s properties (position sums, color sums, and count) into the shared `HashBucket` structure, ensuring thread-safe data aggregation despite concurrent writes.

```

1 __global__ void populateHashMapKernel(
2     const float* x, const float* y, const float* z,
3     const uint8_t* r, const uint8_t* g, const uint8_t* b,
4     HashBucket* table,
5     size_t capacity,
6     size_t numPoints,
7     float minX, float minY, float minZ,
8     float invVoxelSize)
9 {
10     size_t idx = blockIdx.x * blockDim.x + threadIdx.x;
11     if (idx >= numPoints) return;
12
13     // 1. Quantize position to voxel grid and compute Morton key
14     uint32_t ix = (uint32_t)floorf((x[idx] - minX) * invVoxelSize);
15     uint32_t iy = (uint32_t)floorf((y[idx] - minY) * invVoxelSize);
16     uint32_t iz = (uint32_t)floorf((z[idx] - minZ) * invVoxelSize);
17
18     uint64_t mortonCode = mortonEncode(ix, iy, iz);
19
20     // 2. Initial hash slot
21     size_t hashIdx = mortonCode % capacity;
22
23     // 3. Linear probing with atomic CAS and atomic adds
24     for (size_t i = 0; i < capacity; ++i) {
25         size_t currentSlot = (hashIdx + i) % capacity;
26
27         unsigned long long oldKey = table[currentSlot].key;
28
29         // Try to claim an empty slot
30         if (oldKey == EMPTY_KEY) {
31             unsigned long long assumed =

```

```

32         atomicCAS((unsigned long long*)&table[currentSlot].key,
33                   EMPTY_KEY,
34                   (unsigned long long)mortonCode);
35     if (assumed == EMPTY_KEY) {
36         oldKey = mortonCode; // we now own this bucket
37     } else {
38         oldKey = assumed;    // another thread claimed it
39     }
40 }
41
42 // If this bucket belongs to our Morton key, accumulate data
43 if (oldKey == mortonCode) {
44     atomicAdd(&table[currentSlot].sumX, x[idx]);
45     atomicAdd(&table[currentSlot].sumY, y[idx]);
46     atomicAdd(&table[currentSlot].sumZ, z[idx]);
47
48     atomicAdd(&table[currentSlot].sumR, (uint32_t)r[idx]);
49     atomicAdd(&table[currentSlot].sumG, (uint32_t)g[idx]);
50     atomicAdd(&table[currentSlot].sumB, (uint32_t)b[idx]);
51
52     atomicAdd(&table[currentSlot].count, 1);
53     return;
54 }
55
56 // Otherwise: collision with a different key, continue probing
57 }
58 }

```

Listing 9: Kernel inserting points into the GPU hash map (hash voxelizer).

Note that this kernel does *not* use `PointAccum` or `PointAccumOp`; all aggregation happens directly inside the `HashBucket`.

Counting valid voxels. Because open addressing leaves gaps in the hash table, a compaction pass is required. The first step is to count how many buckets are actually occupied:

```

1  __global__ void countValidBucketsKernel(
2      HashBucket* table,
3      size_t capacity,
4      uint32_t* counter)
5  {
6      size_t idx = blockIdx.x * blockDim.x + threadIdx.x;
7      if (idx < capacity) {
8          if (table[idx].key != EMPTY_KEY) {
9              atomicAdd(counter, 1);
10         }
11     }
12 }

```

Listing 10: Kernel counting the number of occupied buckets (hash voxelizer).

Collecting results. A second pass converts the populated buckets into a dense array of voxel representatives. As in the sorting-based implementation, each voxel is represented by its centroid and average color:

```

1  __global__ void collectResultsKernel(
2      HashBucket* table,
3      size_t capacity,
4      Point* output,
5      uint32_t* globalCounter)
6  {
7      size_t idx = blockIdx.x * blockDim.x + threadIdx.x;
8      if (idx >= capacity) return;
9
10     HashBucket bucket = table[idx];
11
12     if (bucket.key != EMPTY_KEY && bucket.count > 0) {
13         uint32_t outIdx = atomicAdd(globalCounter, 1);
14
15         float c = (float)bucket.count;
16

```



```

17     Point p;
18     p.x = bucket.sumX / c;
19     p.y = bucket.sumY / c;
20     p.z = bucket.sumZ / c;
21     p.r = (uint8_t)(bucket.sumR / bucket.count);
22     p.g = (uint8_t)(bucket.sumG / bucket.count);
23     p.b = (uint8_t)(bucket.sumB / bucket.count);
24
25     output[outIdx] = p;
26 }
27 }

```

Listing 11: Kernel converting hash buckets into final voxel points (hash voxelizer).

Integration in the outer loop. The host-side driver for the dynamic hash map voxelizer can then be summarized as:

1. **Capacity selection:** Choose a hash table capacity as a multiple of the number of points (e.g. factor 2-4).
2. **Preprocessing:** Compute the bounding box and `inverse VoxelSize` on the CPU.
3. **Device setup:** Allocate device arrays for point attributes and copy input data.
4. **Hash table initialization:** Allocate and clear `HashBucket` array via `initHashMapKernel`.
5. **Scatter-and-accumulate:** Launch `populateHashMapKernel`, which uses the shared `mortonEncode` to derive keys and accumulates directly into buckets.
6. **Voxel counting:** Use `countValidBucketsKernel` to determine the number of occupied buckets and allocate a dense output array.
7. **Compaction:** Reset the counter and call `collectResultsKernel` to write voxel centroids and colors into the dense output array.
8. **Host transfer:** Copy the compact output back to the CPU and release GPU memory.

In this way, the two voxelizers share a single, consistent Morton encoding implementation, while each maintains its own accumulation and reduction strategy tailored to its parallelization scheme.

4 Results and Analysis

4.1 Timing remarks and methodology

All the following results were obtained on relatively old hardware: a NVIDIA Geforce GTX 1080ti with 11GB of GDDR5X memory, paired with an Intel Core i7-8700K CPU @ 4.8 GHz and 32GB of DDR4 RAM and the pc runs ubuntu 24.04.3 LTS. The point cloud used for testing contains 648433 points, representing a dense scan of an urban environment.

4.2 General Observations

The voxelizers are compared among five voxel sizes (0.25, 0.5, 0.75, 1.0, 1.25). For each voxel size, eight block sizes were tested (1, 2, 4, 8, 16, 32, 64, 256, 512, 1024 threads per block). The hash-table voxelizer was evaluated with three different capacity factors (2, 3, and 4 times the number of input points). The performance metrics recorded include total execution time and a breakdown of time spent in key phases of each algorithm. Every configuration was run 100 times to obtain an average execution time, minimizing the impact of transient system load variations.

4.3 Performance Analysis per Method

4.3.1 Morton-Code Voxelizer

The Morton-based approach demonstrates consistent and predictable performance. Key findings include:

- **Optimal block size:** 4–256 threads per block.
- **Total execution time:** Ranging from 25.18 ms (voxel size 1.25, block size 64) to 40.46 ms (voxel size 0.25, block size 1).
- **Primary bottleneck:** The sorting stage dominates runtime, requiring approximately 0.90–1.02 ms, nearly constant across all tests.

Very small block sizes (1–4 threads) significantly slow down Morton-code generation (0.27–1.07 ms), while larger block sizes stabilize this step to around 0.04 ms.

4.3.2 Hash-Table Voxelizer

The hash-based method demonstrates greater variability, largely influenced by the selected capacity factor of the hash table.

Capacity Factor 2 (highest efficiency)

- **Best overall performance:** 25.67 ms (voxel size 1.0, block size 32).
- **Optimal block sizes:** 8–32 threads per block.
- **Advantages:** Lowest memory overhead and fastest device-to-host transfer times.

Capacity Factor 3 (balanced)

- Execution times are **2–4 ms slower** compared to capacity factor 2.
- Reduced collisions during accumulation, at the cost of increased memory usage.

Capacity Factor 4 (highest overhead)

- **5–10 ms slower** than capacity factor 2.
- Initialization and memory operations increase disproportionately.
- Device-to-host + cleanup time rises up to 5.26 ms (compared to 2.71 ms for CF=2).

4.4 Impact of Voxel Size

Smaller voxel sizes generate a significantly larger number of unique voxel entries. This favors the hash-table method with a low capacity factor, as the Morton approach becomes less efficient when spatial data is highly fragmented.

For larger voxels (e.g., voxel size 1.25), the Morton method benefits from predictable spatial coherence and memory access patterns, making it more efficient than the hash-based approach.

4.5 Impact of Block Size

- **Block sizes 1–4:** Strong performance penalties due to insufficient warp utilization.
- **Block sizes 8–256:** Optimal performance range with minimal overhead.
- **Block sizes 512–1024:** No significant improvements; potentially limited by register pressure.

4.6 Timing Breakdown

4.6.1 Morton Method

Dominant phases:

1. **Sorting:** 0.90–1.02 ms (consistent across tests)
2. **Morton code computation:** 0.04–1.07 ms (strongly dependent on block size)
3. **Point accumulation:** 0.24–1.08 ms

4.6.2 Hash Method

Dominant phases:

1. **Device-to-host transfer + cleanup:** 2.71–5.26 ms (increases with capacity factor)
2. **Initialization:** 0.33–3.71 ms (proportional to hash-table size)
3. **Populate phase:** 0.31–2.46 ms (collision-sensitive)

4.7 Conclusions

4.7.1 Optimal Configurations

The analysis indicates that different voxel sizes benefit from different GPU strategies:

- **Small voxel sizes (0.25–0.5):** Hash-based voxelizer with capacity factor 2 and a block size of 16–32 threads.
- **Medium voxel sizes (0.75–1.0):** Hash-based voxelizer with capacity factor 2 and a block size of 32 threads.
- **Large voxel sizes (1.25 and above):** Morton-code voxelizer with a block size between 64 and 256 threads.

4.7.2 Explanation

These optimal configurations follow from the observed behaviour of both voxelization methods:

1. For small voxel sizes, the hash method with capacity factor 2 achieves the highest throughput due to low memory overhead and efficient device-to-host transfers.
2. For larger voxel sizes, the Morton-based approach becomes more efficient thanks to spatial coherence and predictable memory access patterns.
3. Block sizes between 16 and 64 threads provide the most balanced performance in terms of occupancy and register usage.

4.7.3 Performance Limitations

Despite their strengths, both methods exhibit certain limitations:

- The hash-based method is primarily constrained by memory bandwidth, especially during device-to-host transfers.
- The Morton-based method is strongly dominated by the sorting stage, which forms its main computational bottleneck.
- For both approaches, increasing block sizes beyond 256 threads yields minimal benefits due to hardware saturation.

4.8 gpu vs cpu performance

Voxel Size	CPU (ms)	Morton (ms)	Hash CF2 (ms)	Hash CF3 (ms)	Hash CF4 (ms)
0.25	429.09	36.61	27.70	29.57	29.10
0.5	335.89	29.14	26.20	28.30	28.57
0.75	272.39	27.11	25.89	26.57	27.74
1.0	253.67	25.71	25.67	26.99	27.46
1.25	240.00	25.18	25.85	27.24	28.28

Table 1: CPU vs Minimum GPU Voxelization Time - All Methods

Voxel Size	CPU (ms)	Morton (ms)	Hash CF2 (ms)	Hash CF3 (ms)	Hash CF4 (ms)
0.25	429.09	40.46	35.25	37.48	41.82
0.5	335.89	31.18	33.87	35.53	40.00
0.75	272.39	28.62	33.80	36.09	39.12
1.0	253.67	27.63	33.36	35.26	38.53
1.25	240.00	28.64	34.41	36.85	39.06

Table 2: CPU vs Maximum GPU Voxelization Time - All Methods

The above tables demonstrate substantial performance improvements when using GPU-based voxelization compared to CPU implementations. Across all tested voxel sizes, the best-performing GPU configurations achieve speedup factors ranging from $9.5\times$ to $15.5\times$, with the most significant gains observed at smaller voxel sizes where parallelization is most effective. Even in the worst-case GPU scenarios (maximum execution times), the speedup remains impressive at $6.1\times$ to $12.2\times$ faster than CPU execution. The hash-based method with capacity factor 2 consistently delivers optimal performance, achieving minimum execution times between 25.67ms and 27.70ms compared to CPU times of 240.00ms to 429.09ms. Notably, the speedup factor decreases slightly as voxel size increases, suggesting that GPU implementations particularly excel when processing larger numbers of smaller voxels, where the massive parallelism of the GPU architecture can be fully exploited.

4.9 Figures of the overall gpu times

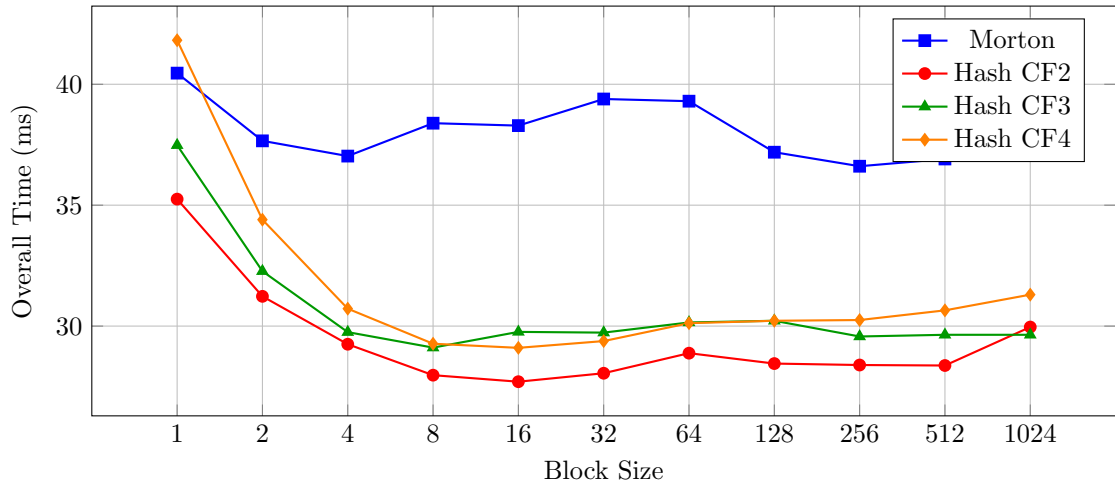


Figure 1: Overall execution time for voxel size 0.25

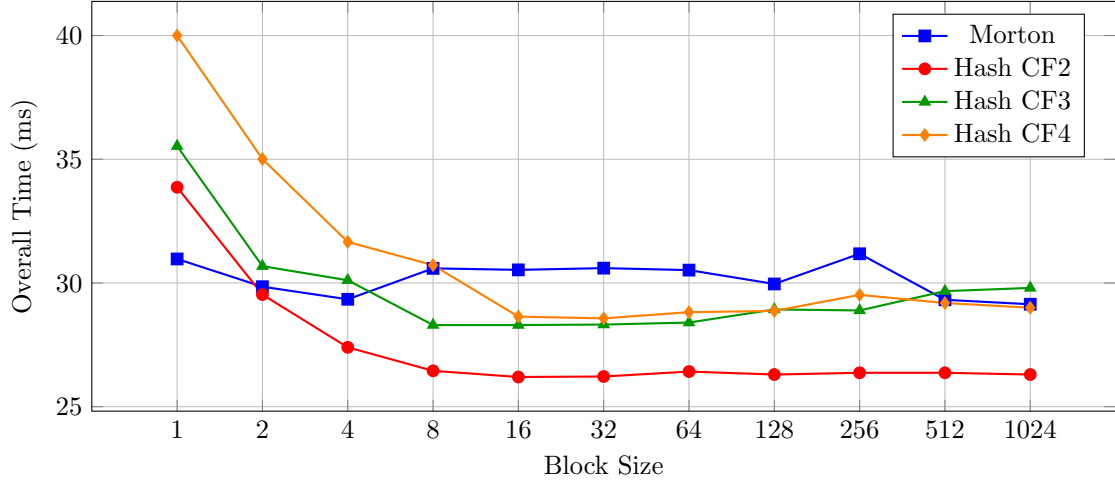


Figure 2: Overall execution time for voxel size 0.5

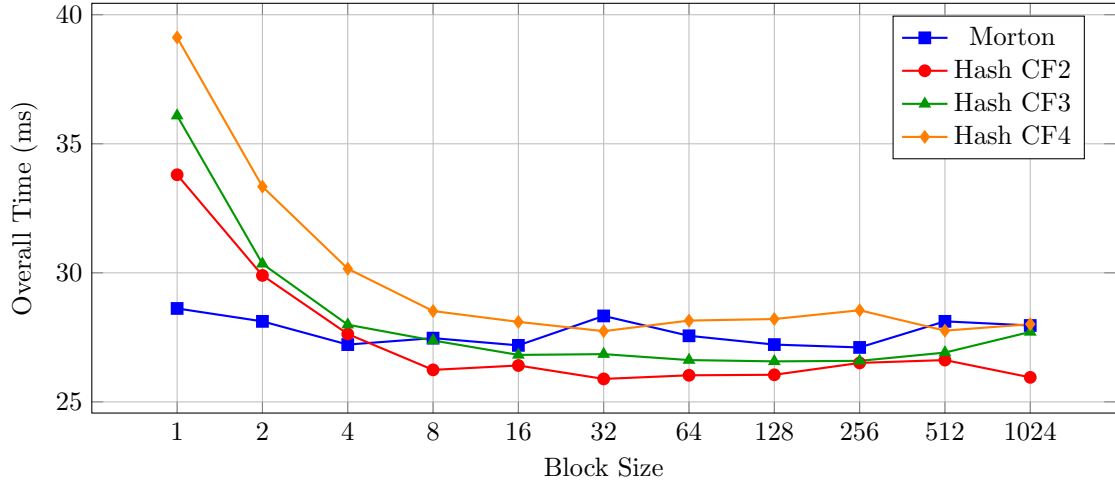


Figure 3: Overall execution time for voxel size 0.75

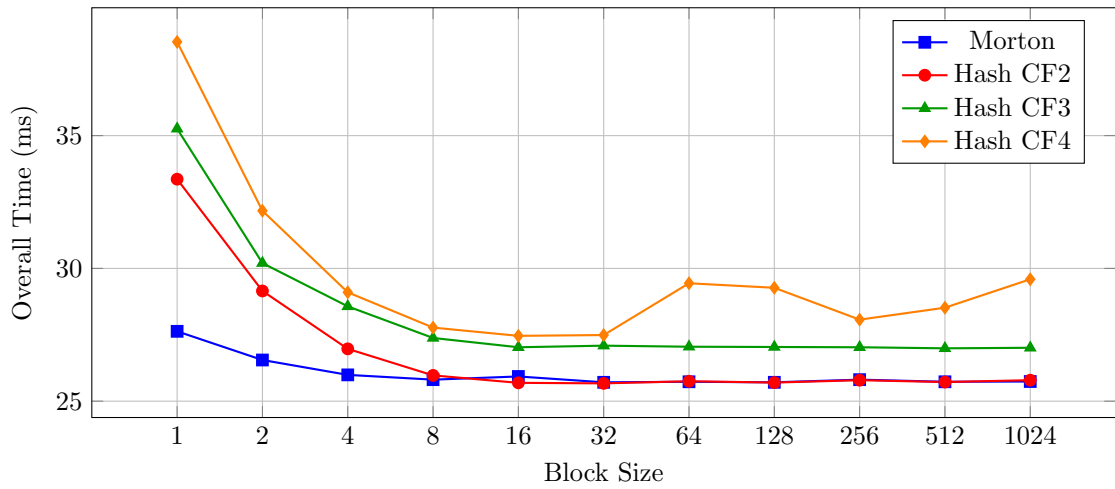


Figure 4: Overall execution time for voxel size 1.0

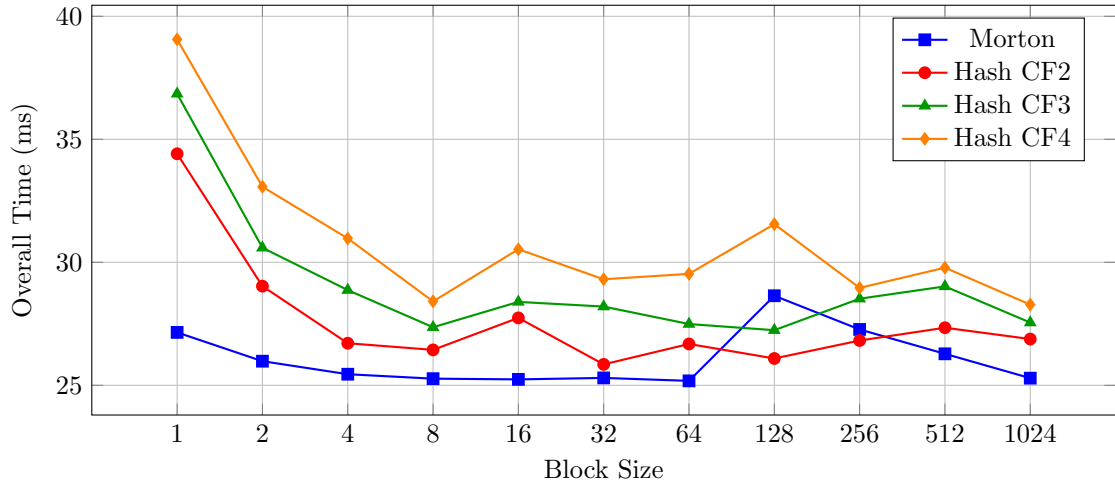


Figure 5: Overall execution time for voxel size 1.25

5 Visualisation

Everything was visualised using the python open3d library. This library allows for easy loading and displaying of point clouds. The original point cloud and the voxelized point cloud were displayed side by side for comparison. **TODOOOO : paste figures**

6 Possible Improvements

The main problem is that both implementations work with a fixed voxel size. This means that if the point cloud is very large but sparse, a lot of memory is wasted on empty space in the hash map implementation. A possible improvement would be to implement an octree structure, which would allow for variable voxel sizes depending on the density of points in a given area.

References

- [1] M. Nießner, M. Zollhöfer, S. Izadi & M. Stamminger, “Real-time 3D Reconstruction at Scale using Voxel Hashing,” ACM Transactions on Graphics (TOG), vol. 32, no. 6, 2013.
- [2] “Morton encoding/decoding through bit interleaving: Implementations”, Forceflow — Jeroen Baert’s Blog, 7 October 2013. Available online: <https://www.forceflow.be/2013/10/07/morton-encodingdecoding-through-bit-interleaving-implementations/> (Accessed: 28 November 2025).
- [3] NVIDIA Corporation, “*Thrust — CUDA Core Compute Libraries (CCCL)*”. Available at: <https://nvidia.github.io/cccl/thrust/> (Accessed: 28 November 2025).
- [4] Wikipedia contributors, “*LAS file format*”, Wikipedia, The Free Encyclopedia. Available at: https://en.wikipedia.org/wiki/LAS_file_format (Accessed: 28 November 2025).
- [5] Ayushi Sharma, “*From Point Clouds to Voxel Grids: A Practical Guide to 3D Data Voxelization*”, Medium, 2021. Available at: <https://medium.com/@ayushi.sharma.3536/from-point-clouds-to-voxel-grids-a-practical-guide-to-3d-data-voxelization-cf5991c1e7bb> (Accessed: 28 November 2025).