

CPU GPU Data Transfer Bandwidth using CUDA

Ali K. Rahimian

November 30, 2025

1 Introduction

This project measures the effective bandwidth between host and device memory using CUDA `cudaMemcpy` transfers. For a range of array sizes, host to device (H2D) and device to host (D2H) copies are timed using CUDA events and the achieved GB/s is reported.

Two host memory types are compared:

- Pageable host buffers allocated with `malloc`.
- Pinned (page locked) host buffers allocated with `cudaHostAlloc`.

Transfers are run for sizes from 1 MB up to 1 GB, and for each size four measurements are obtained: H2D pageable, D2H pageable, H2D pinned, and D2H pinned. The goal is to quantify the impact of host memory type and transfer direction on sustained bandwidth.

2 Methodology

For each transfer size:

1. Allocate a host buffer (pageable or pinned) and a device buffer with `cudaMalloc`.
2. Initialize the host buffer with dummy floating point data.
3. Create CUDA events, perform a warm up copy, then time `num_iters` repeated `cudaMemcpy` calls.
4. Compute effective bandwidth as

$$\text{Bandwidth (GB/s)} = \frac{\text{bytes} \times \text{num_iters}}{\text{time (s)} \times 10^9}.$$

Array sizes are

1, 2, 4, 8, 16, 32, 64, 128, 256, 512, 1024 MB,

and each configuration is repeated with the same number of iterations.

3 Results

3.1 Tables and measurements

Table 1 summarizes the measured bandwidths for all sizes and configurations.

Size (MB)	Pageable		Pinned	
	H2D	D2H	H2D	D2H
1	3.846	2.118	2.506	4.864
2	11.652	6.564	9.892	8.520
4	13.352	7.090	10.056	8.928
8	19.567	12.200	22.272	15.912
16	20.608	12.727	23.510	15.352
32	19.644	13.076	12.188	15.112
64	8.120	9.064	15.961	15.407
128	5.152	13.316	16.913	15.363
256	4.525	13.428	17.027	15.364
512	4.494	13.397	17.120	15.197
1024	13.679	13.396	17.392	15.021

Table 1: Measured H2D and D2H bandwidth (GB/s) for pageable and pinned host memory across different transfer sizes.

3.2 Bandwidth vs array size plot

The CSV output from the executable is postprocessed with a short Python script using `matplotlib` to generate a bandwidth versus size plot with four curves (H2D/D2H and pageable/pinned). A representative plot is shown in Figure 1.

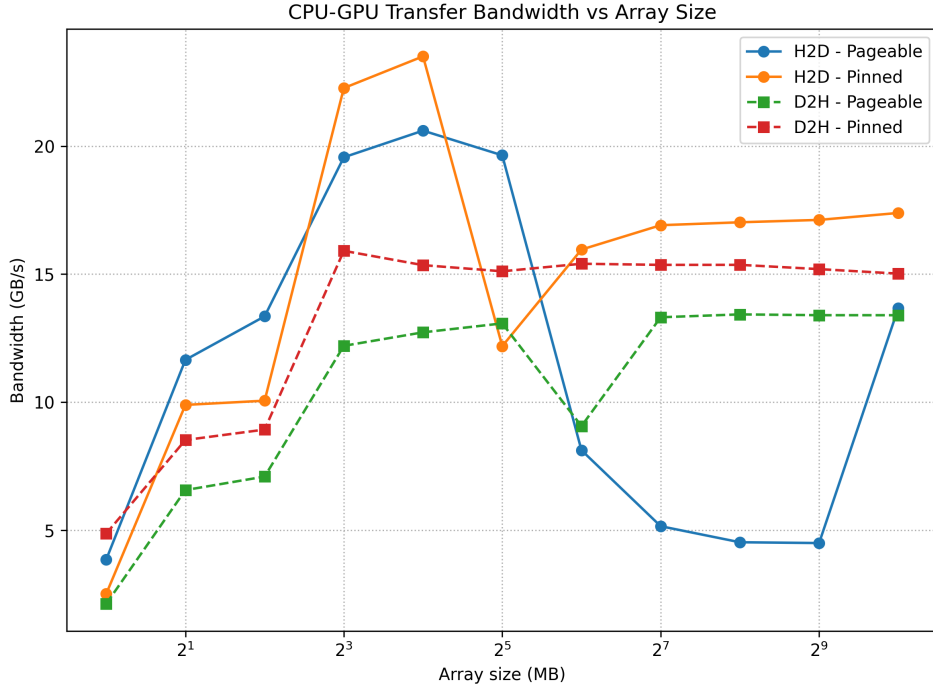


Figure 1: CPU GPU transfer bandwidth vs array size for H2D and D2H, using pageable and pinned host memory.

4 Discussion

Key observations:

- Small transfers (1 to 4 MB) have low and noisy effective bandwidth due to fixed overheads

dominating the timing.

- For medium sizes (8 to 32 MB), bandwidth rises toward a plateau that reflects the interconnect and memory subsystem throughput.
- For large sizes (64 MB and above), pinned memory clearly outperforms pageable memory, especially for H2D.

From the data:

- H2D pinned reaches about 16 to 17 GB/s for large sizes.
- D2H pinned is slightly lower, roughly 15 to 16 GB/s.
- Pageable D2H stabilizes near 13 GB/s, while pageable H2D is lower and more variable.

The advantage of pinned memory is consistent with the elimination of an internal staging copy and direct DMA access to page locked host buffers. The mild asymmetry between H2D and D2H reflects hardware and driver implementation details. For performance critical code, large transfers and pinned memory are preferred; many small transfers should be batched.

5 Source Code

5.1 cuda_utils.cuh

```
1 #ifndef CUDA_UTILS_CUH
2 #define CUDA_UTILS_CUH
3
4 #include <cstdio>
5 #include <cstdlib>
6 #include <cuda_runtime.h>
7
8 #define CHECK_CUDA(call) \
9     do { \
10         cudaError_t err = (call); \
11         if (err != cudaSuccess) { \
12             fprintf(stderr, "CUDA error at %s:%d: %s\n", \
13                 __FILE__, __LINE__, cudaGetErrorString(err)); \
14             std::exit(EXIT_FAILURE); \
15         } \
16     } while (0)
17
18 #endif // CUDA_UTILS_CUH
```

5.2 bandwidth.h

```
1 #ifndef BANDWIDTH_H
2 #define BANDWIDTH_H
3
4 #include <cstdint>
5 #include <vector>
6
7 struct BandwidthResult {
8     double sizeMB; // Array size in MB
9     double h2d_pageable; // GB/s
10    double d2h_pageable; // GB/s
11    double h2d_pinned; // GB/s
12    double d2h_pinned; // GB/s
13 };
14
15 /**
16  * Run bandwidth tests for a list of sizes.
17  *
18  * @param sizes Array sizes in bytes.
19  * @param num_iters Number of repetitions per measurement.
20  * @param results Output vector of BandwidthResult, one per size.
21  */
22 void run_bandwidth_tests(const std::vector<size_t>& sizes,
23                          int num_iters,
24                          std::vector<BandwidthResult>& results);
25
26 #endif // BANDWIDTH_H
```

5.3 bandwidth.cu

```
1 #include "bandwidth.h"
```

```

2 #include "cuda_utils.cuh"
3
4 #include <cstdlib>
5
6 // Measure bandwidth for a single direction and memory pairing
7 static double measure_bandwidth_single_direction(
8     void* dst,
9     const void* src,
10    size_t bytes,
11    cudaMemcpyKind kind,
12    int num_iters)
13 {
14     cudaEvent_t start, stop;
15     CHECK_CUDA(cudaEventCreate(&start));
16     CHECK_CUDA(cudaEventCreate(&stop));
17
18     // Warm up to avoid first use overhead
19     CHECK_CUDA(cudaMemcpy(dst, src, bytes, kind));
20
21     CHECK_CUDA(cudaEventRecord(start, 0));
22     for (int i = 0; i < num_iters; ++i) {
23         CHECK_CUDA(cudaMemcpy(dst, src, bytes, kind));
24     }
25     CHECK_CUDA(cudaEventRecord(stop, 0));
26     CHECK_CUDA(cudaEventSynchronize(stop));
27
28     float ms = 0.0f;
29     CHECK_CUDA(cudaEventElapsedTime(&ms, start, stop));
30
31     CHECK_CUDA(cudaEventDestroy(start));
32     CHECK_CUDA(cudaEventDestroy(stop));
33
34     double seconds = ms / 1000.0;
35     double total_bytes = static_cast<double>(bytes) * num_iters;
36
37     // Bandwidth in GB/s
38     double gb_per_s = total_bytes / (seconds * 1.0e9);
39     return gb_per_s;
40 }
41
42 // Run test for one size, filling in the 4 bandwidth numbers
43 static void run_test_for_size(size_t bytes,
44                             int num_iters,
45                             double& h2d_pageable,
46                             double& d2h_pageable,
47                             double& h2d_pinned,
48                             double& d2h_pinned)
49 {
50     // -----
51     // Pageable host memory
52     // -----
53     {
54         // Host allocation (pageable)
55         float* h_buf = static_cast<float*>(std::malloc(bytes));
56         if (!h_buf) {
57             std::fprintf(stderr,
58                         "Failed to allocate pageable host memory (%zu bytes)\n",
59                         bytes);

```

```

60         std::exit(EXIT_FAILURE);
61     }
62
63     // Initialize host data
64     size_t n = bytes / sizeof(float);
65     for (size_t i = 0; i < n; ++i) {
66         h_buf[i] = static_cast<float>(i);
67     }
68
69     // Device allocation
70     float* d_buf = nullptr;
71     CHECK_CUDA(cudaMalloc(&d_buf, bytes));
72
73     // Measure H2D pageable
74     h2d_pageable = measure_bandwidth_single_direction(
75         d_buf, h_buf, bytes, cudaMemcpyHostToDevice, num_iters);
76
77     // Measure D2H pageable
78     d2h_pageable = measure_bandwidth_single_direction(
79         h_buf, d_buf, bytes, cudaMemcpyDeviceToHost, num_iters);
80
81     CHECK_CUDA(cudaFree(d_buf));
82     std::free(h_buf);
83 }
84
85 // -----
86 // Pinned host memory
87 // -----
88 {
89     // Host allocation (pinned)
90     float* h_buf = nullptr;
91     CHECK_CUDA(cudaHostAlloc(reinterpret_cast<void*>(&h_buf),
92                             bytes,
93                             cudaHostAllocDefault));
94
95     // Initialize host data
96     size_t n = bytes / sizeof(float);
97     for (size_t i = 0; i < n; ++i) {
98         h_buf[i] = static_cast<float>(i);
99     }
100
101     // Device allocation
102     float* d_buf = nullptr;
103     CHECK_CUDA(cudaMalloc(&d_buf, bytes));
104
105     // Measure H2D pinned
106     h2d_pinned = measure_bandwidth_single_direction(
107         d_buf, h_buf, bytes, cudaMemcpyHostToDevice, num_iters);
108
109     // Measure D2H pinned
110     d2h_pinned = measure_bandwidth_single_direction(
111         h_buf, d_buf, bytes, cudaMemcpyDeviceToHost, num_iters);
112
113     CHECK_CUDA(cudaFree(d_buf));
114     CHECK_CUDA(cudaFreeHost(h_buf));
115 }
116 }
117

```

```

118 void run_bandwidth_tests(const std::vector<size_t>& sizes,
119                          int num_iters,
120                          std::vector<BandwidthResult>& results)
121 {
122     results.clear();
123     results.reserve(sizes.size());
124
125     for (size_t bytes : sizes) {
126         double h2d_pageable = 0.0;
127         double d2h_pageable = 0.0;
128         double h2d_pinned = 0.0;
129         double d2h_pinned = 0.0;
130
131         run_test_for_size(bytes, num_iters,
132                          h2d_pageable, d2h_pageable,
133                          h2d_pinned, d2h_pinned);
134
135         BandwidthResult r;
136         r.sizeMB = static_cast<double>(bytes) / (1024.0 * 1024.0);
137         r.h2d_pageable = h2d_pageable;
138         r.d2h_pageable = d2h_pageable;
139         r.h2d_pinned = h2d_pinned;
140         r.d2h_pinned = d2h_pinned;
141
142         results.push_back(r);
143     }
144 }

```

5.4 main.cu

```

1 #include <stdio>
2 #include <vector>
3
4 #include "cuda_utils.cuh"
5 #include "bandwidth.h"
6
7 int main(int argc, char** argv)
8 {
9     // Use device 0 by default
10    CHECK_CUDA(cudaSetDevice(0));
11
12    // Sizes from 1 MB to 1 GB in powers of 2:
13    // 1, 2, 4, ..., 1024 MB
14    const int num_sizes = 11;
15    std::vector<size_t> sizes;
16    sizes.reserve(num_sizes);
17
18    for (int i = 0; i < num_sizes; ++i) {
19        size_t bytes = static_cast<size_t>(1) << (20 + i); // 2^(20+i)
20        sizes.push_back(bytes);
21    }
22
23    // Number of repetitions per measurement
24    int num_iters = 10;
25
26    std::vector<BandwidthResult> results;
27    run_bandwidth_tests(sizes, num_iters, results);

```

```

28
29 // CSV header
30 std::printf("size_MB,h2d_pageable_GB,d2h_pageable_GB,"
31            "h2d_pinned_GB,d2h_pinned_GB\n");
32
33 for (const auto& r : results) {
34     std::printf("%.0f,%.3f,%.3f,%.3f,%.3f\n",
35                 r.sizeMB,
36                 r.h2d_pageable,
37                 r.d2h_pageable,
38                 r.h2d_pinned,
39                 r.d2h_pinned);
40 }
41
42 CHECK_CUDA(cudaDeviceReset());
43 return 0;
44 }

```

5.5 Makefile

```

1 NVCC := nvcc
2 NVCC_FLAGS := -O2
3
4 SRCS := main.cu bandwidth.cu
5 HDRS := bandwidth.h cuda_utils.cuh
6
7 TARGET := bandwidth
8
9 $(TARGET): $(SRCS) $(HDRS)
10  $(NVCC) $(NVCC_FLAGS) $(SRCS) -o $(TARGET)
11
12 clean:
13  rm -f $(TARGET) bandwidth.csv bandwidth_plot.png

```

5.6 plot_bandwidth.py

```

1 import numpy as np
2 import matplotlib.pyplot as plt
3
4 # Load CSV data produced by ./bandwidth > bandwidth.csv
5 data = np.loadtxt('bandwidth.csv', delimiter=',', skiprows=1)
6
7 size_mb = data[:, 0]
8 h2d_pageable = data[:, 1]
9 d2h_pageable = data[:, 2]
10 h2d_pinned = data[:, 3]
11 d2h_pinned = data[:, 4]
12
13 plt.figure(figsize=(8, 6))
14
15 # H2D curves
16 plt.plot(size_mb, h2d_pageable, 'o-', label='H2D - Pageable')
17 plt.plot(size_mb, h2d_pinned, 'o-', label='H2D - Pinned')
18
19 # D2H curves

```



```
20 plt.plot(size_mb, d2h_pageable, 's--', label='D2H - Pageable')
21 plt.plot(size_mb, d2h_pinned, 's--', label='D2H - Pinned')
22
23 plt.xscale('log', base=2)
24 plt.xlabel('Array size (MB)')
25 plt.ylabel('Bandwidth (GB/s)')
26 plt.title('CPU-GPU Transfer Bandwidth vs Array Size')
27 plt.grid(True, which='both', linestyle=':')
28 plt.legend()
29 plt.tight_layout()
30 plt.savefig('bandwidth_plot.png', dpi=300)
31 plt.show()
```