

Tracking.jl: Accelerating multi-antenna GNSS receivers with CUDA

Can Özmaden¹

¹RWTH Aachen University

ABSTRACT

The use of advanced GNSS receiver architectures employing multi-antenna and multi-correlator signal processing poses a high computational demand on Software Defined Radio (SDR) modules. The literature provides promising results in offloading computationally burdensome tasks onto a Graphics Processing Unit (GPU). However, programmers face a myriad of design challenges and low-level optimization strategies while implementing a GPU-accelerated GNSS signal processing chain. This paper addresses the to highlight possible performance improvements of GPU-enabled GNSS receivers based on empirical findings from simulated data. A comparison between GPU algorithms and a parallelized CPU algorithm is carried out on two platforms: on a conventional mid-grade PC and an NVIDIA Jetson embedded system. Algorithms are rapidly prototyped in the Julia programming language and are benchmarked with varying optimization strategy combinations and signal input size. This paper introduces a technique of code replica generation by utilizing the GPU texture memory. The benchmarking results suggest real-time signal processing capabilities of the developed algorithms. The raw measurement data, the source code of the algorithms, and the experiment setup are released under an open-source license for reproducibility.

Keywords

Julia, GNSS, Signal Processing, SDR, CUDA, GPU, Optimization

1. Introduction

Utilization of antenna diversity via antenna arrays is a proven technique for GNSS receivers to mitigate the effects of Radio Interference (RFI), jamming, and spoofing [16]. Furthermore, multi-correlation allows GNSS receivers to alleviate multipath errors [24]. It also provides a robust estimation of a transfer function of a multi-antenna RF front end as shown in [17]. This aids the receiver to compensate for errors caused by the differences between the antenna channels. Moreover, multi-constellation and multi-frequency signal processing are increasingly used for Advanced Receiver Autonomous Integrity Monitoring (ARAIM) to meet strict demands of Safety of Life (SoL) applications [3]. Advanced modern receiver architectures combining the aforementioned techniques pose high demands on the computational resources of a software-defined receiver, making real-time processing a challenge. Most of the computational load can be attributed to the tracking module of a GNSS receiver, specifically to the correlation operation, which is the focus of this paper.

The literature provides promising results on dealing with the computational load by utilizing various parallelization approaches.

These include the use of bit-wise parallelism [14], single instruction multiple data (SIMD) instructions, and multithreading on multicore central processing units (CPUs) [6]. Additionally, the use of hardware acceleration by utilizing application-specific integrated circuits (ASICs) [1], digital signal processors (DSPs) [27], field-programmable arrays (FPGAs)[7] or GPUs [13, 20, 9, 26, 11].

This paper focuses on GPU acceleration. The use of GPUs to accelerate computationally burdensome tasks has become a popular choice in a plethora of applications, such as digital image processing, radar imaging, cryptography, neural networks, and many more. GPUs consist of massive parallel processors able to efficiently deal with a large amount of input data. In comparison with other hardware accelerators, GPUs are affordable, generally easier to program, and readily available in a large amount of existing consumer electronics devices. This makes GPUs highly desirable for SDRs. A significant amount of effort has been made by the vendors and the community to make programming GPU applications easier. Commonly used programming frameworks are NVIDIA's CUDA, AMD's ROCm, and the open-source OpenCL. In this paper, the CUDA framework is used.

Programmers are often challenged with various ways of optimizing the efficiency of GPU applications, especially their execution time. As mentioned above, real-time processing capability is crucial for software-defined GNSS receivers. Therefore, there is a need for a comparison of optimization strategies in designing a GPU-enabled GNSS SDR module. This paper strives to address this need. As a result, guidelines in fulfilling the real-time processing criteria can be provided based on empirical findings. In this paper, the comparison is carried out as a series of benchmarking experiments. Algorithms for the comparison are developed in a rapid prototyping fashion in Julia, a high-level high-performance language that is gaining popularity in the scientific and technical computing communities. GPU programming is performed by utilizing the CUDA.jl package which provides an interface for Julia to wrap low-level functionalities of the CUDA framework without sacrificing performance [2]. The developed algorithms extend an existing multi-antenna multi-correlator GPU-enabled GNSS SDR receiver from [21]. The tested optimization strategies rely purely on CUDA and are therefore not exclusive to Julia.

This paper is organized as follows: in Section 2 the signal model is presented, specifically the multi-correlation operation and the generation of replica signals. In Section 3 the experiment setup is introduced and the data acquisition methodology is described. Section 4 introduces the CUDA programming model and describes the implementation of the GPU algorithms, including a description of the novel texture memory based code replication technique. Section 5 focuses on the analysis of the obtained experimental data, to draw outcomes and guidelines for the best performing optimization

tion strategies. These outcomes are summarized, and a conclusion is presented in Section 6.

2. Signal Model

Following is an introduction to the notation used throughout this paper. Let k denote the k -th satellite from K total satellites, l the l -th correlator out of L total correlators, m the m -th antenna out of M total antennas, n the n -th sample out of N total samples of the received signal. Further, any analytic time-domain signal is denoted as $\underline{x}(t)$, $t \in \mathbb{R}$, and the respective sampled version as $\underline{x}[n]$, $n \in \mathbb{Z}$. Modeling of the signals and the receiver architecture follows and extends those presented in [4, 23].

Multi-Correlator

The output of a multi-correlator $R_{l,m}^{(k)}$ of the k -th satellite, l -th correlator m -th antenna can be expressed as:

$$\underline{r}_{l,m}^{(k)} = \sum_{n=1}^N \underline{r}_{\text{IF}}^{(k)}[n, m] \text{conj}(\hat{s}_{\text{ca}}^{(k)}[n]) \hat{s}_{\text{co},l}^{(k)}[n], \quad (1)$$

where r_{IF} denotes the received signal downconverted to an Intermediate Frequency (IF), \hat{s}_{ca} the carrier replica, \hat{s}_{co} the code replica, and $\text{conj}(\cdot)$ the complex conjugation operation. The correlation has to be performed over N samples and for each of M antennas, L correlators, and K satellites and is a prime example of the need for parallelization. The multiplication of the received signal with the conjugate of the carrier replica is often referred to as "carrier wipe-off". The despreading of the signal by multiplication with the code replica is often referred to as "code wipe-off".

Received Signal

The received signal downconverted to the IF can be expressed as:

$$\underline{r}_{\text{IF}}^{(k)}(t) = \underline{s}_{\text{IF}}^{(k)}(t) + \eta(t), \quad (2)$$

where

$$\underline{s}_{\text{IF}}^{(k)}(t) = c^{(k)}(t) \exp\left\{j \left(2\pi \left(f_{\text{IF}} + f_{\text{d}}^{(k)}\right) t + \phi_0^{(k)}\right)\right\}, \quad (3)$$

is the signal transmitted by the k -th satellite and downconverted to the IF-band, $c^{(k)}$ is the CDMA code of the k -th satellite, f_{IF} denotes the intermediate frequency, $f_{\text{d}}^{(k)}$ the Doppler shift of the k -th satellite relative to the receiver, $\phi_0^{(k)}$ the phase delay, and $\eta(t)$ the additive white Gaussian noise (AWGN) term.

Carrier Replica

The tracking module produces an estimate of the carrier of the received signal. This signal is commonly referred as the carrier replica \hat{s}_{ca} and can be expressed as:

$$\hat{s}_{\text{ca}}^{(k)}[n] = \exp\left\{j \left(2\pi \frac{f_{\text{IF}} + \hat{f}_{\text{d}}^{(k)}}{f_s} n + \hat{\phi}_0^{(k)}\right)\right\}, \quad (4)$$

where $\hat{f}_{\text{d}}^{(k)}$ is an estimate of the Doppler shift of the k -th satellite relative to the receiver, f_s denotes the sampling frequency of the receiver, and $\hat{\phi}_0^{(k)}$ denotes the estimated phase delay in radians.

Code Replica

Alongside the carrier replica generation, the tracking module also generates a code replica signal $\hat{s}_{\text{co},l}$ shifted by δ_l amount of chips

Table 1.: Parameters and their value variation used to execute different benchmarks presented in this paper.

Name	Platform #1	
	Desktop PC	NVIDIA Jetson AGX Xavier
OS	Windows 10 / Linux	Ubuntu 16.04 LTS
CPU Name	6-core Intel Core i5-9600K	8-core NVIDIA Tegra X1
CPU Clock frequency	3.70 GHz	2.27 GHz
GPU Name	NVIDIA GeForce GTX 1050 Ti	NVIDIA GeForce GTX 1080
CUDA Version	v11.5	v10.2
NVIDIA Driver Version	Linux-x86_64 495.46	Linux-ARM64 470.26
GPU Micro-architecture	Pascal (2016)	Tesla (2018)
GPU Clock frequency	1.39 GHz	1.46 GHz
GPU Number of SMs	6	384
GPU Shared memory	49152 bytes	122880 bytes
GPU Thread block size	1024	1024
GPU Grid size	(2147483647, 65535, 65535)	(2147483647, 65535, 65535)

that correspond to Δ_l amount of samples that are required by the l -th correlator. It can be expressed as follows:

$$\hat{s}_{\text{co},l}^{(k)}[n] = c^{(k)} \left[\text{mod} \left(\left\lfloor \frac{f_c + \hat{f}_{\text{c,d}}^{(k)}}{f_s} (n + \Delta_l) + \hat{\tau}^{(k)} \right\rfloor, \text{length}(c^{(k)}) \right) \right] \quad (5)$$

where $\hat{f}_{\text{c,d}}^{(k)}$ is an estimate of the Doppler shift of the code of the k -th satellite relative to the receiver, f_c denotes the code frequency, and $\hat{\tau}^{(k)}$ denotes the estimated phase delay in chips. The calculated code phase needs to be an integer value $z \in \mathbb{Z}$, therefore a flooring operation is needed. Furthermore, since the sample shift Δ_l pushes the earliest and latest chips out of bounds, a mod operation is needed to wrap the out-of-bounds code phase according to the length of the CDMA code. These two operations are computationally prohibitive and are to be avoided if possible.

3. Methodology

Testing Environment

The work in this paper is conducted as a series of experiments on two different platforms provided in Table 1. Platform #1 is a mid-grade desktop Personal Computer (PC) equipped with a 6-core Intel Core i5-9600K 3.70 GHz CPU and an NVIDIA GeForce GTX 1050 Ti GPU. Two Operating Systems (OSes) are installed on separate solid-state memory units: Microsoft Windows 10 Home (Build 19042) and Linux (Endeavour OS, rolling release distribution with a Long Time Support (LTS) 5.15.14-1-lts Linux Kernel). Benchmarks were conducted under both Windows and Linux on Platform #1, without any differences found. For Platform #1, Linux benchmarks are used in this paper.

Platform #2 is an NVIDIA Jetson AGX Xavier Development Kit device. It is equipped with a Tegra System-on-a-Chip (SoC) combining a proprietary 8-core 2.27 GHz NVIDIA ARMv8 CPU and an NVIDIA GPU with Volta microarchitecture.

The two platforms cover interesting aspects of possible implementation scenarios of a GNSS SDR receiver. Platform #1 is most similar to commonplace PCs in research institutes or companies, whereas Platform #2 provides an insight into the performance of a mobile embedded GPU-enabled GNSS SDR, commonly used for autonomous applications.

Signal Generation

Signals of differing lengths and properties are generated according to the parameters provided in Table 2. The sampling frequency is capped at a maximum of 400 MHz. The number of antennas is switched between one and four. The four antenna case is selected to emulate DLR's GALANT receiver with a 2x2 antenna array [7]. The number of correlators is selected to include the classical early-prompt-late correlator, and additionally a correlator with seven taps. In this paper, the focus lies on the optimization of the correlation for one satellite channel. The parallelization over multiple channels can be easily extended as presented in [26]. However, this also requires strict synchronization between the channels for coherent integration over fixed length. Another possibility is asynchronous kernel launches. Both are not explored in this paper for the sake of simplicity in providing optimization guidelines.

Table 2. : Parameters and their value variation used to execute different benchmarks.

	GPS L1 C/A
Signal duration T	1ms
Carrier Doppler f_{IF}	1.5 kHz
Initial carrier phase delay ϕ	0 rad
Initial code phase delay τ	0 chips
Correlator shift δ	0.5 chips
Number of satellites K	1
Number of correlators L	3, 7
Number of antennas M	1, 4
Sampling Frequency T/N	2.048 MHz - 262.144 MHz

The generated data signals assume a constant data bit stream of ones for the entire signal duration T . They are subsequently spread by multiplication with the respective CDMA code of the constellation and upconverted to f_{IF} . Differing from Equation 2, no AWGN is added to the signal, as the time complexity of the operations remains the same. The pseudocode of the described signal generation procedure is provided in Algorithm 1. Depending on the processing unit under test, the received signal is kept in the primary memory of the host, or transferred to the global memory of the GPU. The GPU signal is generated with single precision (from hereto aliased as F32). The CPU signal is also initialized with single-precision floating-point values, to maximize SIMD efficiency. The received signal can be additionally broken into a Quadrature (Q) and an In-phase component (I). These are realized via a structure of arrays holding complex singles (from hereto aliased as \mathbb{C}_{F32}). This ensures a coalesced memory access when implementing I/Q-signal processing, which increases the efficiency of the algorithms, discussed in more detail in Section 4.

Algorithm Naming Scheme

A consistent naming scheme for the algorithms under test is used throughout this paper. The name consists of keywords relating to various properties of the algorithm. Naming starts with a number describing the overall procedural sequence presented as a block diagram in Figure 1. The number following that is connected to the parallel reduction algorithm at use is described in detail in Section 4. A string is added to the number of the reduction algorithm, indicating original implementation by the keyword "pure", a complex reduction extension by "cplx" and a "cplx_multi" keyword indicating the total fusion of all multi-antenna multi-correlator kernel invocations. Additionally, if the kernel utilizes the novel texture

memory code replica generation described in Section 4 a keyword "textmem" is added. To give an example, an algorithm following the 3rd procedure from Figure 1, utilizing a fully fused 3rd reduction algorithm and texture memory, is designated the codename "1_3_cplx_multi_textmem".

Benchmarking

Measurements of execution times are inherently noisy due to OS scheduling events, clock frequency jitter, etc. Therefore, there is a need for a statistical evaluation of the collected timings, i.e. averaging over numerous executions, performing calculations of estimators.

Runtimes in this paper were collected using the BenchmarkTools.jl package. The package tunes the execution parameters, such as the number of executions, automatically. Four estimators are provided as the benchmark result. These are the *Minimum*, *Median*, *Mean*, and *Maximum* estimators. In [5] the authors of the package assess the minimum estimator, rather than median or mean, to be the most robust one, as execution times are heavily right-skewed. This renders the mean estimator heavily influenced by the outliers, and the maximum estimator being purely an outlier. Therefore the presented execution times are taken from the Minimum estimator unless otherwise noted.

Additionally, the GPU benchmarks are externally verified via NVIDIA profiler tools such as Nsight Systems and Nsight Compute. The verification is done by annotating the relevant functions with NVTX ranges and examining reports.

Project Automation

Algorithms under test and benchmarking scripts are all compounded into a Julia module utilizing the DrWatson.jl package [8]. This provides reproducibility of the results found in this paper, as the module keeps a list of all dependencies and their version and can be instantiated via a simple command. Additionally, the algorithms used for the comparison in this paper are tested utilizing the standard Julia testing ecosystem. The repository containing the source code and scripts for the benchmarks acquired and figures found in this paper can be found on Github [28]. Raw data obtained

Algorithm 1 Generate received signal

```

1: procedure GEN_SIGNAL!(signal  $\in \mathbb{C}_{F32}^{N \times M}$ ,  $f_s$ ,  $\phi$ ,  $\tau$ , codes)
2:   for all  $k \in 1 \dots K$  do
3:     for all  $m \in 1 \dots M$  do
4:       code_phases  $\leftarrow \frac{f_c^{(k)}}{f_s} \odot (0 \dots N - 1) \oplus \tau^{(k)}$ 
5:       carrier_phases  $\leftarrow 2\pi \odot \frac{f_{IF}^{(k)}}{f_s} \odot (0 \dots N - 1) \oplus \phi^{(k)}$ 
6:       upsampled_code  $\leftarrow$  codes[mod{[code_phases], length(c)}, k]
7:       signal.re[:, m]  $\leftarrow$  cos.(carrier_phases)  $\odot$  upsampled_code
8:       signal.im[:, m]  $\leftarrow$  sin.(carrier_phases)  $\odot$  upsampled_code
9:     end for
10:   end for
11: end procedure

```

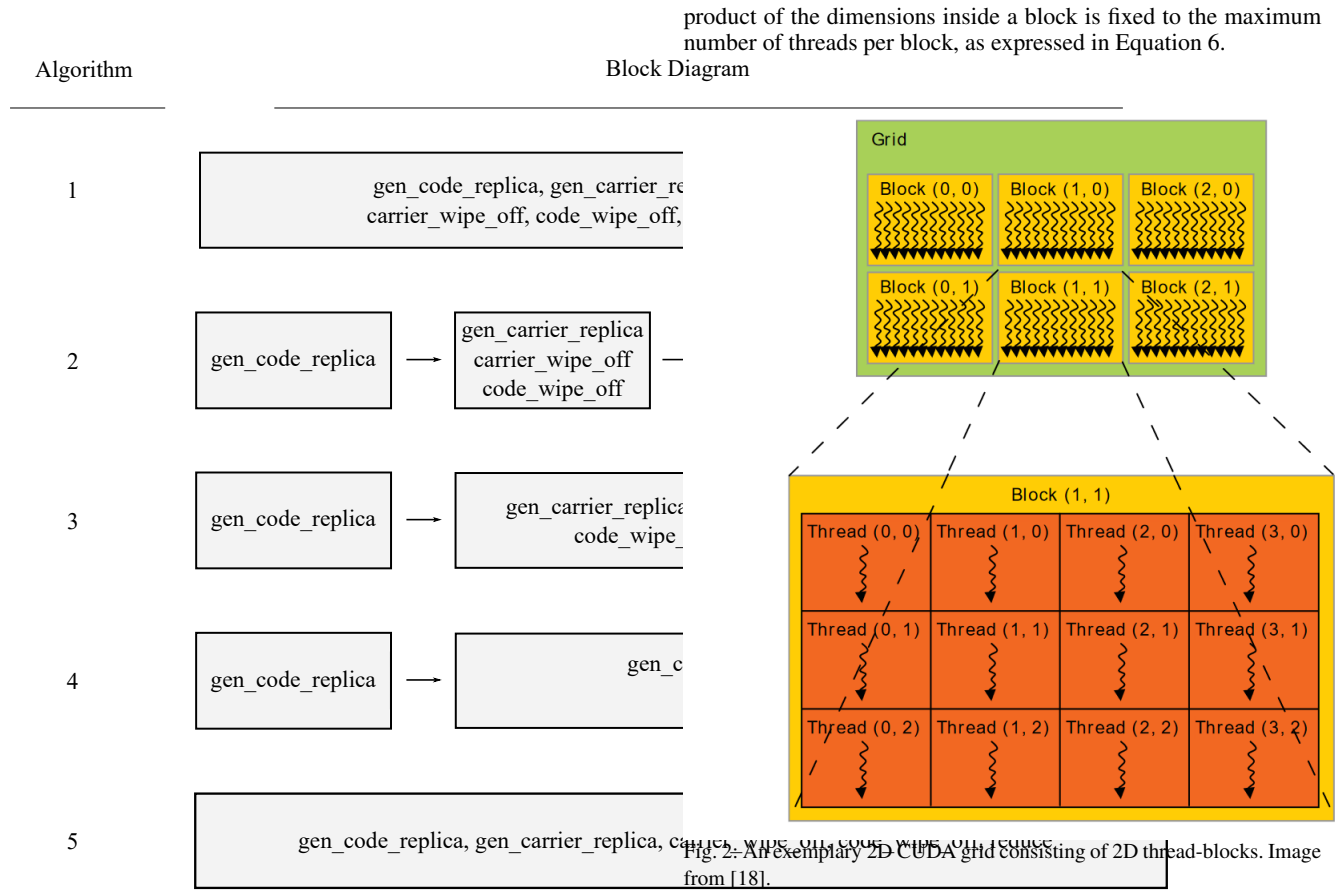


Fig. 1: Block-diagram description of the algorithms being tested in this paper.

from the benchmarks on the platforms described in this paper is also made available [19].

4. Algorithm Design

CUDA

CUDA is a programming framework published by NVIDIA for programming the GPU. Currently, it provides first-grade support for C/C++ and FORTRAN. CUDA is designed to ease the programming of parallel applications on the GPU, by providing an LLVM-based compiler for PTX, the NVIDIA GPU native Instruction Set Architecture (ISA). Albeit, effective CUDA programming requires a solid understanding of the low-level resources of the GPU. In the following, a short explanation of the CUDA framework is given. The information stems from NVIDIA' s "CUDA C++ Programming Guide" [18].

Functions executed on the GPU are referred to as *kernels*. Each kernel is launched on a *grid* consisting of *thread-blocks*. This hierarchy is demonstrated in Figure 2. The smallest units, *threads*, are issued in a group of 32 to perform the same instruction. This grouping is called a *warp*. The next level of thread grouping is called a *block* or a *thread-block*. These differ in their maximum allowed size depending on the hardware. CUDA provides three-dimensional indexing along the grid, and inside an individual thread-block. The

product of the dimensions inside a block is fixed to the maximum number of threads per block, as expressed in Equation 6.

$$\text{max_threads_per_block} = \text{threads}_x \times \text{threads}_y \times \text{threads}_z. \quad (6)$$

Thus, if one were to fix the z and y dimensions of a thread-block, the maximum allowed number of threads in the x dimension can be easily found via:

$$\text{threads}_x = \text{max_threads_per_block} \div (\text{threads}_y \times \text{threads}_z). \quad (7)$$

On the grid level, no limits are tying the x, y , and z dimensions together.

Each thread executes in parallel on the GPU, meaning one of the priorities in parallelization optimization is providing the GPU with enough resources to saturate the number of threads performing calculations in parallel. This is often measured via an *occupancy* metric. Kernels exhibiting high occupancy have the potential to better hide processing latencies, however, they are not per se the best performing as shown in [25]. In line with previous work on GNSS GPU signal processing, this paper strives for optimum occupancy at every kernel launch via a call to a launch configuration subrou-

GPU Memory

The knowledge of various memory resources of the GPU is essential for improving the efficiency of kernels. The memory available to all threads regardless of grouping is called the *global memory* of the GPU. Inside each thread-block threads can access a *shared*

	memory location →							
indices			0	1	2	3		
wrap	2	3	0	1	2	3	0	1
clamp	0	0	0	1	2	3	3	3
mirror	2	1	0	1	2	3	2	1
border	0	0	0	1	2	3	0	0

Fig. 3: Visualization of various addressing modes available for the texture memory in CUDA.

memory. Shared memory is faster than global memory. Additionally, each thread has its local on-chip L2 cache memory for local variables.

Furthermore, there is another type of read-only memory available globally called *texture memory* that possesses some unique properties. Firstly, the speed of accessing the texture memory is increased by caching. Secondly, spatially close memory locations addressed in an unpredictable order get a speed-up. And most importantly for the application in this paper, rules for out-of-bounds accesses and an interpolation method for non-integer indices can be defined upon allocation.

CUDA defines three texture memory addressing modes to define behavior upon an out-of-bounds access: *wrap*, *clamp*, *mirror* and *border*. Three interpolation methods are provided to deal with non-integer indices: *linear*, *bilinear* and *nearest neighbour* (CUDA.jl naming). Additionally, a boolean value called *normalized_coordinates* can be passed to specify if the indices are in the $[0, 1)$ range. To make use of the wrap addressing mode, the indices must be normalized. A simple example is provided in Figure 3 to illustrate the functioning of the described addressing modes.

Parallel Reduction

The cornerstone of the tracking module is the correlation process as expressed in Equation 1. This operation requires each element of the vector to be multiplied element-wise with the others, and the result of these products must be summed over N samples. While vector multiplication is trivially parallelizable, the parallelization of the sum requires some thought. Each partial sum of two elements can be computed independently of the other partial sums, however, some communication between threads is needed to facilitate a step-wise reduction of a vector of length N into a scalar value.

In his seminal work, Harris has presented ways to optimize the parallel reduction in CUDA [12]. He introduces seven kernels, each subsequent improving on the preceding. The reduction strategy involves the use of a binary tree-like summing. Harris' reduction kernel numbering is taken one-to-one in the naming of the algorithms in this paper as described in Section 3.

One of the main issues of parallel reduction is the innate need for synchronization between parallel threads. Before CUDA version 9, it was only possible to synchronize threads across a block. This meant that only a vector of length matching the maximum allowed threads per block could have been reduced in a *single-pass* (sin-

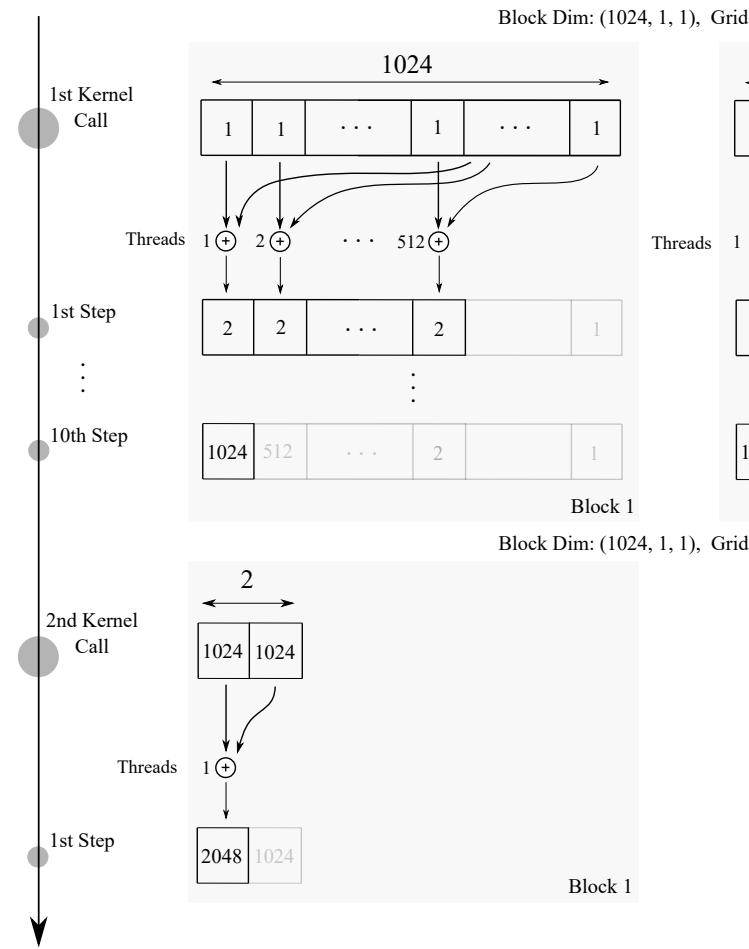


Fig. 4: Visualization of a tree-like reduction algorithm from [12], using a coalesced memory access. The kernel is first launched on a grid of two blocks, which produce their respective partial sums for the 2nd kernel launch to sum.

gle kernel invocation). For greater input sizes, a partial sum vector of length equal to the number of blocks has to be allocated. The second kernel invocation reduces the values of partial sums across blocks. This is called *multi-pass* reduction, in contrast to *single pass*. Figure 4 illustrates the multi-pass reduction over 2048 elements by using the 3rd kernel from [12].

With the introduction of *Cooperative Groups* in CUDA version 9, it has become possible to synchronize GPUs at every level, including thread-blocks across a grid. Therefore, a kernel launched within a cooperative group can reduce across a block and subsequently across all blocks in a single-pass.

Another approach to summing a vector on a GPU is via the use of atomic operations. However, atomic operations come at a cost of not utilizing the full bandwidth of the GPU and therefore must be used sparingly. In this paper atomic reduction is used to eliminate the second reduction kernel call in algorithms 4 and 5 as illustrated in Figure 1.

Kernel Implementations

There are various ways to improve the runtime efficiency of a kernel. Firstly, a distinction has to be made if the kernel is a so-called *compute-bound* or a *memory-bound* kernel. For compute-bound kernels, the mathematical operations performed by each thread have to be optimized. For the memory-bound ones, accesses to the memory resources are to be accelerated.

Four operations in the tracking module have the potential to be parallelized. These are the code replica generation from Equation 5, carrier replica generation from Equation 4, downconversion, and correlation from Equation 1. These can be defined as separate kernels or fused in various ways as shown in Figure 1. The downconversion kernel is the only one of the four that can be described as a compute-bound one. The rest are memory-bound kernels and are to be optimized to reach the peak throughput of the GPU.

As discussed previously, the correlation operation involves two stages, firstly a vector element-wise multiplication, and subsequently a summation over all N elements. As an example, a GPS L1 C/A signal of 1 ms duration is considered. Sampled at 2.048 MHz sampling frequency this produces 2048 I/Q samples ($N = 2048$). With an assumption of early-prompt-late correlation ($L = 3$) and four antennas ($M = 4$) this results in 48 reduction kernel invocations for each 1ms chunk. This is firstly due to the maximum amount of threads in a block, which is 1024 for both platforms from Table 1. Therefore, a single-pass reduction with the algorithms from [12] is not possible. Secondly, the reduction has to be performed on both the I and Q components of the signal.

In this paper, kernel calls are reduced by implementing extended versions of the original reduction algorithm. These are the so-called "cplx" and "cplx_multi" kernels. The "cplx" kernel assigns each thread to deal with both in-phase and quadrature components of the signal and has to perform two passes for each $m \in M$ and $l \in L$. Whereas, the "cplx_multi" reduction kernel is totally parallelized and performs two passes only once per $k \in K$.

The CUDA kernel currently implemented in [21] is a variant of the first algorithm from Figure 1, albeit without the texture memory code generation. Additionally, this is the only kernel that possesses a multi-dimensional grid. The newly devised ones in this paper are all developed to be one-dimensional kernels, assigning more workload to each thread.

Code Replica Generation from Texture Memory

Generation of the code replica is computationally prohibitive due to the modding and flooring operations (see Equation 5). The proposed receiver architecture loads all CDMA codes $c^{(k)}$ of K satellites directly into the texture memory of the GPU in their original length at the beginning of the operation. Later, a code replica can be generated as shown in Algorithm 3. This has the advantage of the sparing of the mod and floor operations, due to the selected memory addressing modes.

Authors of [22] have recently introduced a texture-memory-based algorithm for BeiDou signal generation by loading codes into the texture memory. However, the specific addressing and filtering modes used are not mentioned. The speed-up is therefore only achieved via caching. Authors of [15] implement the code and carrier replica generation via tables stored in texture memory. The use of addressing modes is mentioned but the algorithm is not presented.

In this paper, the technique of code replica generation from texture memory is specifically applied to a multi-antenna GNSS receiver. Additionally, global memory and texture memory are used in coop-

Algorithm 2 Complex multi reduction kernel extended from Harris #3

```

1: procedure REDUCE_CPLX_MULTI!(output, input)
2:   tid  $\leftarrow$  threadIdx.x
3:   iq_offset  $\leftarrow$  blockDim.x
4:   shmem := DynamicSharedMemory  $\in \mathbb{F}_{32}^{(2N) \times M \times L}$ 

5:   # Load input into shared memory
6:   n = blockDim.x  $\times$  blockDim.x + threadIdx.x
7:   for all  $m \in M, l \in L$  do
8:     if  $n \leq \text{length}(\text{input})$  then
9:       shmem[tid, m, l]  $\leftarrow$  input.re[n, m, l]
10:      shmem[tid + iq_offset, m, l]  $\leftarrow$  input.im[n, m, l]
11:    end if
12:  sync_threads

13:  # Perform tree-like reduction in shared memory
14:  for  $s = \text{blockDim.x} \div 2, s \neq 0, s = s \div 2$  do
15:    if  $\text{tid} - 1 < s$  then
16:      shmem[tid, m, l]  $\leftarrow$  shmem[tid +
17:      s, m, l]
18:      shmem[tid + iq_offset, m, l]  $\leftarrow$  shmem[tid +
19:      s + iq_offset, m, l]
20:    sync_threads
21:  end if
22:  end for

23:  # First thread holds the result
24:  if  $\text{tid} == 1$  then
25:    output.re[blockIdx.x, m, l]  $\leftarrow$ 
26:    shmem[1, m, l]
27:    output.im[blockIdx.x, m, l]  $\leftarrow$  shmem[1 +
28:    iq_offset, m, l]
29:  end if
30:  end for
31: end procedure

```

Table 3. : Parameters for testing the reduction algorithms.

Parameter Name	Values
Number of elements, N	2048 - 32768
Number of antennas, M	1, 4
Number of correlators, L	3, 7

eration, as the mutable code replica array is contained in the global memory of the GPU, initiated originally from the chip-length codes in the texture memory.

5. Analysis

Reduction

An evaluation of the "cplx_multi" reduction introduced in Section 4 under Algorithm 2 is carried out. Parameters used in the experiment to assess the developed algorithm are provided in Table 3. The results are provided in Figure 5. One can observe "cplx_multi" outperforming "cplx" and "pure", due to the fusion of the kernel invocations. Therefore only "cplx_multi" kernels are presented in this paper.

Algorithm 3 Generate code replica from texture memory

```

1: # Allocate a texture memory with the following addressing
   modes
2: codes := TextureMemoryArray(warp, nearest_neighbour, normalized_coordinates)

3: # Load the PRN codes into the texture memory of the GPU
4: for all  $k \in 1 \dots K$  do
5:   codes[:,  $k$ ]  $\leftarrow c^{(k)}$ 
6: end for

7: # Allocate a global memory code replica with an appended
   length of  $\Delta$  to accomodate shifts
8:  $\hat{s}_{co} := \text{GlobalMemoryArray} \in \mathbb{C}_{F_{32}}^{N+\Delta}$ , where  $\Delta := \max(\Delta_l)$ 

9: # Call for each  $k$ 
10: procedure GEN_CODE_REPLICA_KERNEL!( $\hat{s}_{co} \in \text{TextureMemoryArray}\{\mathbb{C}_{F_{32}}^{N+\Delta}\}, k, f_s, f_c, \text{length}(c^{(k)}), \tau, \text{codes}$ )
11:   for all  $n \in 1 \dots N$  do
12:      $\hat{s}_{co} \leftarrow \text{codes}\left[\left(\frac{f_c}{f_s}(n - \Delta) + \tau\right) \div \text{length}(c^{(k)}), k\right]$ 
13:   end for
14: end procedure

```

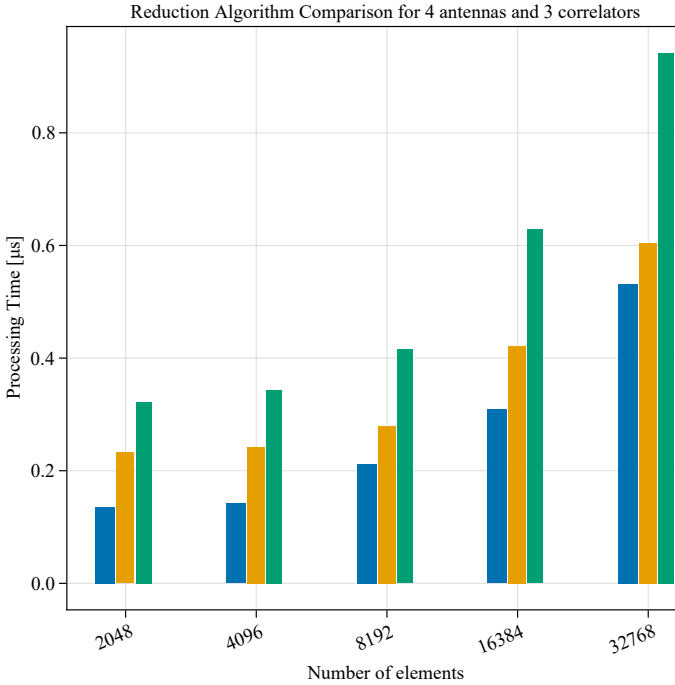


Fig. 5: Runtime analysis of the three reduction algorithms: "pure", "cplx", and "cplx_multi"

Texture memory

The generation of the code replica from the texture memory of the GPU requires the use of the nearest neighbor addressing fil-

Table 4.: Statistical analysis of the relative code phase error between the global memory and texture memory code replication algorithms.

Minimum relative error	0%
Mean relative error	0.03%
Median relative error	0.02%
Maximum relative error	3.17%

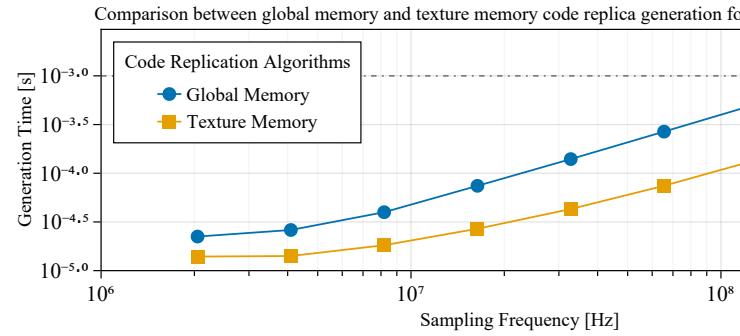
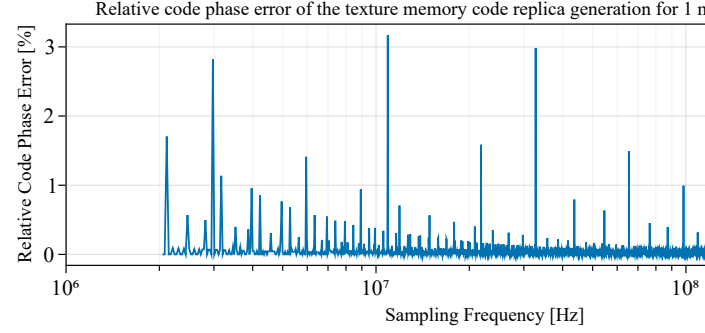


Fig. 6: Above: Code phase error due to the discrepancy between global memory and texture memory code replication algorithms. Below: runtime analysis between the global and texture memory code replication algorithms. The dashed line indicates the real-time bound.

tering and wrap addressing mode. This, however, differs in its implementation from the modding and flooring of code phases introduced in Equation 5. The discrepancy results in a slightly different code phase being calculated by the algorithms generating the code replica from global memory and texture memory. An analysis of the introduced error by this algorithm is carried out. The results can be seen in Figure 6, along with a timing experiment against the global memory counterpart. As one can observe, the discrepancy results in a negligible relative code phase error, and the speed-up is significant. Short statistical summary is found in Table 4. One can conclude, the speed-up of using the code replication via texture memory is a well-founded trade-off.

Algorithms

Due to the multitude of various kernel combinations, only a few selected ones are presented in this paper. The experiment setup and related functions exist in the repository [28]. For the sake of brevity, only kernels with the 4th reduction algorithm are shown in this paper. They are decently performing and work effortlessly on non-power-of-two input sizes.

Upon close analysis of Figures 7 and 8, one can conclude the results suggest that differences between the architectures of the GPUs can not be ignored. The algorithms performing best on Platform #1

are the worst-performers for Platform #2, and vice-versa. The results also suggest merit in exploring parallelization on the CPU, as the CPU SIMD algorithms from [21] are found to consistently outperform the GPU on Platform #1. This is, however, not always the case on Platform #2. A possible explanation for this is the fact that Platform #2 possesses an ARM chip, with SIMD capabilities unmatched to the Desktop-grade Intel Core i5-9600K. Another factor is also the higher CPU clock frequency of Platform #1. The relative mismatch between the power of the GPU and the power of the CPU on Platform #2, leaning towards the GPU, also plays a role. The CUDA algorithm currently implemented in [21] is being outperformed by the algorithms devised in this paper. Overall it can be said that the 4th algorithm under the name "4_4_cplx_multi_textmem" is the best performing for most of the time on Platform #1. Whereas algorithm 2 under the name "2_4_cplx_multi_textmem" performs the best on Platform #2, for cases utilizing under seven correlators. This is somewhat surprising, as it suggests launch configurations of individual kernels playing a bigger role on the Jetson device, and/or that Platform #2 has shorter API overhead. Moreover, this could be attributed to the somewhat slower atomic operations on Platform #2, as these are utilized in algorithms 4 and 5.

The results suggest real-time operation capabilities of the developed algorithms, even under introduced multi-antenna and multi-correlator load. Most algorithms can be executed in real-time under 20 MHz sampling frequency both on Platform #1 and #2. Some algorithms reach around 40 MHz processing capability with 4 antennas or even 100 MHz with a single antenna and an early-prompt-late correlator. Only the case of processing GPS L5 I5 signals with 4 antennas and 7 correlators remains an unsolved challenge.

6. Conclusion

In this paper, a series of benchmarks is carried out across different GPU algorithms with varying signal input sizes, ranging in the sampling frequency, number of antennas, and number of correlators from the GPS constellation at two frequency bands and respective codes: GPS L1 C/A, GPS L5 I5. The GPU algorithms differ in their implementation of the parallel reduction, use of texture memory for code replication, and the fusion of subsequent kernels. Benchmarking measurements are collected via a statistical benchmarking package, accounting for the differences in timings due to system noise. For the GPU algorithms, the measurement results are verified via external NVIDIA profiling tools. Benchmarking is carried out under the recommended guidelines for limiting system noise [10].

The algorithms are implemented in Julia, an open-source, dynamic, just-in-time compiled, high-performance high-level language. GPU kernel functions utilize the CUDA.jl package, which provides an interface for Julia to wrap functionalities of CUDA C at no performance cost. The correctness of the results calculated by the algorithms is verified via the testing ecosystem of Julia. The outcomes of this paper, however, are not limited to the Julia language, as these rely purely on CUDA.

A comparison between the algorithms has shown differing results for the two platforms under test. For the Desktop PC, an implementation with the code replication kernel separate from a fused carrier wipe-off and single-pass reduction kernel shows the best performance for lower sampling frequencies of the receiver (referred to as "4_4_cplx_multi_textmem" in this paper). With the increasing sampling frequency, however, the difference between the tested algorithms diminishes. The SIMD-based CPU algorithms outperform the GPU algorithms on a Desktop PC under test (referred to

Desktop PC, Intel Core i5-9600K @ 3.70 GHz, NVIDIA

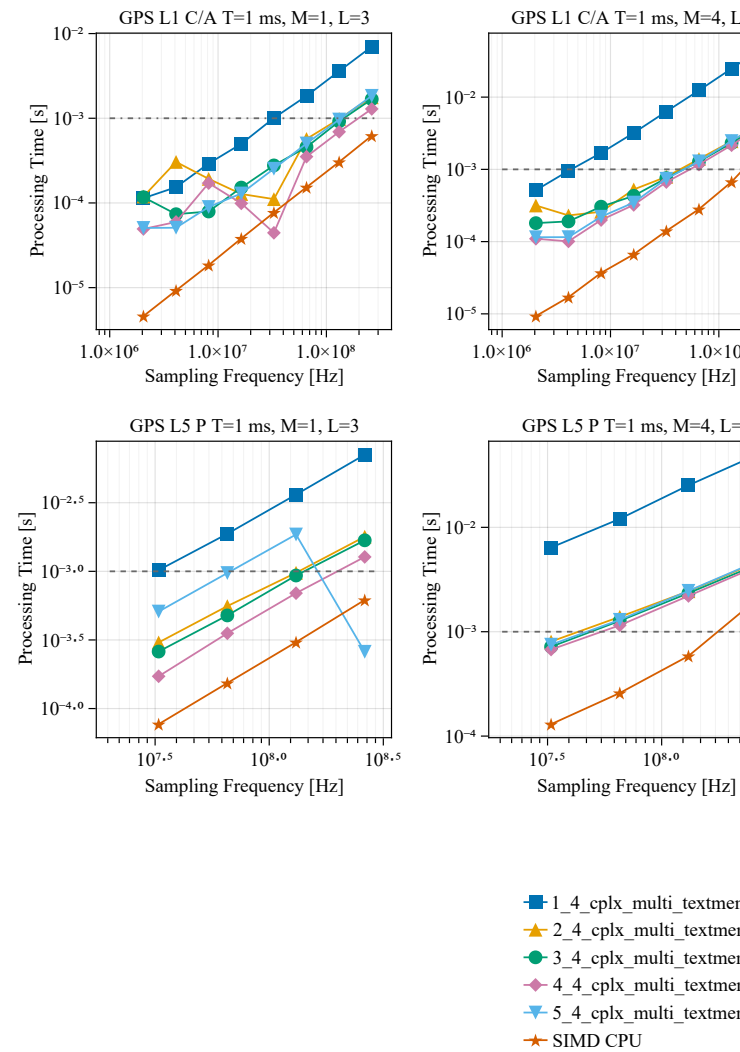


Fig. 7: Processing time of correlation of a 1ms GPS signal on Platform #1. The dashed line indicates the real-time processing bound.

as Platform #1), albeit fall behind on the NVIDIA Jetson device (referred to as Platform #2). The fully separate kernel algorithm performed the best on the NVIDIA Jetson device (referred to as "2_4_cplx_multi_textmem"). All the developed algorithms show real-time capabilities under the assumption of 4 antennas and 20 MHz sampling frequency, with some algorithms reaching the 100 MHz mark for single antenna operation.

A repository containing the experiment setup and the source code of the algorithms is made available on Github [28]. Raw data obtained from the experiments on the two platforms described in this paper is also made available under a Creative Commons license [19].

To the best of the authors' knowledge, this is the first publication describing the code replication algorithm on the GPU utilizing both

texture memory and global memory. Additionally, it is the first to utilize the Julia programming language to implement CUDA GNSS signal processing algorithms.

7. References

- [1] Nikola Basta, Achim Dreher, Stefano Caizzzone, Matteo Sgammini, Felix Antreich, Götz Kappen, Safwat Irteza, Ralf Stephan, Matthias A Hein, Eric Schäfer, et al. System concept of a compact multi-antenna GNSS receiver. In *2012 The 7th German Microwave Conference*, pages 1–4. IEEE, 2012.
- [2] Tim Besard, Christophe Foket, and Bjorn De Sutter. Effective Extensible Programming: Unleashing Julia on GPUs. *IEEE Transactions on Parallel and Distributed Systems*, 2018. doi:10.1109/TPDS.2018.2872064. 1712.03112.
- [3] Juan Blanch, Todd Walter, Per Enge, Stefan Wallner, Francisco Amarillo Fernandez, Riccardo Dellago, Rigas Ioannides, Ignacio Fernandez Hernandez, Boubeker Belabbas, Alexandru Spletter, and Markus Rippl. Critical Elements for a Multi-Constellation Advanced RAIM. *NAVIGATION*, 60(1):53–69, 2013. doi:https://doi.org/10.1002/navi.29. https://onlinelibrary.wiley.com/doi/pdf/10.1002/navi.29.
- [4] Kai Borre, Dennis M Akos, Nicolaj Bertelsen, Peter Rinder, and Søren Holdt Jensen. *A software-defined GPS and Galileo receiver: a single-frequency approach*. Springer Science & Business Media, 2007.
- [5] Jiahao Chen and Jarrett Revels. Robust benchmarking in noisy environments. *CoRR*, abs/1608.04295, 2016. 1608.04295.
- [6] Yu-Hsuan Chen, Jyh-Ching Juang, Jiwon Seo, Sherman Lo, Dennis M. Akos, David S. De Lorenzo, and Per Enge. Design and Implementation of Real-Time Software Radio for Anti-Interference GPS/WAAS Sensors. *Sensors*, 12(10):13417–13440, 2012. doi:10.3390/s121013417.
- [7] Manuel Cuntz, Lukasz Greda, Marcos Heckler, Andriy Konovaltsev, Michael Meurer, and Lothar Kurz. "GALANT - Architecture of a Real-Time Safety of Life Receiver". In *ION GNSS 2009*, September 2009.
- [8] George Datseris, Jonas Isensee, Sebastian Pech, and Tamás Gál. DrWatson: the perfect sidekick for your scientific inquiries. *Journal of Open Source Software*, 5(54):2673, 2020. doi:10.21105/joss.02673.
- [9] Carles Fernández-Prades, Javier Arribas, and Pau Closas. Accelerating GNSS software receivers. In *Proceedings of the 29th International Technical Meeting of The Satellite Division of the Institute of Navigation (ION GNSS+ 2016)*, pages 44–61, 2016.
- [10] Bill Fiser and Sebastian Jodłowski. Best Practices When Benchmarking CUDA Applications, 2019.
- [11] Chengjun Guo, Bingyan Xu, and Zhong Tian. Research on multi-constellation GNSS compatible acquisition strategy based on GPU high-performance operation. *EURASIP Journal on Wireless Communications and Networking*, 2018(1):1–10, 2018.
- [12] Mark Harris et al. Optimizing parallel reduction in CUDA. *Nvidia developer technology*, 2(4):70, 2007.
- [13] Kamran Karimi, Aleks G Pamir, and M Haris Afzal. Accelerating a cloud-based software GNSS receiver. *International Journal of Grid and High Performance Computing (IJGHPC)*, 6(3):17–33, 2014.
- [14] B.M. Ledvina, M.L. Psiaki, S.P. Powell, and P.M. Kintner. Bit-wise parallel algorithms for efficient software correlation applied to a GPS software receiver. *IEEE Transactions on Wireless Communications*, 3(5):1469–1473, 2004. doi:10.1109/TWC.2004.833467.
- [15] Qiushi Li, Zheng Yao, Hong Li, and Mingquan Lu. A CUDA-based real-time software GNSS IF signal simulator. In *China Satellite Navigation Conference (CSNC) 2012 Proceedings*, pages 359–369. Springer, 2012.
- [16] D-J Moelker, Edwin van der Pol, and Yeheskel Bar-Ness. Adaptive antenna arrays for interference cancellation in GPS and GLONASS receivers. In *Proceedings of Position, Location and Navigation Symposium-PLANS'96*, pages 191–198. IEEE, 1996.
- [17] Michael Niestroj, Marius Brachvogel, Soeren Zorn, and M Meurer. Estimation of antenna array manifolds based on sparse measurements. In *Proceedings of the 31st International Technical Meeting of the Satellite Division of The Institute of Navigation (ION GNSS+ 2018)*, pages 4004–4011, 2018.
- [18] NVIDIA. CUDA C++ Programming Guide v11.6.0, January 2022.
- [19] Can Özmaden. GPUAcceleratedTracking Raw Data, January 2022. doi:10.5281/zenodo.5933726.
- [20] Kwi Woo Park, Sangwoo Lee, Min Joon Lee, Sunwoo Kim, and Chansik Park. An accelerated signal tracking module using a heterogeneous multi-GPU platform for real-time GNSS software receiver. In *2015 IEEE Global Conference on Signal and Information Processing (GlobalSIP)*, pages 1412–1416. IEEE, 2015.
- [21] Soeren Schoenbrod, Michael Niestroj, Erik Deinzer, Can Ozmaden, and Kristoffer Carlsson. JuliaGNSS/Tracking.jl: v0.14.10, January 2022. doi:10.5281/zenodo.5870363.
- [22] Pengliang Shi, Shunxiao Wu, Tian Zeng, and Rui Xue. Satellite Navigation Simulation Signal Generation Method Based on GPU Acceleration. *Journal of Mathematics and Informatics*, 16:20, March 2021.
- [23] Peter JG Teunissen and Oliver Montenbruck. *Springer handbook of global navigation satellite systems*, volume 1. Springer, 2017.
- [24] R.D.J. van Nee. The Multipath Estimating Delay Lock Loop. In *IEEE Second International Symposium on Spread Spectrum Techniques and Applications*, pages 39–42, 1992. doi:10.1109/ISSSTA.1992.665623.
- [25] Vasily Volkov. Better performance at lower occupancy. In *Proceedings of the GPU technology conference, GTC*, volume 10, page 16. San Jose, CA, 2010.
- [26] Liangchun Xu, Nesreen I Ziedan, Xiaoji Niu, and Wenfei Guo. Correlation acceleration in GNSS software receivers using a CUDA-enabled GPU. *GPS solutions*, 21(1):225–236, 2017.
- [27] Qingxi Zeng, Qing Wang, Shuguo Pan, and Chuanjun Li. A GPS L1 Software Receiver Implementation on a DSP Platform. In *2008 First International Conference on Intelligent Networks and Intelligent Systems*, pages 612–615, 2008. doi:10.1109/ICINIS.2008.99.
- [28] Can Özmaden. ozmaden/GPUAcceleratedTracking: v1.0, January 2022. doi:10.5281/zenodo.5933659.

NVIDIA Jetson AGX Xavier Development Kit, Tegra SoC, ARMv8 @ 2.27 GHz, Volta GPU

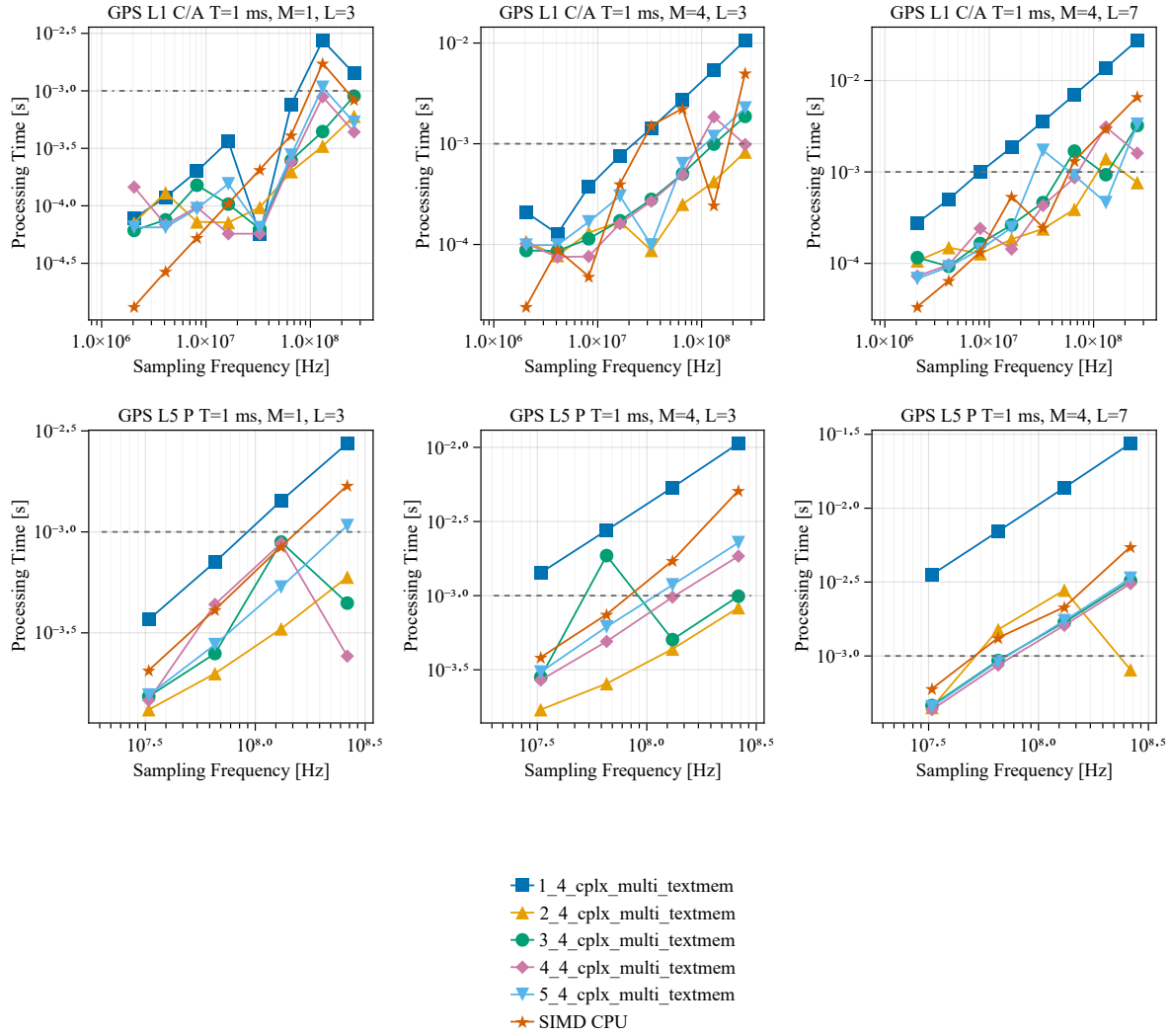


Fig. 8: Processing time of correlation of a 1ms GPS signal on Platform #2.
The dashed line indicates the real-time processing bound.