



SQUARE KILOMETRE ARRAY SOUTH AFRICA

INTERNSHIP TECHNICAL REPORT

GPU-accelerated Inverse Polyphase Filter bank

Author:

Benjamin HUGO

Department of Computer Science
University of Cape Town
bennahugo@aol.com

Advisor:

Ludwig Schwardt

REVISION 2

February 10, 2014

The results obtained in this report were made possible through the use of the ICTS High Performance HEX cluster of the University of Cape Town (UCT) `hex.uct.ac.za`. The author wishes to thank Andrew Lewis of UCT ICTS for the technical support he provided during this project.

Glossary

DFT	Discrete Fourier Transform. 1, 3–5, 7
FFT	Fast Fourier Transform is a fast algorithm for computing the DFT using the Cooley–Tukey algorithm amongst others. 7, 10
FIR	Finite Impulse Response. 4, 5, 8
FPGA	Field Programmable Gate Array. 6
lo	Local Oscillator. 4
PFB	Polyphase Filter Bank. 4, 5
PR	Perfect Reconstruction. 5, 6
VLBI	Very Long Baseline Interferometry. 4

Contents

1	Introduction	4
2	Analysis and design	5
2.1	Validation of the forward and inverse processes	5
2.2	Validation of the GPU version	6
2.3	Optimizations on the GPU-accelerated synthesis filter bank	7
3	Results and discussion	8
3.1	Reconstruction quality	8
3.2	GPU Computational performance	8
4	Future work	10

List of Figures

1	DFT Leakage	5
2	Correlation between synthesis filter bank and input	6
3	GPU performance	9

List of Tables

2	Validation of the GPU version	7
3	Reconstruction quality	9

1 Introduction

At its core Very Long Baseline Interferometry (VLBI) is an interferometry process where the output from several radio antennae are combined, to form an equivalent output to a telescope of the size equal to the distance between the furthest two antennae in the VLBI array. These arrays spans over very large areas or entire continents. For a comprehensive overview of the technique the reader is referred to the detailed introduction by Middelberg et al. [2]. Several such extended arrays exist including the European VLBI network and the Very Long Baseline Array. The hope is to include the KAT-7 in future VLBI observations.

Traditionally the process required raw data to be dumped to a storage medium, say tape, and physically shipped to a correlator where the data from several telescopes could be combined. The new trend in VLBI is to perform real-time correlation between antennae using high-speed internet connections and is known as 'eVLBI'. Both the Australian Long Baseline Array and the European VLBI network have performed successful eVLBI observations.

Ultimately the problem being investigated (at least in part by this report) can be boiled down to converting data sent over the SPEAD protocol (employed internally by the KAT-7 array) to the VDIF format. This conversion process includes a necessary first step, where the current sampling rate of 800 MHz¹ is reduced to 128 MHz through *Digital Downconversion*. The basic operation involves the following three steps:

1. **Mixing.** Where the signal is shifted down to *baseband* (frequency 0 of the spectrum), by mixing the signal with a Local Oscillator (lo). The lo is simply a generated sine wave with its frequency equal to the lower end of the sub-spectrum to be extracted. Mixing is achieved through an element-wise multiplication of the original signal with the lo tone. In essence mixing is not a *linear* (refer to [3, ch. 5]) process. If f was a single channel in the frequency domain then mixing produces replicated channels at $f - f_{lo}$ and $f + f_{lo}$.
2. **Filtering.** In order to eliminate aliasing in the frequency domain, due to mixing and frequencies above the new sampling rate, a Finite Impulse Response (FIR) filter with its cutoff set at the rate $\frac{1}{2}f_{decimated}$ is used to comply with the Nyquist sampling theorem (see [3, ch. 3]).
3. **Interpolation and Decimation.** The reader is referred to <http://www.dspguru.com/dsp/faqs/multirate/basics> for an overview of the process.

There is, however, a further complication to deal with, before this downconversion process can begin: the KAT-7 beamformer produces a series of frequency spectra. If no filtering was applied to these spectra, the undo operation would only have involved applying the inverse Discrete Fourier Transform (DFT). However, the original voltage data went through a filter-bank operation, in a process known as the Polyphase Filter Bank (PFB). This method is also known as the Weighted Window Overlap method and is necessary for computing Short Time Discrete Fourier Transforms on subsections of very long signals. The reader is referred to https://casper.berkeley.edu/wiki/The_Polyphase_Filter_Bank_Technique for a detailed description on the forward PFB process.

This report first discusses the validation techniques we employed to test both the construction of the analysis and synthesis filter banks of the PFB and inverse PFB, along with the validation of our GPU code. We then provide an overview of our design and optimization techniques for the GPU version of the filter bank, before discussing the quality of the reconstructed signal and GPU performance.

¹According to the KAT-7 Data Interface document

2 Analysis and design

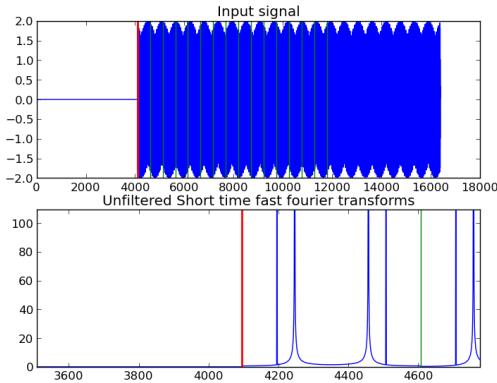
2.1 Validation of the forward and inverse processes

The forward PFB process can be thought of as a very basic analysis filter bank, whereas the inverse process can be thought of as a synthesis filter bank. The analysis filter bank uses a Hamming-windowed FIR with $P = 8$ banks of $N = 1024$ elements, with no up- or down-sampling stages². Our implementation focuses on 8-bit-sized samples produced by the KAT-7 beamformer³, but could easily be adapted to other sample-sizes, as well as prototype filter lengths.

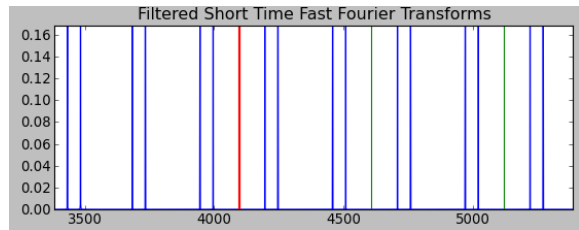
The key differences between the analysis and synthesis processes are:

1. If $H[n]$ is the prototype FIR filter and the analysis filter bank uses the subfilters $H_1[n], H_2[n] \dots H_P[n]$, each of length N , then the synthesis filter bank uses the subfilters $\bar{H}_1[N - n - 1], \bar{H}_2[N - n - 1] \dots \bar{H}_P[N - n - 1]$ ⁴.
2. The commutator is flipped around (as shown on page 10 of Zhou's report) and therefore the subfilters should be processed in reverse order.

We've implemented both the basic analysis and synthesis filter banks, as well as the necessary tools to analyse the performance of this basic filter bank construction. As shown in figure 1 the analysis filter bank successfully limits leakage of any frequencies that lie between neighboring bin centers of the DFT. This confirms that our basic analysis filter bank construction is correct.



(a) Unfiltered analysis filter bank



(b) Filtered analysis filter bank

Figure 1: Example of DFT leakage. It is clear that the unfiltered analysis filter bank suffers severely when containing frequencies even slightly off the bin centers of the DFT. Its filtered counterpart produces crisp binning of both centered and uncentered frequencies.

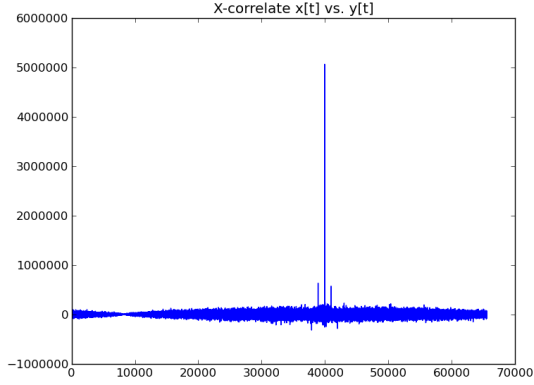
The synthesis filter bank construction was tested through cross correlation with the original input signal. A surefire way of determining if there is a correlation between the input and the output of the sythesis filter bank, is to use Gaussian noise as input signal. Perfect Reconstruction (PR) requires that the signal only varies in the following two ways:

1. The output signal may be shifted by some Δ time steps.

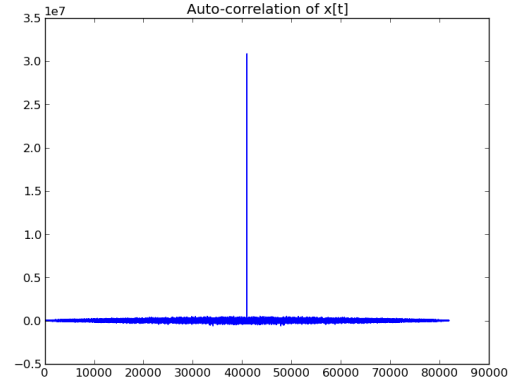
²According to information provided by Jason Manley

³According to the KAT-7 Data Interface document

⁴Findings published in an in-house final technical report titled 'A review of polyphase filter banks and their application'. Daniel Zhou. Air Force Research Laboratory, Information Directorate, Rome Research Site, Rome, New York



(a) Cross correlation between $x[t]$ and $y[t]$



(b) Auto correlation of $x[t]$

Figure 2: It is clear that even for Gaussian noise the synthesis filter bank achieves a limited reconstruction - the correlation between the reconstructed- and input signals is about 6 times weaker than the auto correlation of the input signal. The reader should also notice that it appears that $\Delta = -N(P - 1)$. This observation was made for trials with input signals of different lengths, but we have not rigorously proven this, and it should not be taken as a fact. Future expansion, testing and debugging should use the cross correlation to determine Δ .

2. The output signal may be scaled by a constant.

The signal shift may be found by the following method (where $x[t]$ is the input to the analysis filter and $y[t]$ is the output of the synthesis filter):

$$\Delta = \underset{t}{\operatorname{argmax}}(x[t] \star y[t]) \quad (1)$$

According to Parishwad Vaidyanathan [4] this basic filter bank cannot achieve PR and have to be modified in order to do so. This involves modifying both the analysis and synthesis filter banks. In our situation this is not plausible, since only the basic form of analysis filter bank is implemented on Field Programmable Gate Arrays (FPGA) in the current production environment. This is due to their simple construction and *fast evaluation times* [4]. Vaidyanathan goes on to state, however, that the synthesis filter bank requires subfilters of much greater length to obtain reasonably good reconstruction and is therefore computationally prohibiting. That said, as the reader can see from figure 2, $x[t] \star y[t]$ shows that our basic synthesis filter bank achieves limited reconstruction.

2.2 Validation of the GPU version

We will, for now, only focus on the most basic filter bank construction to determine what throughput rates can be achieved, as well as the quality of the reconstruction. The full source code along with test and analysis suite should be included with this report. With this suite the C++ CUDA GPU code can be validated against a serial python implementation. The reader should have GCC 4.6 installed, along with the NVIDIA NVCC compiler toolkit 5.0 redistributable. Although the code should work on any platform, we have only tested this under recent GNU/Linux kernels. Another requirement is that the user have a GPU of Compute Capability ≥ 2.0 installed on their system. In order to run the test suite the user should have Python 2.7 along with the SciLab suite of packages installed (including NumPy, SciPy and Matplotlib).

In order to validate that the GPU version is in fact producing valid results we've tested it against the

Metric	Python Synthesis Bank	GPU Synthesis Bank
\bar{x}	-0.000307	-0.000307
σ	2.969016	2.969015
Mean Squared Error	0.000001	

Table 2: Results from validating the GPU version against the serial python version of the synthesis filter bank

output of our Python synthesis filter bank. The Mean Squared Error is defined as:

$$MSE(input_1, input_2) := \frac{\sum_{k=0}^{L-1} (input_1[k] - input_2[k])^2}{L} \text{ where } ||input_1|| = ||input_2|| = L \quad (2)$$

By generating 0.976 GiB of Gaussian noise we determined the characteristics of the two output files, as shown in table 2

This subtle difference between the output files may be attributed to both floating-point rounding errors and the fact that we have discarded the last non-redundant sample of FFT output (explained shortly). Different FFT implementations may behave differently under this situation and may produce slightly different results. It should be noted that here the GPU has processed the file in several batches as we will explain next.

2.3 Optimizations on the GPU-accelerated synthesis filter bank

The GPU version has been constructed with the requirements of real-time operation in mind. Although it has not been integrated into a SPEAD receiver, the device simply has to be initialized once. The GPU will allocate all the memory it requires during this initialization process, after which the user can send batches, of any size, off to the GPU for processing. The user can send these batches on a rolling bases, as long as one batch is fully processed before the next batch is started. The GPU inverts all the DFTs in a combined operation and keeps track of the last $P \times N$ N -sized inverse DFTs. This is achieved by keeping a persistent buffer between batches. The only hard constraint placed on this batching process is that the size of the batch is an integral multiple of $\frac{N}{2}$ (the length of each analysis spectrum as mentioned earlier). At this point we wish to make the reader aware of subtlety in the output data received from the beamformers. Although a real-valued inverse DFT will require $\frac{N}{2} + 1$ non-redundant samples to perform a computation, the last frequency bin is discarded before we receive it. The author has assumed $X_{filtered}[f = \frac{N}{2} + 1] = 0$ in his construction of the analysis and synthesis filterbanks.

Care has been taken in optimizing the GPU code according to the architectural requirements of GPUs:

1. Memory copies to and from the device is done in batches to improve latency
2. Memory accesses have been coalesced, including necessary padding to align memory to the 128 byte memory boundaries of Compute Capability 2.0. This constant can be adjusted for newer architectures if necessary. At this point we should warn the reader that N must be a multiple of 128 for this coalescing scheme to function as intended.
3. Memory transfers to and from the GPU have been optimized by pinning host memory to primary memory so that it does not get paged out to disk by the memory controller of the host operating system.
4. DFTs are computed in batches using the NVIDIA CUFFT libraries using FFTs.
5. CUDA streams are used to run kernels and memory copies asynchronously whenever possible. This adds an additional level of parallelism to the solution.

6. We investigated using the cached texture memory of the GPU to store the prototype FIR filter coefficients. This did not provide greater throughputs compared to coalesced global memory accesses. We note for future reference that it is probably not worthwhile copying these filter taps to constant device memory, as this memory only provides better throughputs when multiple threads (usually either a half-warp or warp of 32 threads) access the same element simultaneously. This is not true in our situation, and using constant memory will result in serialization of the entire group of threads!
7. In compliance with the limitations of compute capability 2.0⁵ devices we split the number of blocks being executed in each grid between the Dim_x and Dim_y dimensions of the grid.
8. As it stands the computations of the basic synthesis filter bank did not require the use of shared memory. If kernels are added at a latter stage care should be taken to avoid bank-conflicts between warps of threads. The reader is referred to GPU Gems 3 [1, ch. 3] for a short introduction to this architectural constraint.

The user should manually tweak the memory allocation to fulfill their needs. All these constants are defined in the header file `inv_pfb.h`. It should be pointed out that the card should be completely filled with data in order to achieve good occupancy and offset the comparatively slow memory transfers through the PCI express bus.

3 Results and discussion

3.1 Reconstruction quality

In order to asses the quality of reconstruction we used Gaussian noise as input to the analysis filter bank and measured the reconstruction factor. Our results corroborates the literature surrounding filterbank construction, but does show that we achieve limited reconstruction. That said, there is a reasonable correlation between the input and the output of the synthesis filter bank. This can be confirmed by visual inspection of the (time-shifted) output of the synthesis filter bank.

The reconstruction factor is calculated in the following manor ($x[t]$ is the input Gaussian noise and $y[t]$ is the output of the synthesis filter bank).

1. Assume $y[t] = c(x[t]c_{ax} + n[t]c_{an})$ where $c_{ax}^2 + c_{an}^2 = 1$, ie. the reconstructed signal will consist of the scaled sum of scaled input and scaled noise.
2. An estimate of c is given by $c^{est} = \frac{\sigma_{y[t]}}{\sigma_{x[t]}}$
3. An estimate of c_{ax} is given by $c_{ax}^{est} = \max(x[t] \star y[t]) \times (c^{est})^{-1} \times (\max(x[t] \star x[t]))^{-1}$
4. The reconstruction ratio is defined as $\frac{\sqrt{2}c_{ax}^{est}}{\sqrt{1+(c_{ax}^{est})^2}}$

Table 3 shows that, although the synthesis filterbank has degraded a significant portion of the signal, strong sources will still be observed. Faint sources will require a combination of longer observation and/or longer synthesis filters.

3.2 GPU Computational performance

We have recorded performance measurements on a Tesla M2090 and a single core (Python implementation) of a AMD Opteron 6274 clocked at 2200 MHz with 128 GiB primary memory, clocked at 1333 MHz.

⁵Compute Capability 2.0 accepts a maximum of 65536 blocks per grid dimension

Variable	Value
$ x[t] $	81920
$ y[t] $	$81920 - 2NP = 65536$
$\sigma_{x[t]}$	25.944
$\sigma_{y[t]}$	10.940
$\max(x[t] \star y[t])$	10577811
$\max(x[t] \star x[t])$	55151821
Recovery factor	0.59
Signal degradation	-2.32dB

Table 3: Reconstruction quality obtained by the basic synthesis bank

These timings include memory transfer costs between the host and the GPU. It would be considered an unfair comparison if they were neglected.

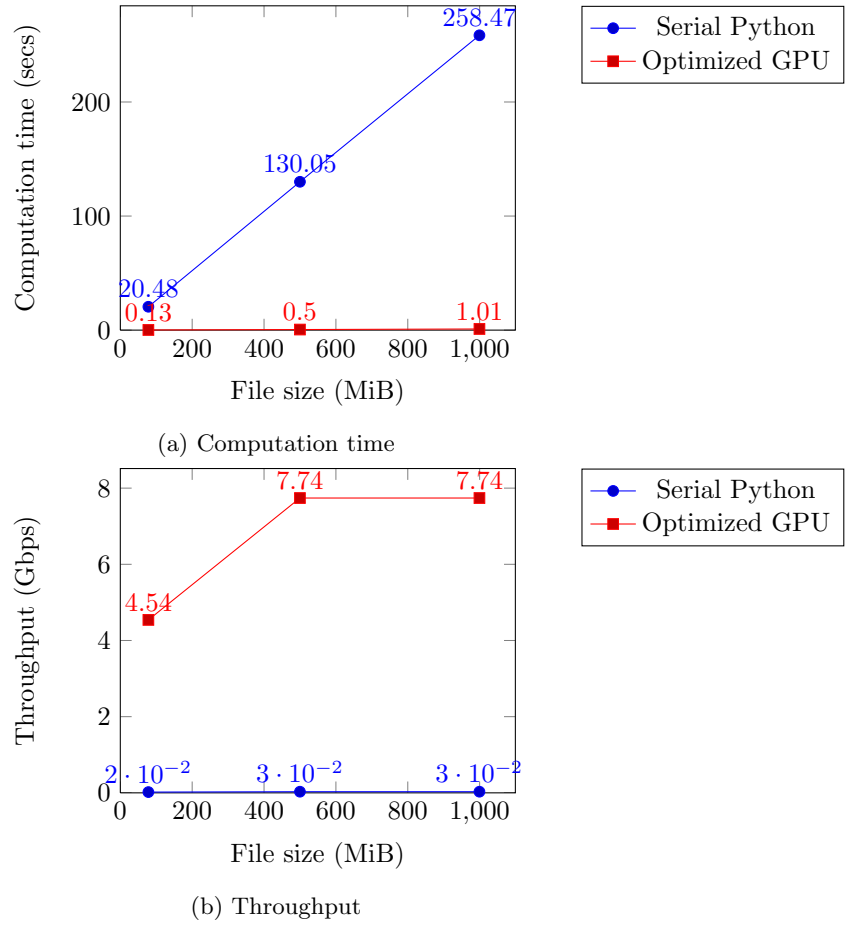


Figure 3: These results indicate that the GPU significantly outperforms its serial counterpart. Additionally, when the GPU version is executed on a much larger file, of almost 4 GiB, the throughput increases to 8.20 Gbps - more than a 273x speedup! These rates well exceeds the requirements for real-time synthesis. It does, however, show that the GPU memory should be saturated with data in order to offset the memory transfer between the host and device.

Although this speedup seems remarkably good the author wants to caution the reader that, ideally, the optimized GPU version should be compared to an parallel optimized CPU version. Not only can we easily parallelize this algorithm using OpenMP, but there is also opportunity here to exploit the full parallel

capabilities of modern CPUs, by making use of the special vectorized instruction sets like Intel/AMD Streaming SIMD Extensions or Intel Advanced Vector eXtensions. The latter can perform up to 8 32-bit floating-point operations per core simultaneously, by loading 8 elements into special 256-bit register at a time. This approach, coupled with higher CPU clock speeds may reduce this speedup significantly if executed on an 8 or 16 core server CPU, and by extension the latest Intel Xeon Phi co-processors.

4 Future work

There is a clear need for an investigation into extending the lengths of the subfilters in order to obtain better a signal reconstruction. The filterbank operation has complexity of $O(n)$. Therefore it is reasonable to assume that if the subfilters are expanded by a constant factor the computation time will also increase in a linear fashion. The author notes that this will require a slight adjustment in both the Python and GPU implementations where the number of FFT bins and subfilter lengths are not assumed to be equal, as per the basic construction.

References

- [1] Mark Harris, Shubhabrata Sengupta, and John D Owens. Parallel prefix sum (scan) with cuda. *GPU gems*, 3(39):851–876, 2007.
- [2] Enno Middelberg and Uwe Bach. High resolution radio astronomy using very long baseline interferometry. *Reports on Progress in Physics*, 71(6):066901, 2008.
- [3] Steven W Smith et al. The scientist and engineer’s guide to digital signal processing. 1997.
- [4] Parishwad P Vaidyanathan. Multirate digital filters, filter banks, polyphase networks, and applications: A tutorial. *Proceedings of the IEEE*, 78(1):56–93, 1990.