

**Automatic Mapping of Real Time Radio Astronomy Signal Processing
Pipelines onto Heterogeneous Clusters**

by

Terry Filiba

A dissertation submitted in partial satisfaction of the
requirements for the degree of
Doctor of Philosophy

in

Electrical Engineering and Computer Sciences

in the

Graduate Division

of the

University of California, Berkeley

Committee in charge:

Professor John Wawrzynek, Co-chair
Dan Werthimer, Co-chair
Professor Jan Rabaey
Associate Professor Aaron Parsons

Summer 2013

The dissertation of Terry Filiba, titled Automatic Mapping of Real Time Radio Astronomy Signal Processing Pipelines onto Heterogeneous Clusters, is approved:

Co-chair _____

Date _____

Co-chair _____

Date _____

Date _____

Date _____

University of California, Berkeley

**Automatic Mapping of Real Time Radio Astronomy Signal Processing
Pipelines onto Heterogeneous Clusters**

Copyright 2013
by
Terry Filiba

Abstract

Automatic Mapping of Real Time Radio Astronomy Signal Processing Pipelines onto
Heterogeneous Clusters

by

Terry Filiba

Doctor of Philosophy in Electrical Engineering and Computer Sciences

University of California, Berkeley

Professor John Wawrynek, Co-chair

Dan Werthimer, Co-chair

Traditional radio astronomy instrumentation relies on custom built designs, specialized for each science application. Traditional high performance computing (HPC) uses general purpose clusters and tools to parallelize the each algorithm across a cluster. In real time radio astronomy processing, a simple CPU/GPU cluster alone is insufficient to process the data. Instead, digitizing and initial processing of high bandwidth data received from a single antenna is often done in FPGA as it is infeasible to get the data into a single server.

I propose to develop a universal architecture where each problem is partitioned across a heterogeneous cluster, taking advantage of the strengths different technologies have to offer. I propose we take an HPC approach to instrument development with a heterogeneous cluster that has both FPGAs and traditional servers. This cluster can be reprogrammed as necessary in the same way an HPC cluster is used to run many different applications on the same hardware.

The challenge in this heterogeneous of approach is partitioning the problem. Normal HPC uses a homogeneous cluster where the nodes are interchangeable. In a heterogeneous cluster, there is an additional issue of determining how to partition the problem across different types of hardware. I propose to design a tool that automatically determines how to partition instruments for radio astronomy. In order to do this work, I will need to model the platforms and based on a description of the final instrument, generate a processing pipeline. The partitioning needs to be done using a variety of techniques to assess the hardware. A static model of the hardware is useful to determine the amount of processing available in different types of hardware. Dynamic benchmarking would also be needed to deal with varying server architectures and determine how much processing and bandwidth the cpu/gpu servers can handle. Finally, to capture any overlooked subtleties or deal with things the tools cannot handle, the user will be able to input hints as to how the instrument should be generated.

The development of this tool will be driven by 2 instruments. First, the design of the GBT spectrometer, a spectrometer designed to support many different modes using the same cluster. By using the tool to design this spectrometer, additional modes can be easily added and if the cluster is expanded the each mode can be redesigned to do additional processing that takes advantage of the extra hardware. I will also be working on the LEDA correlator. This is a low bandwidth, large N correlator, which is an ideal application for heterogeneous clusters.

We can assess the performance of this automatic partitioning tool in a number of ways. First, this tool should significantly reduce NRE and time to science. By automatically generating the instrument the need for engineers who understand both science goals and programming is removed. However, this benefit should not come with a large increase in cost. The instruments produced by this tool will be compared to optimized implementations with the same parameters on the basis of hardware utilization and power consumption.

i

For my parents

Contents

| | |
|----------------------------------------------------------------------------------|-----------|
| Contents | ii |
| List of Figures | iv |
| List of Tables | vi |
| 1 Introduction | 1 |
| 1.1 Motivation | 1 |
| 1.2 Challenges | 2 |
| 1.3 Optimal Rearrangement of Cluster-based Astronomy Signal Processing | 3 |
| 2 Real Time Radio Astronomy Algorithms | 5 |
| 2.1 Spectroscopy | 6 |
| 2.2 Pulsar Processing | 7 |
| 2.3 Beamforming | 8 |
| 2.4 Correlation | 9 |
| 3 Related Work | 12 |
| 3.1 Radio Astronomy | 12 |
| 3.2 Tuning | 28 |
| 4 High Level Toolflow | 30 |
| 4.1 ORCAS Goals | 30 |
| 4.2 Instrument Definition | 32 |
| 4.3 Dataflow Model | 33 |
| 4.4 Mapping | 38 |
| 4.5 Code Generation | 41 |
| 5 Algorithm Partitioning | 42 |
| 5.1 Variables | 42 |
| 5.2 Constraints | 42 |
| 5.3 ILP Implementation | 46 |
| 5.4 Performance Modeling | 47 |

| | | |
|----------|---------------------------------------------------|-----------|
| 5.5 | Cost Modeling | 50 |
| 5.6 | Optimization | 50 |
| 5.7 | Final Mapping | 53 |
| 6 | Analysis | 55 |
| 6.1 | Spectrometer Case Study | 55 |
| 6.2 | High Resolution Spectrometer Case Study | 57 |
| 6.3 | FX Correlator Case Study | 59 |
| 7 | Conclusions | 62 |
| 7.1 | Future Work | 63 |
| | Bibliography | 64 |

List of Figures

| | | |
|------|-------------------------------------------------------------------------|----|
| 1.1 | The VLA and Arecibo Telescopes | 1 |
| 1.2 | Technology Spectrum | 2 |
| 2.1 | TODO | 7 |
| 2.2 | Telescope Locations in the Very Long Baseline Array | 11 |
| 3.1 | xGPU Performance | 13 |
| 3.2 | Map of CASPER Collaborators | 14 |
| 3.3 | ROACH Board | 16 |
| 3.4 | Adder Tree Simulink Simulink Diagram | 17 |
| 3.5 | Adder Tree Simulink Mask Script | 18 |
| 3.6 | CASPER FFT Library | 19 |
| 3.7 | CASPER FFT Options Menu | 20 |
| 3.8 | A comparison of FFT and PFB response | 20 |
| 3.9 | CASPER Library IO Blocks | 21 |
| 3.10 | SERENDIP V.v Block Diagram | 22 |
| 3.11 | Arecibo ALFA Feed | 22 |
| 3.12 | SERENDIP V.v Hardware installed at Arecibo Observatory | 23 |
| 3.13 | CASPER Correlator Block Diagram | 24 |
| 3.14 | PASP Dataflow | 24 |
| 3.15 | PASP Interface | 25 |
| 3.16 | PASP Low Level Implementation | 26 |
| 3.17 | TODO | 27 |
| 3.18 | Example high level spectrometer architecture | 28 |
| 4.1 | The Four Stage ORCAS Toolflow | 31 |
| 4.2 | ORCAS Toolflow Instrument Definition | 32 |
| 4.3 | The ORCAS Toolflow: Dataflow Model | 34 |
| 4.4 | Example FX Correlator Dataflow Model Demonstrating Blocktypes | 35 |
| 4.5 | TODO | 36 |
| 4.6 | TODO | 36 |
| 4.7 | The ORCAS Toolflow: Mapped Dataflow Representation | 38 |

| | | |
|-----|-------------------------------------------------------------------------------------------------------------------------------|----|
| 4.8 | Two potential mappings for the FX Correlator | 40 |
| 5.1 | Full-Crossbar Interconnect | 44 |
| 5.2 | PFB FIR and FFT benchmark data on the Virtex V SX95T | 48 |
| 5.3 | PFB FIR and FFT benchmark data on the GTX 580 | 49 |
| 5.4 | TODO | 52 |
| 5.5 | ORCAS Output | 54 |
| 6.1 | General Spectrometer Dataflow Model | 56 |
| 6.2 | Example High Resolution Spectrometer Dataflow Model | 58 |
| 6.3 | FX Correlator F-Engine Model | 59 |
| 6.4 | FX Correlator X-Engine Model | 60 |
| 6.5 | Example FX Correlator Dataflow Model | 60 |
| 7.1 | A potential architecture for multiple scientific instruments simultaneously processing data from the same telescope | 62 |

List of Tables

| | | |
|-----|-------------------------------------------------------------------------|----|
| 3.1 | CASPER ADC Boards | 15 |
| 5.1 | Comparative Resource Utilization of a 32k Channel 800 MHz FFT | 50 |
| 5.2 | Monetary and Power Costs for Common CASPER Platforms | 50 |
| 5.3 | GPU Board Costs | 50 |

Acknowledgments

TODO

Chapter 1

Introduction

Radio astronomers are trying to answer a very diverse set of questions.

1.1 Motivation



Figure 1.1: The VLA and Arecibo Telescopes

Traditionally, observatories designed custom instruments that would run on one telescope and solve one problem. This custom approach was the only way to get the requisite processing power to analyze the radio signals, but it resulted in costly designs, because the boards, backplanes, chips, protocols, and software all needed to be designed from scratch. To make matters worse, this approach resulted in a very long design cycle, requiring 5-10

years of development before an instrument could be deployed at a telescope and the time the instrument was released, the hardware would be out of date.

Due to their custom implementations, these instruments also lacked flexibility. Each instrument was designed specifically for a single purpose. A hardware upgrade or algorithm modification would require a complete redesign of the instrument, and another long design cycle.

While these older designs needed to trade off flexibility for performance, newer technology can offer both performance and flexibility. Programmable devices such as FPGAs, GPUs and even CPUs can provide enough processing power to keep up with the data from many new telescopes. These devices make it easy to reprogram existing hardware to support newer algorithms, and, since they are programmed using portable languages, provide a quick path to upgrade hardware without redesigning the entire instrument.

With a huge range of technology available, choosing what hardware to use to build an instrument is a challenge.

1.2 Challenges

Scientific Challenges

Technological Challenges

When building a large instrument, there is a wide variety of technology to choose from. Unlike many applications, where the goal is to provide the fastest possible implementation, in real time radio astronomy instrumentation there is a performance target. Understanding the tradeoffs between different implementations is key to cost-effective design.

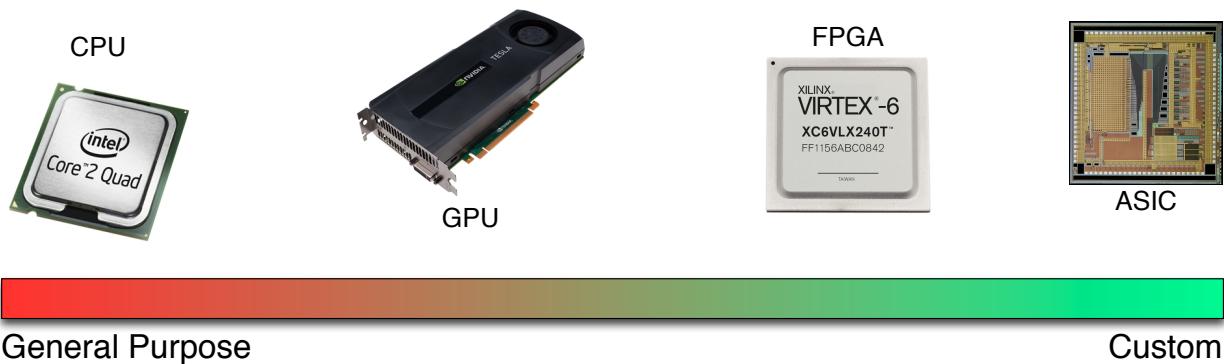


Figure 1.2: Technology Spectrum

An instrument designed typically has four choices of hardware for their algorithm implementations. They can use CPUs, GPUs, FPGAs, and ASICs. Figure 1.2 shows each of these as a spectrum from general purpose to custom. The CPUs and, to a lesser extend

GPUs, provide a very general purpose design experience, while FPGAs and ASICs require custom designs. This spectrum also can be used to generalize a number of other properties of the presented hardware. Platforms further right require higher design times and design complexity and less flexibility in the final design. And because the chips will be deployed in relatively low volume, the platforms on the right will also cost more money. The tradeoff is in power and performance. The benefit that comes with the platforms on the left is higher performance and lower power consumption.

To better understand these tradeoffs, I compare the two platforms in the center of the spectrum, the GPU and FPGA. NVIDIA GPUs can be programmed in CUDA. Its C-like structure makes it easy for CPU programmers to pick up and provides a lot of flexibility. CUDA allows the programmer to use conditional and Iterative programming the same way they would use C, easing the transition into GPU programming. FPGAs are programmed using a hardware description language, or HDL, which is less flexible and harder to learn than CUDA. In the FPGA computing model the data is streaming and everything is happening at the same time, requiring specialized programming. On the other hand, FPGAs offer an order of magnitude improvement in performance and power consumption. Each GPU requires hundreds of watts of power while an FPGA can be powered with tens of watts and a GPU can only process hundreds of megahertz of antenna data while an FPGA is capable of processing multiple gigahertz.

Although it may be easy to enumerate the differences between these platforms, understanding which is best is more difficult. In most cases, a heterogeneous approach is best. The best mix of platforms will ultimately depend on the desired instrument, available technology, as well as the price and power restrictions. This means that a brand new implementation must be designed for every new instrument.

1.3 Optimal Rearrangement of Cluster-based Astronomy Signal Processing

This dissertation explores an automated approach to instrument design using a tool called Optimal Rearrangement of Cluster-based Astronomy Signal Processing or ORCAS. This tool was designed to provide cost-optimal mappings in a short amount of time. ORCAS allows the developer to define the instrument at a high level and a cost function defining some cost the developer is aiming to minimize. The cost function can represent something as simple as the price of the instrument or can be used to represent more complex parameters of the instrument such as design time.

ORCAS assesses the performance of different types of technologies in two ways. First, by using benchmarks of existing implementations of instrument building blocks such as FFTs and FIR filters we can directly assess the performance. Second, if a benchmark is not available, the tool can use a performance model instead, giving an estimate of how the block will perform. This makes it easy to keep up with improving technology and library

performance without rewriting the entire tool every time a new library version or board is released.

We evaluate ORCAS by benchmarking existing libraries, such as CUFFT and the CASPER library, and use those benchmarks to produce mappings for 3 types of instruments. Each mapping is compared to existing implementations to understand how this automated approach compares to the old approach of hand optimized mapping.

The remainder of the dissertation is organized as follows. Chapter 2 provides a more in-depth look at the algorithms commonly used in radio astronomy and their applications. Chapter 3 describes related work, done by myself and others, that preceded this work. Chapter 4 presents a high level description of the tool and explains how the tool goes from a description of the instrument to a fully mapped algorithm. The algorithm used to map the instrument is fully specified in chapter 5. Chapter 6 describes three instrument case studies and compares them to existing instruments. And finally, in Chapter 7, I present my conclusions and some opportunities for future work.

Chapter 2

Real Time Radio Astronomy Algorithms

Radio astronomy simply refers to the type of science that can be done by observing astronomical objects at radio wavelengths, rather than a specific scientific goal. There is a huge variety of different experiments, such as studying the formation and structure of galaxies stars and black holes and searching for gravity waves, traces of the first stars (TODO: ref PAPER/LEDA), or aliens (TODO: ref SETI). But, despite this variety, the small number of algorithms detailed in this chapter serve as the first step in processing the data for many such projects.

Radio telescopes produce very high amounts of data (TODO: add estimates, how much data?). The reason for this high influx of data is twofold. First, to enable new science, new radio antenna observe increasingly higher bandwidths. Second, to sate the need for larger collecting area, rather than designing a single large dish, many new telescopes are being designed as antenna arrays, where the data from multiple antennas is combined to act as a single large dish. While it may be cheaper and easier to get a larger collecting area using small dishes rather than a single large dish, this adds additional complexity processing the data coming off the telescope. Rather than processing a single stream of data, now the instrument must process and combine multiple streams to make the array seem like a single large dish. As the size of the arrays and bandwidths for single dishes simultaneously increase, the data produced cannot be feasibly recorded in real-time. To cope with the progress in science and antenna technology, there is a constant need for new systems to process, rather than record, this data in real time. Each of these instruments begin processing the data immediately after it is digitized, and need to reduce the data without loosing scientific information. Once the data is partially processed, and reduced to a manageable bandwidth, it can be stored and processed further offline.

There are a small number of real time algorithms commonly used to reduce the data. In this work, we specifically focus on spectroscopy, pulsar processing, beamforming and correlation. Spectroscopy and pulsar processing are both spectral methods of analyzing data from a single beam. They both can be used on antenna arrays but would either need

to treat the array as separate dishes rather than one large dish, or the data from the dishes would need to be combined into a single beam, which can be done using the third algorithm on the list, beamforming. The last two techniques, beamforming and correlation, are both ways of combining data from multiple antennas. Beamforming combines the data into a single beam by delaying and summing the data from each antenna. Correlation doesn't aim to form a single beam, but instead is the first step towards getting an image of the sky.

2.1 Spectroscopy

A spectrometer is simply an instrument that produces an integrated, or averaged, frequency domain spectrum from a time domain signal. A real time spectrometer works by constantly computing a spectrum over short windows of data (channelization), then each channel is summed for a predetermined amount of time to compute the average power in that channel (accumulation). Figure x shows a block diagram for a simple spectrometer design. After digitization, there is the channelization step, where the signal is processed by a digital filter bank, and then the channels are accumulated.

High resolution spectroscopy

Increasing resolution often requires an increase in complexity in the spectrometer design. Once the number of required channels is sufficiently high, it becomes infeasible to compute the spectrum using a single filter bank. To cope with this, the channelization is done in two steps as shown in figure (TODO: add hi res spectrometer block diagram). In the first step, the signal is divided into coarse channels using a filter bank. At this point the channels are much wider than intended and can't be accumulated yet. After coarse channelization, the spectrometer treats the data from a single channel as time domain data and passes it through a filter bank again. This step breaks up the wide channel into a number of smaller channels. At this point the data can be accumulated, since it has the desired resolution.

Applications

This high resolution spectroscopy technique is used in the SERENDIP V.v (Search for Extraterrestrial Radio Emissions from Nearby Developed Intelligent Populations) project. This project is part of the SETI (Search for Extraterrestrial Intelligence) effort to detect extraterrestrial intelligence. The SERENDIP V.v project is a commensal survey at the Arecibo observatory, meaning anytime the ALFA receiver is used for any observation, the SERENDIP spectrometer will also process and record data.

The project is focused on detecting strong narrow band radio signals, requiring a high resolution spectrometer installed at the observing telescope to analyze the data. The spectrometer should resolve channels of less than 2 Hz, so that natural astronomical signals that typically have a wider bandwidth can easily be distinguished from narrower, possibly

extraterrestrial, signals. The SERENDIP V.v spectrometer meets this by providing 128 million channels across 200MHz of bandwidth, for an resolution of 1.5 Hz per channel. Since it would be infeasible to channelize a 200MHz into 128 million channels using a single filter bank, the SERENDIP V.v spectrometer uses the high resolution architecture described in figure x.

2.2 Pulsar Processing

A pulsar processor is an instrument designed specifically to observe transient events, such as pulsars. A pulsar is a rotating neutron star that emits an electromagnetic beam. When the beam sweeps past Earth, due to the rotation of the star, it is observed as a pulse of wideband noise. As the pulse travels through the interstellar medium (the matter filling interstellar space), the pulsar gets dispersed, meaning the low frequencies arrive before high frequencies, despite the fact that they were emitted at the same time.

A spectrometer, as described in Section 2.1, that accumulates the spectrum would smear the pulse, so there will need to be a few adjustments to the spectroscopy algorithm to make it suitable for processing transient events. In the case of pulsars, the algorithm starts with a high-resolution spectrometer without an accumulator. Instead, the algorithm becomes specialized to detect this type of quickly occurring event.

The high resolution data is then sent to a process called dedispersion, which undoes the dispersion caused by the ISM, realigning the pulse. There are 2 techniques to do this. First, the pulse can be dedispersed by shifting the frequency channels by different amounts to compensate for the different de-

lays as shown in figure x, in a process called *incoherent dedispersion*. This process can't be used to reconstruct the original pulse, but due to its relatively low compute cost, is a useful algorithm to search for new pulsars. The second technique, *coherent dedispersion*, models the effect of the ISM as a convolving filter. To remove this effect, the signal is deconvolved with the model. This is more compute intensive than incoherent dedispersion, but can recover the original pulse.

After dedispersion, there is still a lot of data and the pulse has very low SNR. The next step in processing is *folding*, or adds together many pulses, reducing the amount of data and improving the SNR. At this point, the data has been significantly reduced and can be recorded.

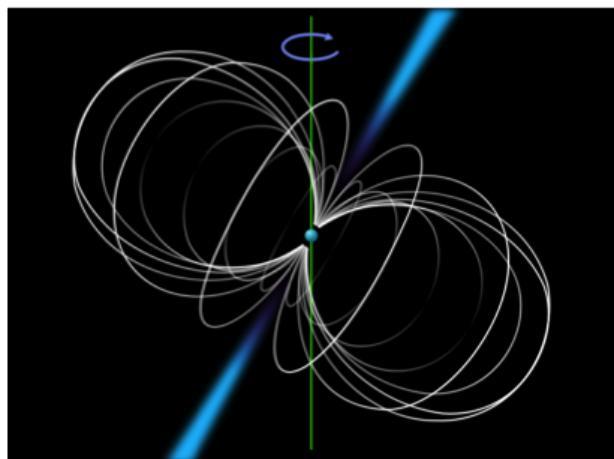


Figure 2.1: TODO

Applications

Some pulsars serve as very accurate astronomical time keepers. There are millisecond pulsars with extremely stable periods, allowing any perturbations in the observed pulse period to be attributed to some external effect. This makes pulsars extremely useful for conducting relativity experiments. One such example is the North American Nanohertz Observatory for Gravitational Waves, or NANOGrav, Project, which uses pulsars in an attempt to make the first detection of gravitational waves [7]. The project will try to detect gravitational waves by observing an array of pulsars, measuring the effect of the waves passing between Earth and the pulsars as changes in observed pulse arrival times.

[1]

2.3 Beamforming

Beamforming is a technique for combining data from an array of antennas. The beamformer combines the data from multiple antennas into a single beam pointed at a single point in the sky. This is achieved by delaying the signal from each antenna by a different amount and then summing the delayed signals. As illustrated in figure x, the signal from the intended source will not arrive at all the antennas at the same time. The delay compensates for the disparity between the arrival times, and once the signals are summed it creates constructive interference in the direction of the source, and destructive interference in other directions. These delays can be changed to point the beam at a different source or to track a single source moving across the sky.

BF Beamforming

Since the digitized signal from a radio telescopes receiver is discrete, and the amount of delay might not be an integer multiple of the sample period s , the delay d is applied to the signal, $f[n]$, in 2 steps. The delay in clock cycles is represented as d/s and broken up into its integer and fractional parts. First, the integer part is applied as a coarse delay, $n_f = \lfloor d/s \rfloor$. This can simply be implemented by buffering the signal. Applying the fractional delay $p \bmod s$ is more complex. Since there is no observed data at that time, the signal fractional samples are calculated by convolving the signal with an interpolation filter b . Applying these 2 steps, results in a new delayed signal $f[n] = f[n - n_f] * b_{fi}$

In practice, it is common to calculate the spectrum of the beamformed signal. A typical 2 element beamformer, with discrete input signals $f[n]$ and $g[n]$ is trying to calculate the spectrum of beam i , H_i , using coarse delays n_{fi} n_{gi} and delay filters b_{fi} and b_{gi} , as follows:

$$H_i(x) = FT(f[n - n_{fi}] * b_{fi} + g[n - n_{gi}] * b_{gi})$$

We describe this as BF beamforming because the delay operation (B), happens before the FT operation (F).

FB Beamforming

In the case where many beams are required, each beam needs a different FIR filter for every antenna, and a separate Fourier transform. This is called FB beamforming because the fourier transform operations (F) occur before the delay (B). Using linearity of the Fourier transform and the convolution theorem, the computation can be rearranged as follows:

$$H_i[x] = F[x] \cdot B'_{fi} + G[x] \cdot B'_{gi}$$

Where $F[x]$ and $G[x]$ are the fourier transforms of the signals f and g . The signal delay becomes a phase shift in frequency domain that is represented by B'_{fi} and B'_{gi} . In this case, there is a FT for each antenna, rather than one FT per beam. When forming multiple beams, there is no need to recalculate the Fourier transform of the signals. So, as the number of beams increases, it can be advantageous to use this algorithm rather than the BF algorithm described in the previous section.

Applications

[9]

[9]

2.4 Correlation

Aperture synthesis is another technique to combine the data from many antennas. The goal is to form an image of the sky, using 2 steps, correlation followed by imaging. In the correlation step, the cross-correlation of each pair of antennas is calculated. Once the cross-correlation is calculated, it is possible to accumulate the data, greatly reducing the amount of data that needs to be processed in the imaging step.

[27]

The uv plane represents the Fourier transform of the 2-dimensional sky image. The cross-correlation of an antenna pair represents a point in the uv plane, called a *visibility*. Since the visibilities are not evenly or continuously sampled, the imager must interpolate points on the uv plane in an evenly spaced grid so the FFT algorithm, which relies on even spacings, can be applied to the data. The imager must also account for the fact that the response of the telescope distorts the image, undoing the effect using iterative algorithms like CLEAN or Maximum Entropy. Once this is done, a two dimensional inverse Fourier transform can be applied to recover the sky image.

Typically, large telescope arrays do correlation in real time and finish the imaging offline, but there are a few notable exceptions. One of these is the VLBA, or Very Long Baseline Array. When designing antenna arrays, one important parameter is the distance between the antennas, or *baselines*. Arrays with longer baselines provide images with better angular resolution. The VLBA is an extremely high resolution array, achieving this resolution by

using telescopes on opposite ends of the Earth. The distance between the antennas, while useful for science, creates a logistical issue for any real time processing that depends on data from different antennas. Instead of correlating in real time, the VLBA records the digitized data directly from the telescopes at each site, without any reduction. Later, the recorded data is flown to a central location, where it gets correlated.

The cross-correlation of two signals, $f[n]$ and $g[n]$ is defined as:

$$h[n] = f[n] \star g[n] = \sum_{i=0}^m f^*[i]g[n+i]$$

Like in beamforming, it is often desirable to calculate a spectrum, so the correlation step typically also calculates the spectrum of the cross-correlations. So, the final result of the correlator produces the spectrum of the visibilities, $H[x]$.

$$H[x] = FT(f[n] \star g[n])$$

XF Correlation

An XF, or lag, correlator calculates the spectrum by calculating the cross-correlation first (X), followed by a Fourier transform (F).

The calculation of the cross-correlation is very similar to a convolution, in fact the cross-correlation of 2 signals can be re-expressed as a convolution.

For an telescope with n antennas, this algorithm requires $O(n^2)$ cross-correlations, followed by $O(n^2)$ Fourier transform operations.

FX Correlation

[3]

The cross-correlation theorem further extends the parallel between correlation and convolution, relating the Fourier transform of the cross-correlation of two signals to the Fourier transforms of the original signals.

$$\text{signal} = \text{otherstuff}$$

$$H_i[x] = F^*[x] \cdot G[x]$$

Applications

The Very Long Baseline Interferometry or VLBI technique uses an array of antennas that are extremely far apart to achieve very high angular resolution

Event Horizon Telescope ([9])

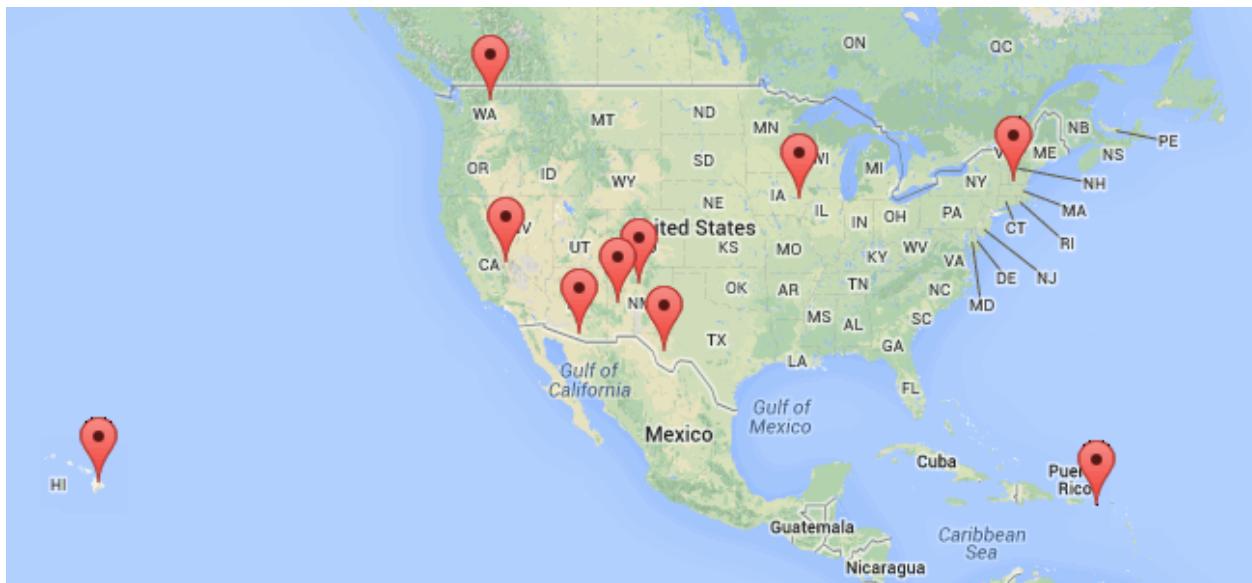


Figure 2.2: Telescope Locations in the Very Long Baseline Array

Chapter 3

Related Work

3.1 Radio Astronomy

The need for high bandwidth processing manifests in many different radio astronomy applications. Keeping up with increasing computation demands has often resulted in the specialized design of spectrometers.

Distributed FX Correlator (DiFX)

The DiFX Correlator is a scalable software implementation of an FX Correlator [6]. DiFX was designed as a software correlator that targets CPUs in order to maintain flexibility in the design. The DiFX correlator was originally developed to do VLBI (Very Long Baseline Interferometry).

Unfortunately, the flexibility of the software implementation comes at a high cost. Nearly 100 nodes are required to cross correlate 64 MHz of bandwidth from 10 antennas.

LOFAR

Like DiFX, the correlator for Low-Frequency Array for Radio Astronomy (or LOFAR) was also built using CPUs [29]. This project used a Blue Gene BG/P to implement a 64 station correlator. The implementation got very high performance out of the cluster, 96% of the peak, but required an entire cluster and required finely tuned code to achieve that performance.

xGPU

xGPU is a CUDA package that implements the cross-correlation

[4]

[30]

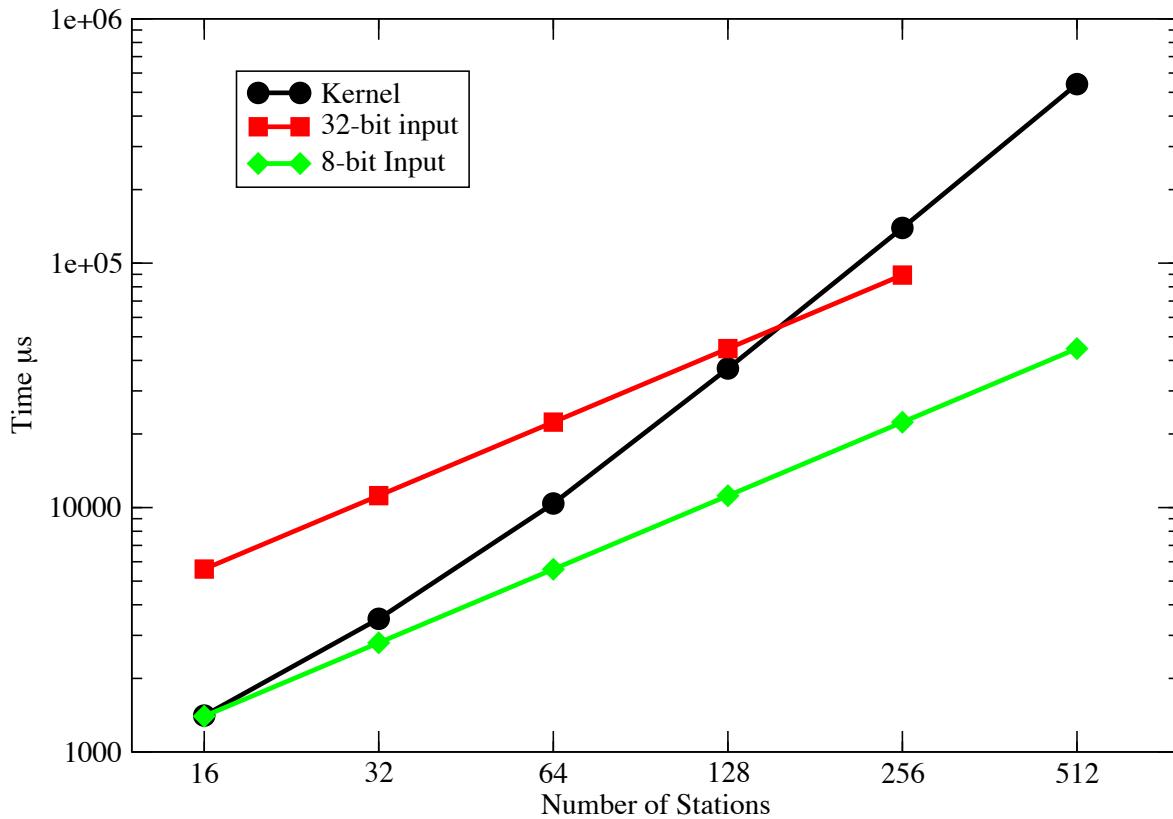


Figure 3.1: xGPU Performance (reprinted from Clark, La Plante, and Greenhill [4])

CASPER

At the Collaboration for Astronomy Signal Processing and Electronics Research (CASPER), we are developing common libraries and hardware that can be used by radio astronomers developing instrumentation. The mission statement, available on the group website casper.berkeley.edu, concisely summarizes the goals of the group:

The primary goal of CASPER is to streamline and simplify the design flow of radio astronomy instrumentation by promoting design reuse through the development of platform-independent, open-source hardware and software. Our aim is to couple the real-time streaming performance of application-specific hardware with the design simplicity of general-purpose software. By providing parameterized, platform-independent "gateware" libraries that run on reconfigurable, modular hardware building blocks, we abstract away low-level implementation details and allow astronomers to rapidly design and deploy new instruments.

The collaboration, started at UC Berkeley by Dan Werthimer, Don Backer, and Mel Wright, has spread to include astronomers all over the world, as shown in Figure 3.2. The group collaborates with many large digital engineering groups including groups from the Jet Propulsion Laboratory (JPL), the National Radio Astronomy Observatory (NRAO), the Arecibo Observatory, SKA South Africa, and the Giant Metrewave Radio Telescope (GMRT) in India. The diversity and far reach in this collaboration ensures that the tools continue to be general purpose and widely accessible.



Figure 3.2: Map of CASPER Collaborators

The CASPER FPGA libraries were developed to mitigate the need to redevelop common signal processing blocks for every new instrument [19]. The library provides parameterized blocks such as FFTs, digital down-converters, and FIR filters that are the necessary building blocks for a wide range of instruments. Coupled with open source FPGA boards, such as the ROACH (Reconfigurable Open Architecture Computing Hardware), the CASPER libraries provide a useful toolbox for radio astronomy instrumentation development. The follow sections describe each of these and how they have been used to create a number of different instruments.

CASPER Hardware

The CASPER group provides set of modular FPGA boards and ADCs that are designed specifically to deal with the high bandwidth requirements of real time radio astronomy signal processing. Since each FPGA boards meets the needs of the radio astronomy community as a whole rather than a single application, the group releases a small number of boards and

| ADC Name | Max sample rate | Streams | Bitwidth |
|-----------------------|-----------------|---------|----------|
| 64ADCx64-12 | 50 Msps | 64 | 12 bits |
| ADC4x250-8 or QuadADC | 250 Msps | 4 | 8 bits |
| KatADC | 1.5 Gsps | 1 | 8 bits |
| ADC1X2200-10 | 2.2 Gsps | 1 | 10 bits |
| ADC1x10000-4 | 10 Gsps | 1 | 4 bits |

Table 3.1: CASPER ADC Boards

typically only releases a new board to take advantage of improving technology. By releasing a handful of boards, the CASPER group ensures that users get proven and tested hardware for their new instruments. The CASPER library and software, discussed below, also make it easy to upgrade the hardware since a CASPER design easily be recompiled to work with a different board and the software interface and the signal processing model are standardized. And, since the boards communicate over a common set of industry standard protocols, they can be upgraded one by one, or all at once.

Each FPGA board implements a number of high speed interfaces to send and receive data. The Z-DOK+ connectors are primarily used to interface the boards to high speed ADCs and DACs. This common Z-DOK+ interconnect implemented on nearly all CASPER ADC boards allows the astronomer to choose an ADC to match their scientific goals. The diversity of available boards is shown in Table 3.1.

Each board can also send or receive data using ethernet. The boards are designed to communicate using common protocols, like 10 Gigabit Ethernet, so a board can be upgraded without modifying how it communicates with the rest of the cluster. The use of an industry standard protocol also makes communication to non-CASPER boards simple, allowing an FPGA to create UDP packets that eventually get received by a CPU or GPU-based server, making the CASPER boards an ideal component in a heterogeneous cluster. This makes it simple to design a cluster, and allows continuous upgrades as technology improves, as the signal processing model and communication model do not change between boards.

Currently, CASPER has designed the IBOB, or Interconnect Break-out Board, ROACH and ROACH2 boards. The IBOB was originally designed as an interface to digitize antenna data and send it to a server, but the powerful Virtex-II Pro make it very useful for channelizing data as well. The ROACH board is shown in Figure 3.3 with an iADC board connected via Z-DOK+ and an ethernet cable to get data off the board. This board is pin-compatible with two large FPGA chips, the Virtex 5 SX95T and LX100T, with the SX95T being preferred for its large number of DSP-48s. The ROACH has 2 Z-DOK+ connectors and supports up to 4 10-gigabit Ethernet links. The ROACH2 is an upgrade to the ROACH board to take advantage of the newer Xilinx chips. The Virtex 5 chips have been upgraded to the Virtex 6 SX475T, and this board supports up to 8 10-gigabit Ethernet links.

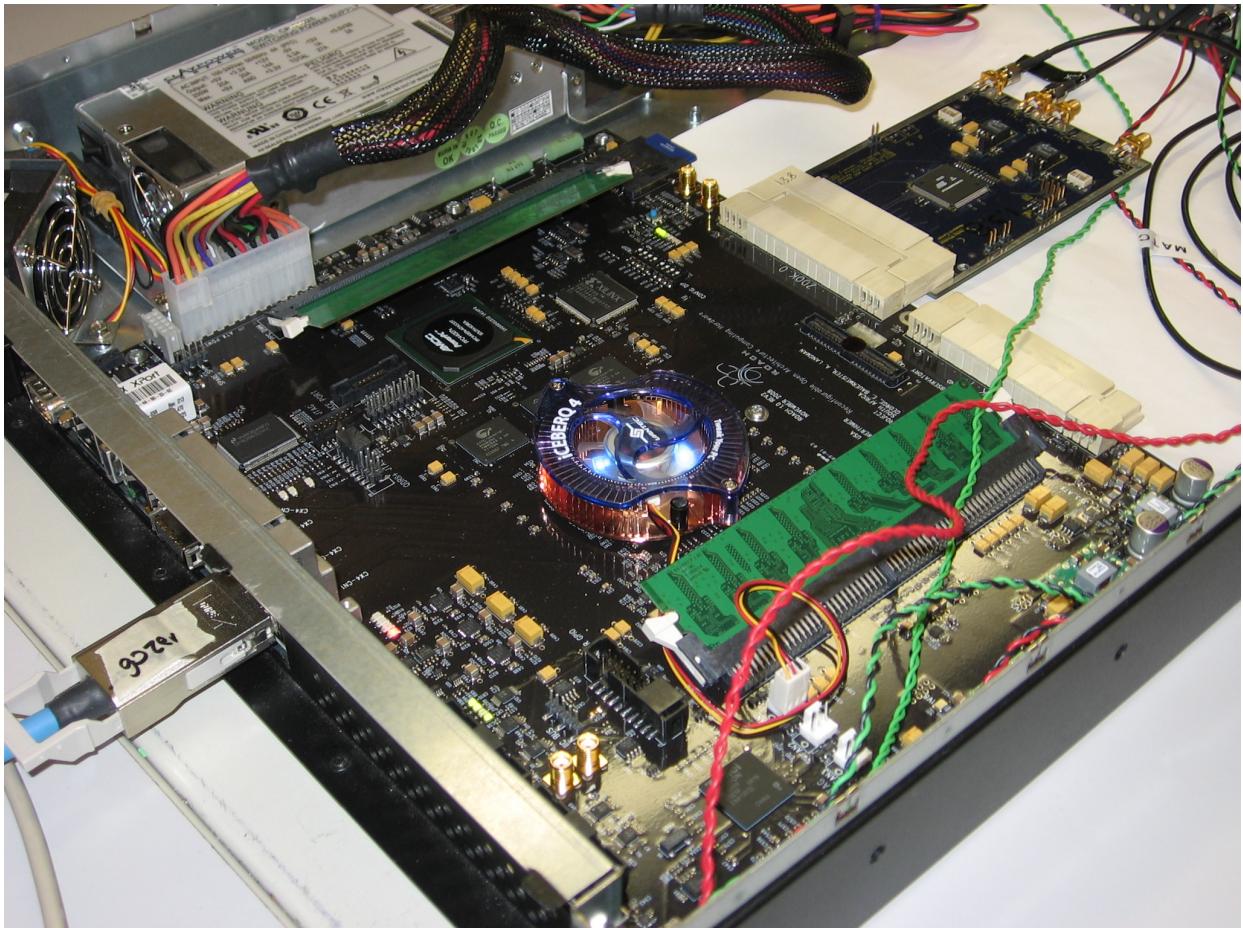


Figure 3.3: ROACH Board

CASPER Library

In addition to hardware, CASPER develops a DSP library that can be compiled into an FPGA bitstream. The library is implemented in Simulink, which allows for both simulation and, using Xilinx System Generator, compilation to FPGA code. The Simulink implementation is a key part of the success of this project. Using a visual programming language makes it easy to snap together blocks without getting into the low level details of FPGA programming.

The library blocks are built out of Xilinx System Generator primitive blocks. These primitives make it possible to retarget Simulink designs to different Xilinx FPGAs without changing the original implementation. Figure 3.4 shows a simple example of one of the library blocks, an adder tree. The adder tree is used to implement pipelined, parallel addition using a binary tree. The top left is the library with the adder tree block, and the bottom right is the implementation of the adder tree block with three inputs.

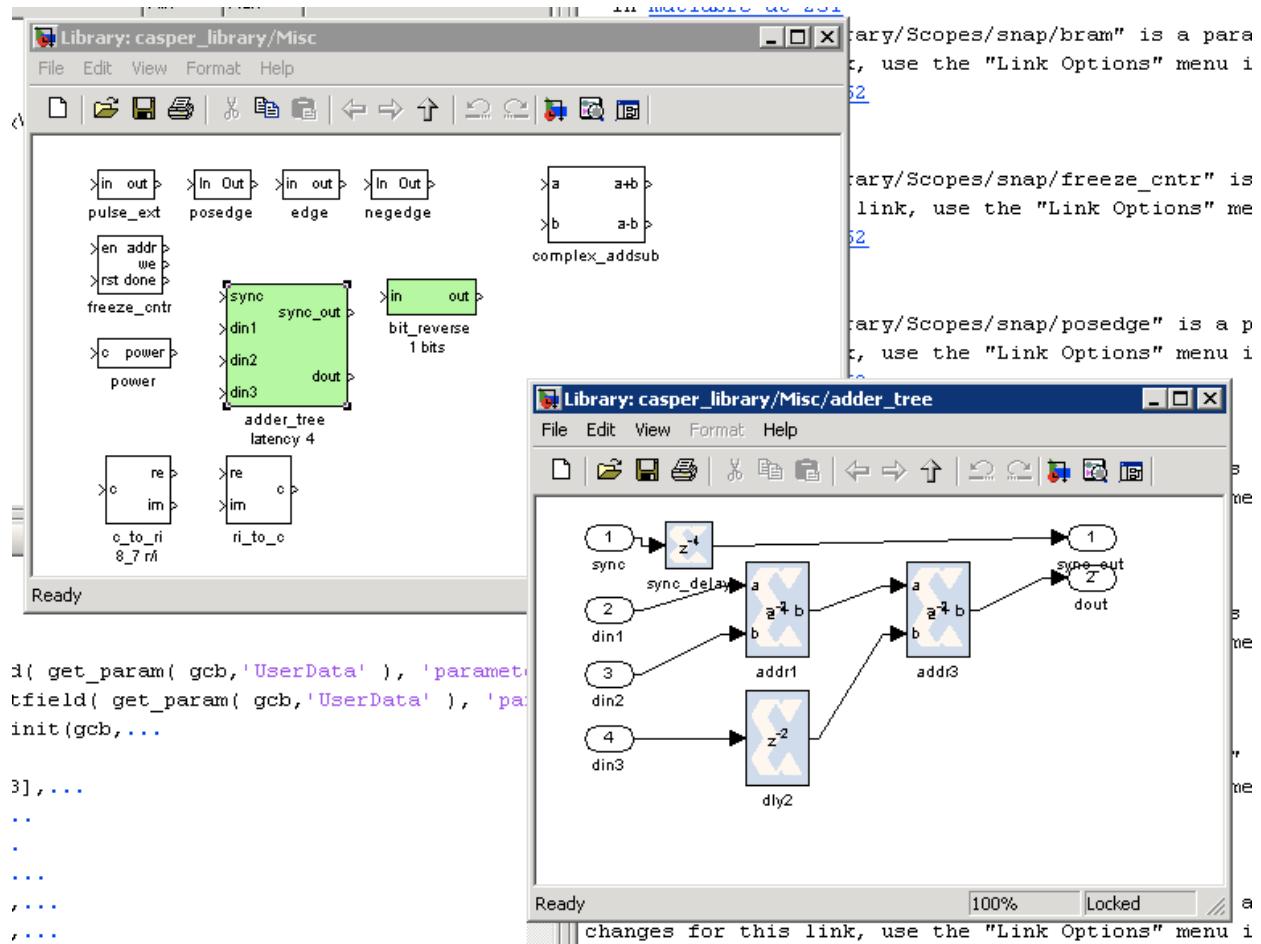


Figure 3.4: Adder Tree Simulink Diagram

In order to configure the number of inputs, Matlab scripts are used to redraw the implementations. Figure 3.5 shows a snippet of code from the script used to redraw the adder tree block. Every parameterized CASPER block uses at least one of these scripts to redraw the underlying implementation when the parameters are modified.

The CASPER library includes commonly used DSP blocks in radio astronomy instruments. For example, the CASPER library provides FFTs, FIR filters, accumulators, digital downconverters, digital mixers which can be linked together to make an instrument. Each block is parameterized, making them useful for a variety of instruments. Figure 3.6 shows the FFT library, which contains different types of FFTs, including blocks that can process multiple samples in parallel and blocks that are optimized to compute the FFT of a real signal.

Figure 3.7 shows the options menu for one of the FFT blocks. In order to support a variety of instruments, this block can be reconfigured to support different FFT lengths.

```

33
34 function adder_tree_init(blk,varargin)
35
36 - check_mask_type(blk, 'adder_tree');
37
38 - defaults = {'latency', 2};
39 - if same_state(blk, 'defaults', defaults, varargin{:}), return, end
40 - munge_block(blk, varargin{:});
41 - n_inputs = get_var('n_inputs', 'defaults', defaults, varargin{:});
42 - latency = get_var('latency', 'defaults', defaults, varargin{:});
43
44 - stages = ceil(log2(n_inputs));
45
46 - delete_lines(blk);
47
48 % Take care of sync
49 - reuse_block(blk, 'sync', 'built-in/import', 'Position', [30 10 60 25], 'Port', '1');
50 - reuse_block(blk, 'sync_delay', 'xbsIndex_r4/Delay', 'latency', num2str(stages*latency), ...
51 - 'Position', [30+50 10 60+50 40]);
52 - reuse_block(blk, 'sync_out', 'built-in/outport', 'Position', [30+(stages+1)*100 10 60+(stages+1)*100 25], ...
53 - 'Port', '1');
54 - add_line(blk, 'sync/1', 'sync_delay/1');
55 - add_line(blk, 'sync_delay/1', 'sync_out/1');
56
57 % Take care of adder tree
58 - for i=1:n_inputs,
59 -     reuse_block(blk, ['din',num2str(i)], 'built-in/import', 'Position', [30 i*40+20 60 35+40*i]);
60 - end
61 - reuse_block(blk, 'dout', 'built-in/outport', 'Position', [30+(stages+1)*100 40 60+(stages+1)*100 55]);
62
63 % If nothing to add, connect in to out
64 - if stages==0
65 -     add_line(blk,'din1/1','dout/1');
66 - else
67 -     % Make adder tree
68 -     cur_n = n_inputs;
69 -     stage = 0;
70 -     blk_cnt = 0;
71 -     blks = {};
72 -     while cur_n > 1,
73 -         n_adds = floor(cur_n / 2);

```

Figure 3.5: Adder Tree Simulink Mask Script

There are a number of other parameters provided like input bit width, which helps support a number of different ADCs or preprocessing algorithms, and FPGA-specific parameters like add latency, and multiply latency which have no effect on the result of the computation but change how the FFT gets mapped into hardware.

The CASPER library also provides an FIR filter than can be coupled with the FFT to create a polyphase filter bank or PFB. Figure 3.8 shows a comparison between the FFT and PFB response. The FFT response (on the left) has a lot of spectral leakage while the PFB (on the right) has a much sharper filter shape and a better frequency response. Despite the additional FPGA resources required by the FIR filter before the FFT, many instruments

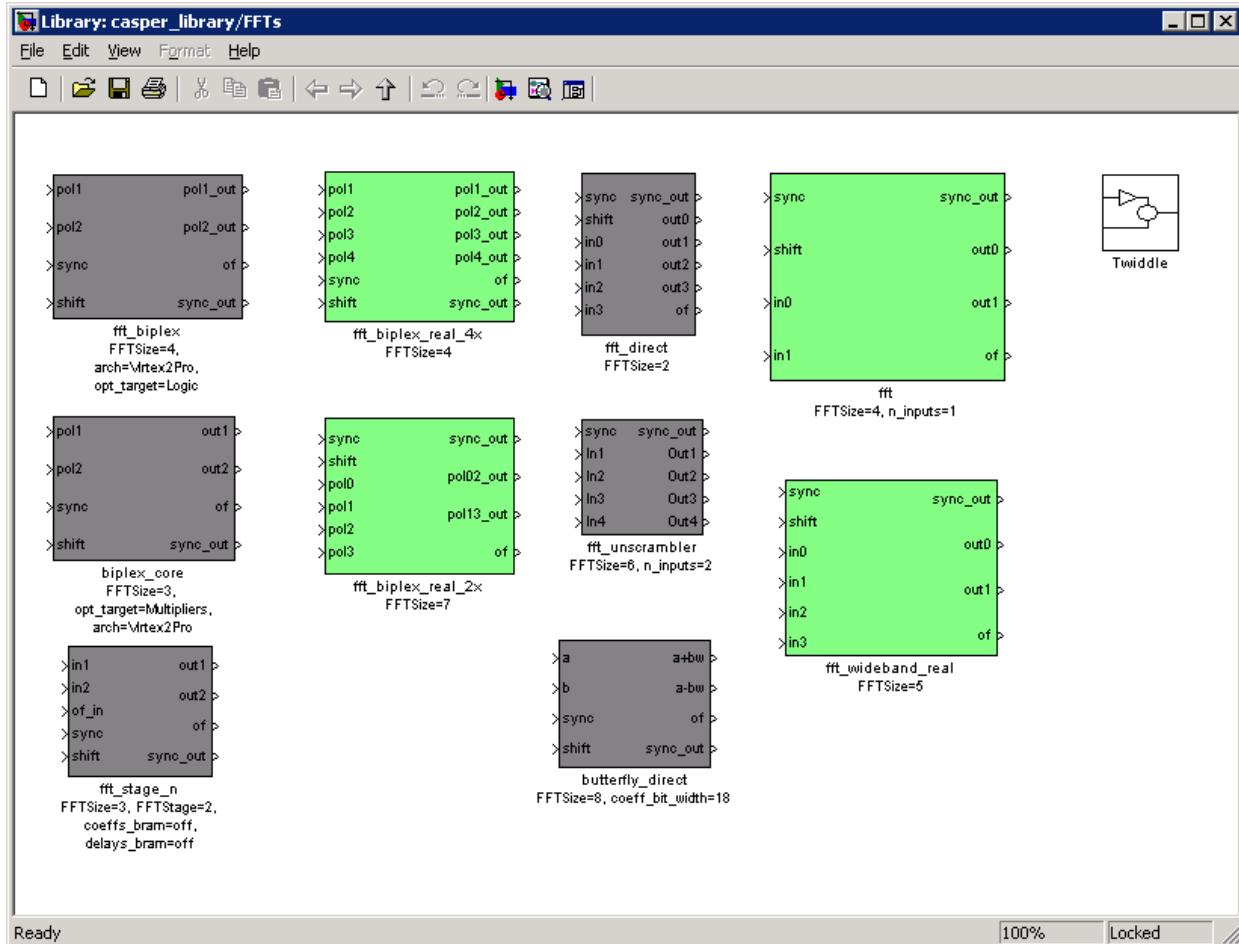


Figure 3.6: CASPER FFT Library

implement a PFB rather than an FFT because of the superior frequency response.

The CASPER library also provides a set of blocks to abstract away the implementation of I/O interfaces, which are called yellow blocks. Figure 3.9 provides some examples of what those library blocks look like. Each block provides input and output ports that correspond to the data it can send or receive. For example, the adc block has an output ports labeled $o0, o1, \dots, o7$ that represent the data the FPGA receives from the iADC board. In addition to this, simulation ports are provided to allow the user to proved test signals mimicking the outside world. In the case of the ADC, the user could provide a sine wave into the *sim_in* port to observe how a sine wave would be processed by the system. During compilation, the CASPER XPS toolflow automatically processes the blocks and ensures the wires are connected to the correct pins.

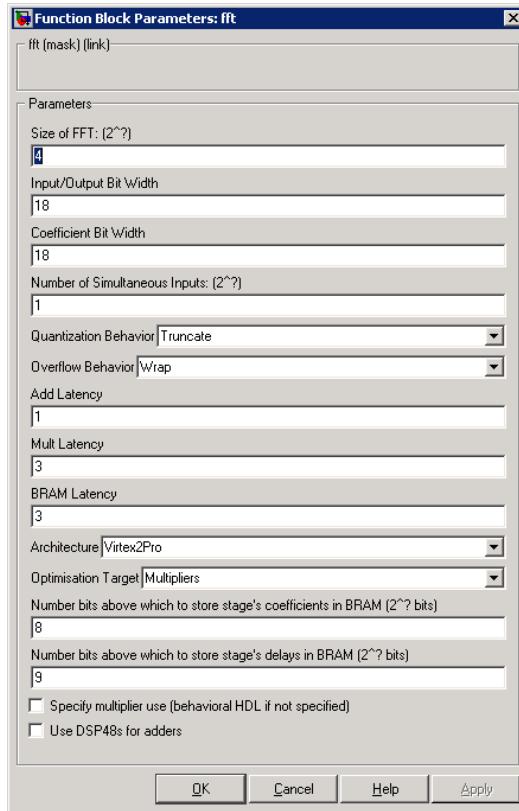


Figure 3.7: CASPER FFT Options Menu

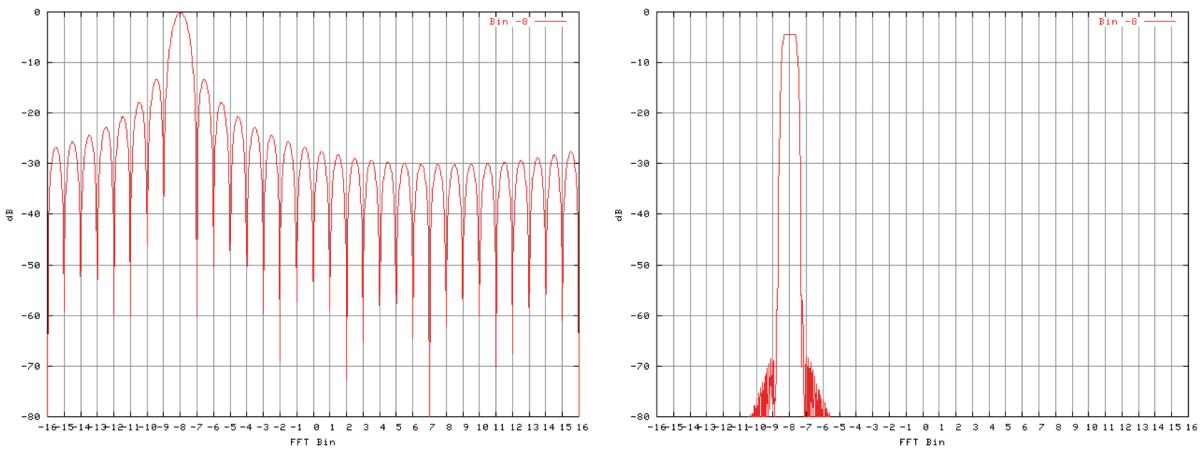


Figure 3.8: A comparison of FFT and PFB response

CASPER Software

To simplify the use of the FPGA further, the CASPER boards run a modified version of Linux directly on the board. Using this Linux environment, called the Berkeley Os for

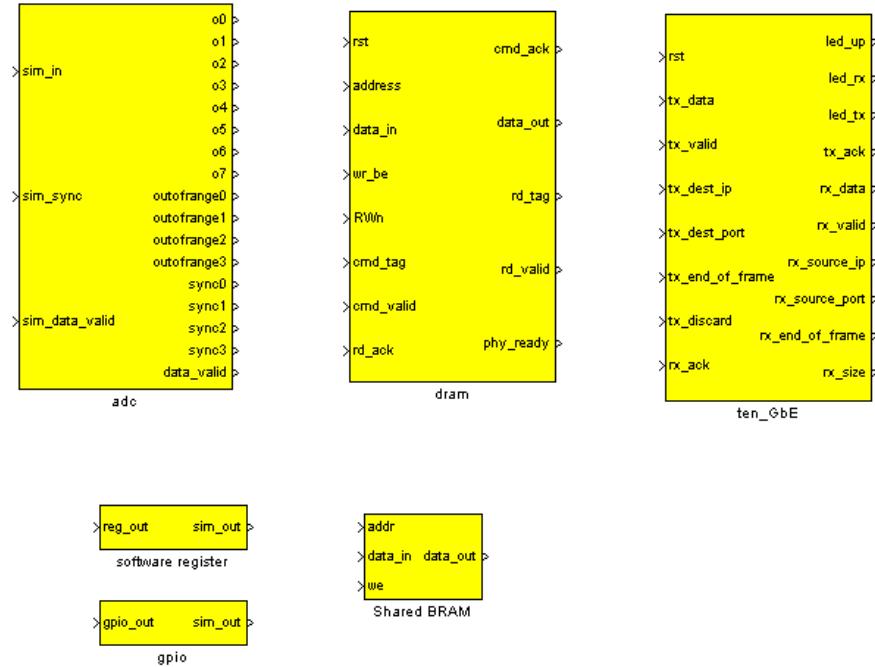


Figure 3.9: CASPER Library IO Blocks

ReProgrammable Hardware or BORPH, programming the FPGA is as simple as running an executable on the command line [25]. Once the board has been programmed, BORPH can communicate with the chip using an interface where some components on the FPGA like registers or memory appear as files in the operating system. These files can be accessed using normal file I/O, making it simple to send control signals or monitor the status of the chip, greatly simplifying debugging and command and control.

SERENDIP V.v

SERENDIP, or the Search for Extraterrestrial Radio Emissions from Nearby Developed Intelligent Populations, aims to find extra terrestrial intelligence by analyzing radio signals. SERENDIP V.v is a high resolution SETI spectrometer that was deployed at Arecibo Observatory in 2009 as part of this project. The installed instrument was built using an iADC, iBOB and BEE2 board. The instrument breaks up a single 200MHz stream into 128 million frequency channels, achieving a resolution of about 1Hz.

Figure 3.10 shows the dataflow for the spectrometer. The ADC data is fed into an IBOB where it is processed by a decimating downconverter.

The SERENDIP V.v instrument processes data from the Arecibo L-band Feed Array. The array has 7 dual-polarization beams, shown in Figure 3.11. The spectrometer

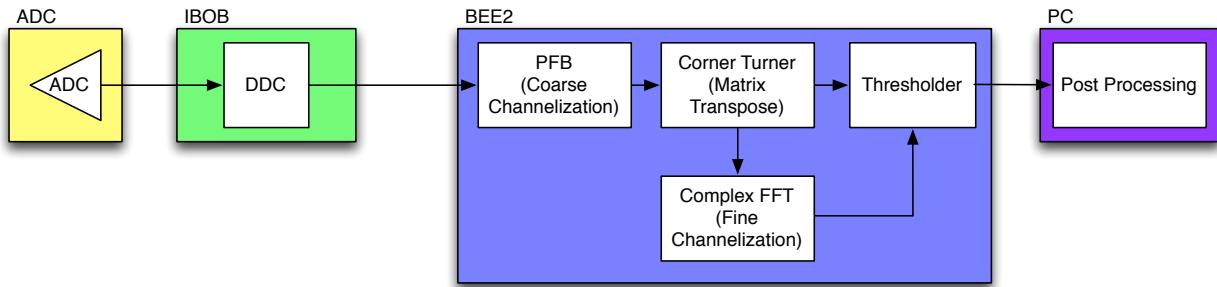


Figure 3.10: SERENDIP V.v Block Diagram

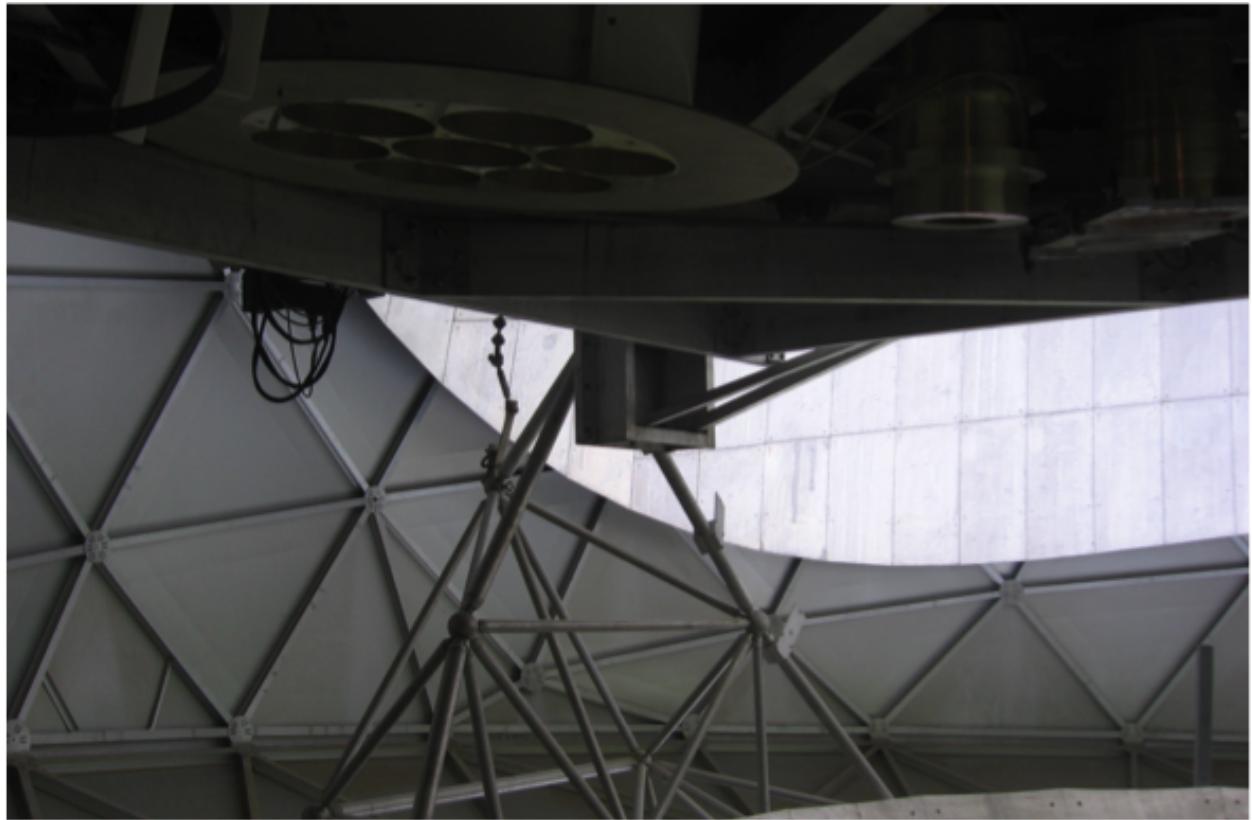


Figure 3.11: Arecibo ALFA Feed

Figure 3.12 shows the installed instrument at Arecibo Observatory.
[24]

CASPER Correlator

[20]



Figure 3.12: SERENDIP V.v Hardware installed at Arecibo Observatory

Packetized Astronomy Signal Processor

Many radio astronomy instruments channelize the antenna data as soon as it is digitized. PASP implements a parameterized, FPGA-based channelizer, shown in Figure 3.14. The FPGA interfaces to an ADC board that simultaneously digitizes 2 signals. The samples are sent into a polyphase filter bank (PFB), consisting of an FIR filter and an FFT, which breaks up the entire bandwidth sampled by the ADC into smaller subbands. After dividing up the subbands, each band is rescaled. This step allows us to compensate for the shape of the analog filter feeding data into the ADC. After rescaling, the FPGA forms packets where each packet contains data from a single subband. The packets are sent out over CX4 ports to a 10 gigabit Ethernet switch.

PASP is designed for flexibility. Building on the CASPER goal to automate the design of commonly used signal processing elements such as FFTs and digital downconverters, PASP automatically designs an entire FPGA instrument using only a few parameters. The user can input the desired number of subbands, CPU/GPU cluster size, and packet size and a new design is automatically generated in Simulink. Figure 3.15

The interface to PASP is a simple menu that consists of four parameters. From the menu,

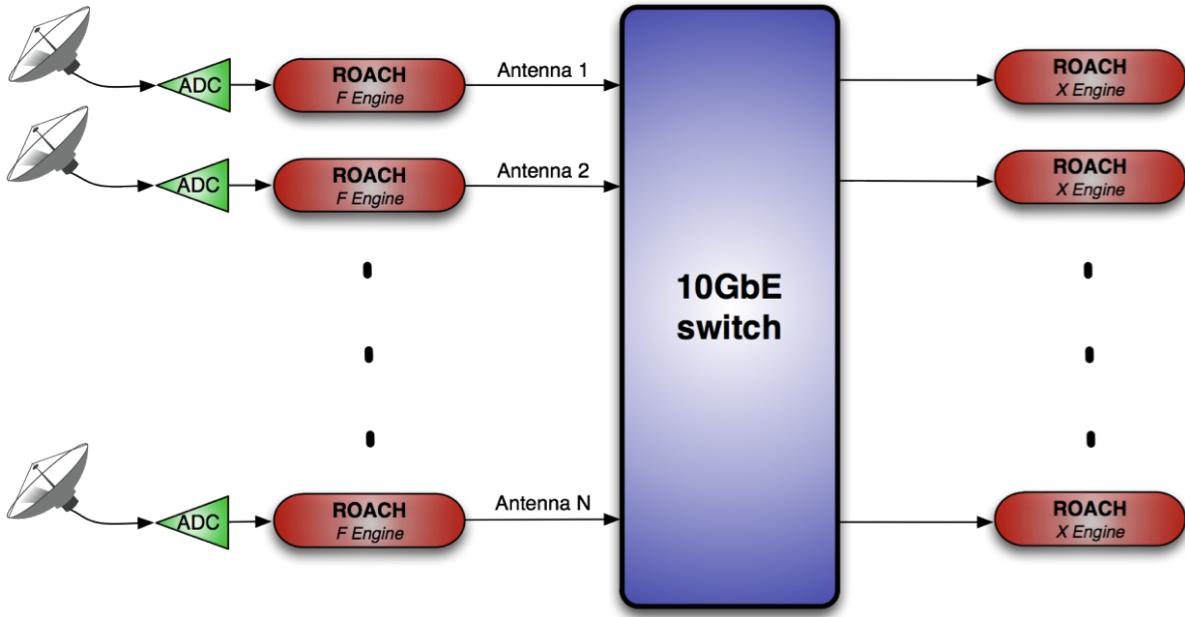


Figure 3.13: CASPER Correlator Block Diagram

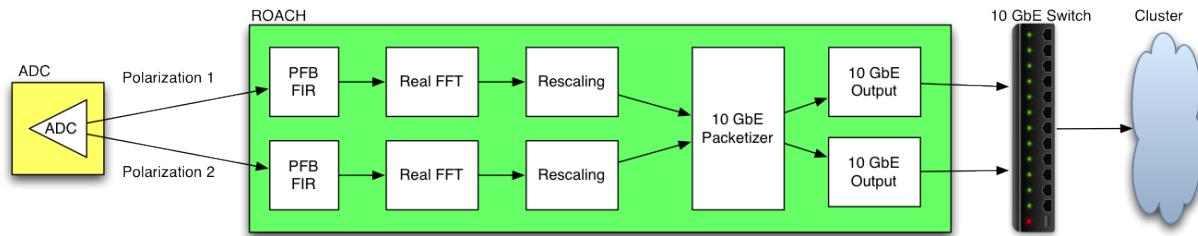


Figure 3.14: PASP Dataflow

Number of IPs determines the number of IP addresses the instrument will need to send data to, defined by the size of the cluster it must communicate with. *Samples per Packet* defines the packet size, adjustable to ensure that the packet rate is low enough for the server to handle. And *Number of Channels* defines the FFT length. The last parameter *Reorder in QDR*, allows the packet buffering to be handled by off chip QDR memory. This is only a parameter and not the default because it is only supported by certain boards.

Changing any parameter causes the entire low-level design to be redrawn. Figure 3.16 shows the low level implementation of the instrument. Although this diagram is complex, it is generated automatically by PASP and the user never needs to look at this design to implement a working instrument.

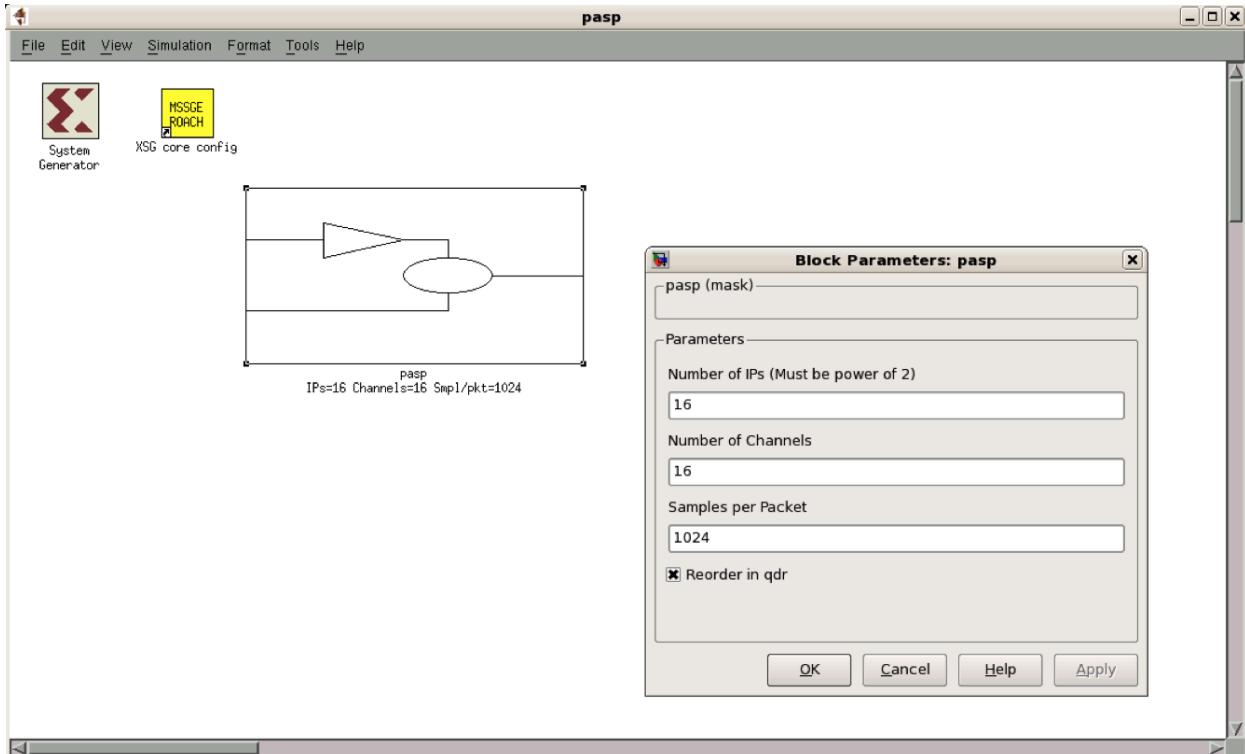


Figure 3.15: PASP Interface

PASP has a wide variety of potential applications due to the flexibility of the server software. In the next sections, we describe a few specific applications than can make use of this package.

Pulsar Processing

This instrument was originally designed to do pulsar science. The fast channelization on the FPGA with no data reduction makes it an ideal pulsar spectrometer, since no information is lost before sending the data to the servers. GPUs provide a good platform for pulsar processing algorithms such as coherent dedispersion [23], which can easily be used as the processing function for the server software distributed in our package. Similar to SETI instruments, pulsar instruments designed using this package can also keep up with improvements in technology with a simple recompile.

This design has been used in pulsar processors that have been deployed at the Parkes Observatory, Nanay Decimetric Radio Telescope, and Effelsberg Radio Telescope.

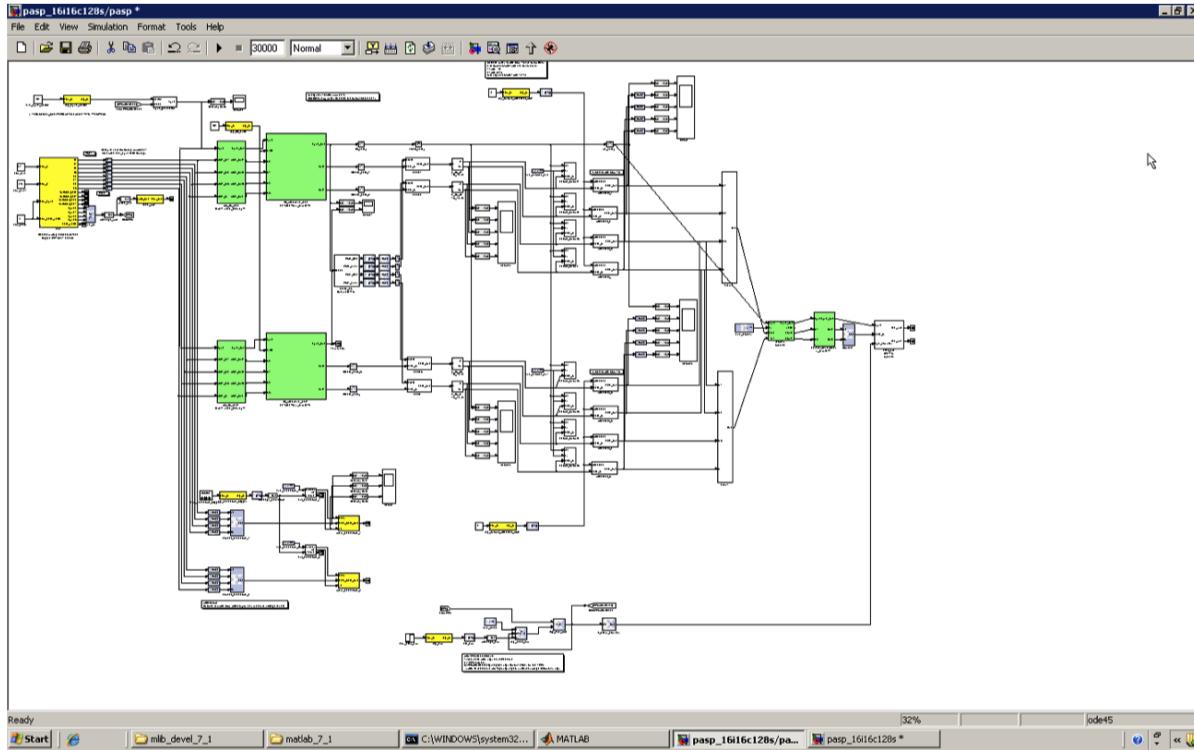


Figure 3.16: PASP Low Level Implementation

Heterogeneous Radio SETI Spectrometer

In the search for extraterrestrial intelligence (SETI), the ability to keep up with changes in technology allows searching instrumentation to stay on the leading edge of sensitivity. SETI aims to process the maximum bandwidth possible with very high resolution spectroscopy. This instrument allows SETI projects to easily keep up with improvements on the telescope and increasing computational power. An increase in detector bandwidth, improving the breadth of the search, can be processed simply recompiling the FPGA design and distributing the extra subbands to new servers. As computation improves, the instrument can be reconfigured to send more bandwidth to each computer, reducing the required cluster size, or improve the resolution of the instrument by doing a larger FFT on the server.

The design of HRSS is based on the Serendip V.v design presented in Section 3.1. Figure 3.17 shows how in HRSS the IBOB and BEE2 are replaced with a single chip FPGA board such as an IBOB or ROACH, and the fine channelization, thresholding and post processing are done on a GPU.

Unlike Serendip V.v, HRSS is a software package to automatically generate spectrometers with minimal user input. We have automated the design of the spectrometer, creating a parameterized spectrometer that only requires a recompile to implement a change in specification. This spectrometer combines FPGAs running PASP with GPUs, doing coarse

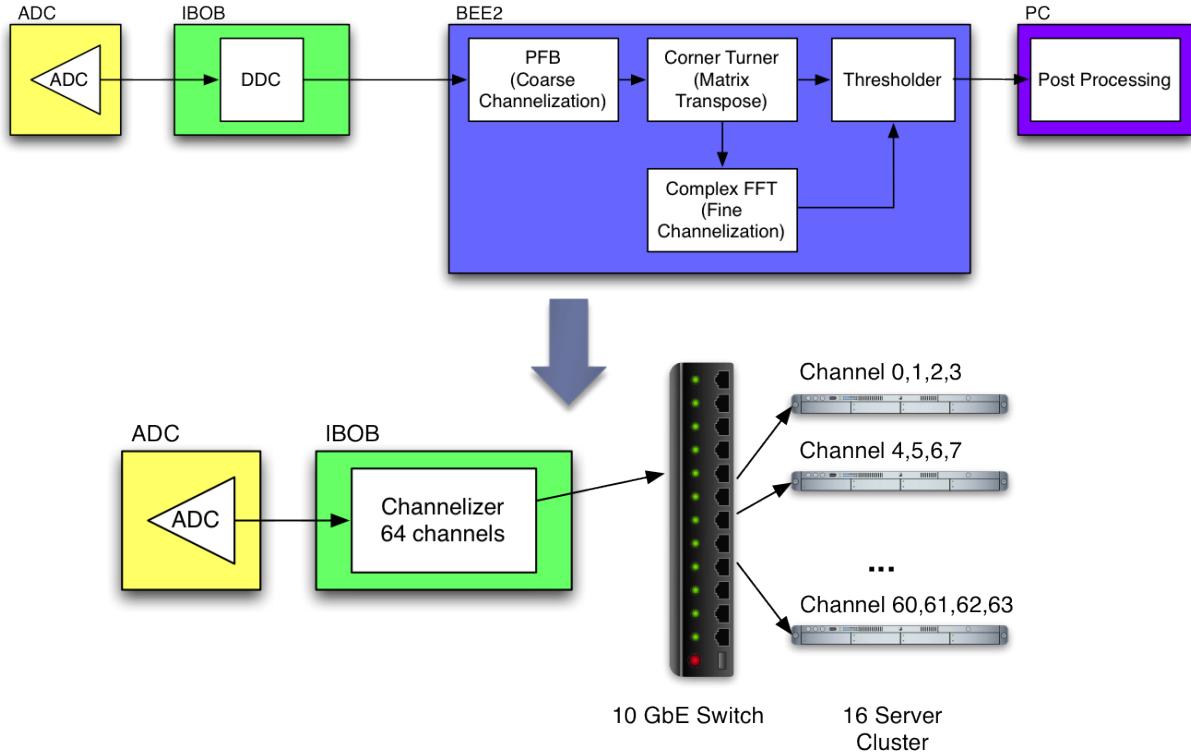


Figure 3.17: TODO

channelization on the FPGA and sending each subband to the GPUs for further processing. The server software is designed for flexibility, allowing astronomers to easily modify the processing algorithm run on the GPU and customize the instrument to fit their science goals.

The software package is comprised of two parts, the PASP design that runs on the FPGA and the server software. Like PASP, the server software is parameterized, supporting a variety of spectral resolutions and integration times. This software receives data over an Ethernet port and transfers it from the CPU to the GPU. The GPU runs an FFT and then sends the data back to the CPU to be recorded. The GPU software, like the GPU benchmark, uses the CUFFT library to run FFT. This allows for rapid deployment of a working spectrometer that is configured to take full advantage of available computing resources.

Figure 3.18 shows a high level view of a spectrometer that could be designed with this package. In this example, a ROACH board divides the input band into 64 subbands and sends them out to a 16 server cluster. An ADC is used to digitize data from the telescope and connects to the ROACH board via Z-DOK connectors. The digitized data is split into 64 subbands and sent through a 10 gigabit Ethernet switch. Each server in the cluster receives and processes 4 subbands.

The server software was designed so other applications could easily be implemented on

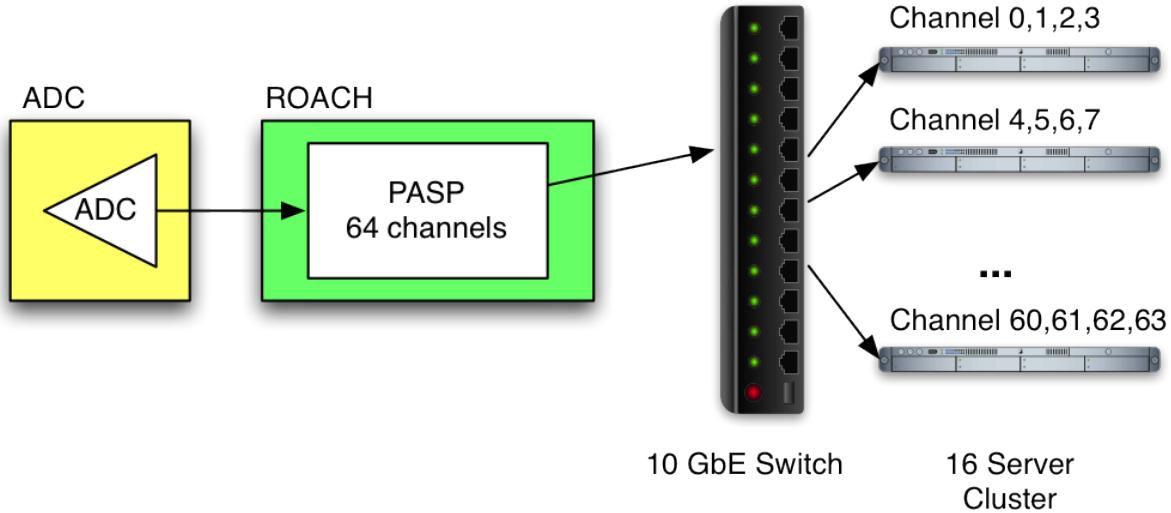


Figure 3.18: Example high level spectrometer architecture

the GPU without altering or rewriting the receive code that interprets the packet headers and transfers data to the GPU. Once the data is on the GPU, the software calls a process function and passes it a pointer to the GPU data. An initialization function is called before the data processing begins to do any setup needed by the processing function, and a corresponding destroy function cleans up once the processing is complete. In the spectroscopy software included in the package, the initialization function creates the FFT plan, the processing function calls CUFFT, and the destroy function deletes the FFT plan. Modifying the application run on the GPU simply requires a redefinition of these three functions. Using this interface, we successfully replaced the CUFFT processing with software developed for SETI searches designed by Kondo et al. [17].

3.2 Tuning

Metropolis

The Metropolis project focused on mapping algorithms onto embedded systems [5]. The tool provides a framework for an abstract block based description of the algorithm. This description makes it easy to stitch algorithms together without specifying the eventual hardware implementation, providing a simple path to simulation and algorithm development that is separate from the implemented design. Then, the tool automatically maps the description onto an existing heterogeneous embedded system.

[2] [8]

The mapping generated by the tool is simply a schedule, specifying where and when each part of the computation gets executed. The tool seeks to optimize performance of the algorithm, ensuring the generated schedule runs as fast as possible on the hardware provided. This technique requires a fixed hardware model and uses the existing hardware to optimize performance. While this work is useful when a hardware model exists, it does not provide any flexibility in the hardware model while mapping the algorithm. So, when it is necessary to design the hardware to begin with, this does not solve the problem.

Additionally, this type of solution is ill-suited to mapping the algorithms required to do real-time radio astronomy signal processing. This tool assumes it must schedule a discrete task onto a fixed piece of hardware and attempts maximize the performance of the task. In the applications described in Chapter 2, the computation should always be running and needs to meet some minimum performance target. Once the performance target is met, it is better to have a tool that will minimize other costs like power or amount of hardware, rather attempting to improve the runtime of the algorithm.

Scheduling

An integer linear programming model for mapping applications on hybrid systems

[11]

[26]

[28]

[14]

[22]

Chapter 4

High Level Toolflow

Instrument design is often done by building the instrument from scratch. This work extends the CASPER philosophy, demonstrating that entire heterogeneous instruments can be designed with minimal user input. Rather than designing a completely different instrument for every different specification, this software package is parameterized so new designs can be generated quickly.

In this chapter, I introduce a toolflow I developed called ORCAS or Optimal Rearrangement of Cluster-based Astronomy Signal processing. ORCAS automatically sifts through different designs and chooses the best platforms based on the final cost of the instrument. This automatic approach makes it easy to keep up with constant changes in technology and algorithmic implementations. ORCAS is released as an open source Python library, freely available on GitHub at <https://github.com/tfiliba/orcas>.

4.1 ORCAS Goals

ORCAS extends the CASPER philosophy discussed in Section 3.1 by providing an end-to-end toolflow that allows the user to go from a high level description of the instrument to a low level mapping automatically. While the CASPER toolflow has proven successful in the FPGA domain, it's not always clear that FPGAs are the most cost-effective implementation for a given instrument. The ORCAS captures the successful elements of the CASPER tools while extending its reach to the heterogeneous domain. One of the most important elements of the CASPER tools is the library of DSP blocks. The CASPER library blocks are constantly being updated, ensuring that the library can always be compiled on the latest technology. By incorporating existing libraries implemented for CPUs, such as FFTW, of GPUs, like cuBLAS, xGPU, and CUFFT, ORCAS can also keep up with changing technology, without the need for constant updates to the tool itself. The flexibility of the libraries also allows us to buy technology at the last minute, ensuring we get the cheapest price available for an instrument that will provide the needed performance.

In addition to the elements drawn from the CASPER library, ORCAS emphasizes a

twofold approach to cost reduction in instrument design. The first cost savings is in the design time. The tool is designed to map designs to hardware within a few hours. This fast mapping makes it possible to quickly assess performance on new or nonexistent technology and allows the designer to assess the impact of potential optimizations, making it easy to test a number of different design ideas without spending time testing different implementations. By incorporating tested libraries, there is much less debugging needed in the eventual implementations.

The second cost ORCAS is able to reduce is the cost of the instrument itself. While the implementation of a design has been greatly simplified by the availability of tools described in Section 3.1, none of those tools help determine which type of hardware is best. The extensive library allows the tool to support whatever hardware is cheapest, allowing the choice of hardware to be completely determined based on cost. This cost could be defined in a number of ways, depending on which parameters the instrument designer needs to minimize. Obviously the cost could refer to the monetary price of the hardware, but the tool can also be used to minimize the amount of power needed by the instrument, the physical size of the instrument, and other user defined costs.

ORCAS is designed to be accessible to both computer experts and radio astronomers, or other domain experts. For a domain expert, ORCAS allows the user to choose a predefined instrument type and select high level parameters without worrying how it will ultimately map to hardware.

The computer expert can use the tool to define a new instrument, by specifying the algorithm, and is able to provide optimization. For example, if an engineer wrote an optimized FFT algorithm, the tool will be able to incorporate that into the final optimized result. In the end, regardless of who is using this tool, a cost-optimal mapping of the instrument gets produced.

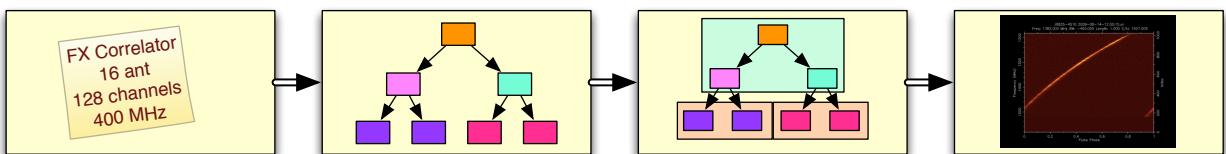


Figure 4.1: The Four Stage ORCAS Toolflow

These goals are achieved through a four stage tool that allows the user to go from a high level description of an instrument to an optimized cluster design. Each stage of the ORCAS toolflow is represented in Figure 4.1. The first stage is instrument definition, which is described in Section 4.2. This stage is designed to fit the needs of the domain expert. It allows the radio astronomer to create a predefined instrument using a handful of parameters. The instrument definition is converted into a dataflow model, as described in Section 4.3. The second stage, the dataflow model, is aimed at the computer expert, since it allows for a more detailed and flexible definition of the algorithm than the previous stage and

provides the means to include optimized blocks. The dataflow model represents an abstract definition of the algorithm without taking into account the eventual hardware target. The next stage, appropriately named mapping, maps the dataflow model to specific hardware. In the mapping stage, each block defined by the dataflow model gets mapped to a specific piece of hardware. This stage takes into account hardware and network limitations to produce a cost optimal mapping of the original dataflow. Mapping is discussed in Section 4.4 and Chapter 5 provides more details about the algorithm used to produce a mapped instrument and the performance of that algorithm. Finally, once each block has been mapped to a piece of hardware, the code can be stitched together into a working instrument, as presented in Section 4.5. The rest of the chapter will describe each stage of the toolflow in more detail and explain how they are designed to meet the needs of the users they are targeting.

4.2 Instrument Definition

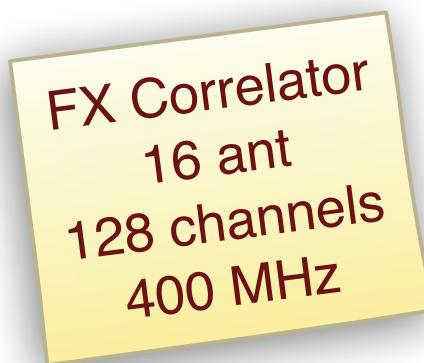


Figure 4.2: ORCAS Toolflow Instrument Definition

In the first step in the ORCAS toolflow, the user must describe the instrument using high level parameters. These parameters should all be relevant to the astronomer and abstract away any implementation details that do not pertain to the scientific goals. While it would be easy to expose many of the low level parameters at this top layer, this would force the domain expert to become a computer expert as well, exactly the scenario this tool is aiming to avoid.

The instrument description as represented in Figure 4.2 fits on a small sticky note. The idea that the parameters should be so few that they fit on a single sticky note was the driving force behind the instrument definition. An instrument designer who finds that he or she needs additional control beyond what is provided by the instrument description always has the flexibility to work with the dataflow model directly, where many low level parameters are

exposed. Additional instruments or new parameters for existing instruments can be added by defining how those parameters affect the final dataflow diagram.

In software, each instrument type is represented by a Python class. The instrument parameters used when an instrument object is instantiated. For example, to create the instrument specified on the sticky note in Figure 4.2 would simply require the following function call (variables are defined for clarity):

```
antennas = 16
channels = 128
bandwidth = 0.4 #defined in GHz

#create the instrument
myfxcorrelator = FXCorrelator(antennas, channels, bandwidth)
```

Design space exploration

In addition to defining a single instrument, an astronomer can use this tool to explore different implementations and assess the tradeoff in cost vs additional processing. Rather than specifying single values for the instrument parameters, the astronomer can choose a range to search through and have it generate a design, and an associated cost, for each value. This proves valuable when the specification of an instrument isn't fully defined. In the case of a high resolution spectrometer, the astronomer could adjust the number of coarse and fine channels, keeping the spectral resolution the same, to find the cheapest design that achieves that resolution. For polyphase filter banks, the number of taps in the FIR can be varied to assess the increase in cost associated with a better filter response. Even non-numerical parameters can be varied, such as the FIR window type, to see how it will affect the eventual design.

4.3 Dataflow Model

The second step of the toolflow translates the instrument definition into a dataflow model. The dataflow model represents the instrument as a set of high level blocks (such as FFTs and FIRs) and input and output connections. Figure 4.3 shows an abstract example of this type of dataflow.

Computational Blocks

The computational blocks are a set of algorithmic building blocks. These blocks are similar to the set of blocks provided by the CASPER library or CUFFT, allowing for parameterization, but unlike the afore mentioned libraries, these blocks do not necessarily include an implementation of the function they claim to provide. Instead, each block provides some method to asses the performance of that function on each supported platform.

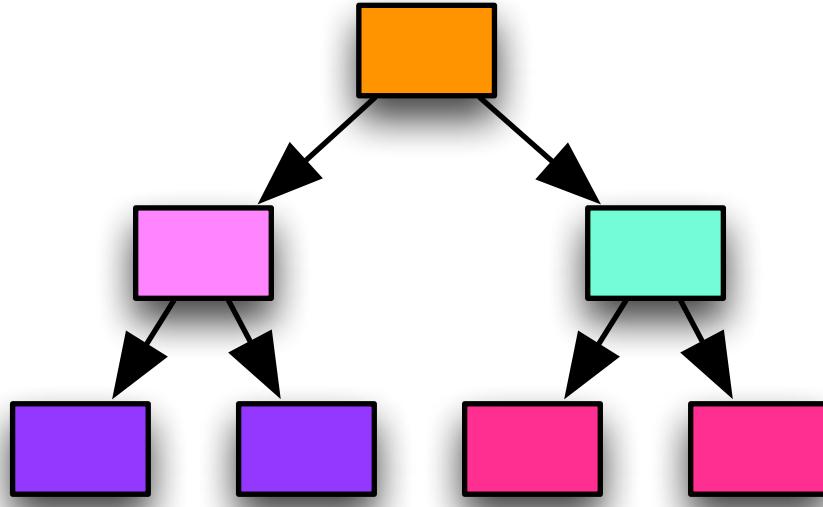


Figure 4.3: The ORCAS Toolflow: Dataflow Model

One way to assess performance is to simply include the code, generating a model by synthesizing or running the code and getting real performance data. But this is not always necessary or desirable. Many algorithm designers provide benchmarks that can be used directly without needing to download or run the code [12] [15]. The predictable scaling of many DSP functions can also make it easy to create a model for resource utilization [21]. And finally, the user can always estimate the performance of a block, a useful option when attempting to assess the cost savings of a potential optimization or explore mapping to hardware that does not exist.

Connection Types

In order to reduce the total number of connections that need to be defined, blocks that receive data from and sends data to the same set of block are put into a group. Figure 6.5 shows an example FX Correlator dataflow model, where blocks grouped together have the same color. In this example, the FFTs all take data from an FIR block and must send data to the XEng blocks. The name of the group of blocks is referred to as the "block type".

Knowing which block types must communicate with each other is not sufficient to specify the dataflow. Suppose we know we have 2 types of blocks: A , and B and blocks of type A must send their data to blocks of type B . We also need to definite the communication patterns between A blocks and B blocks. This could happen in 2 ways, 'one-to-one' and 'all-to-all'.

A 'one-to-one' connection is where every A block communicates with exactly one B block,

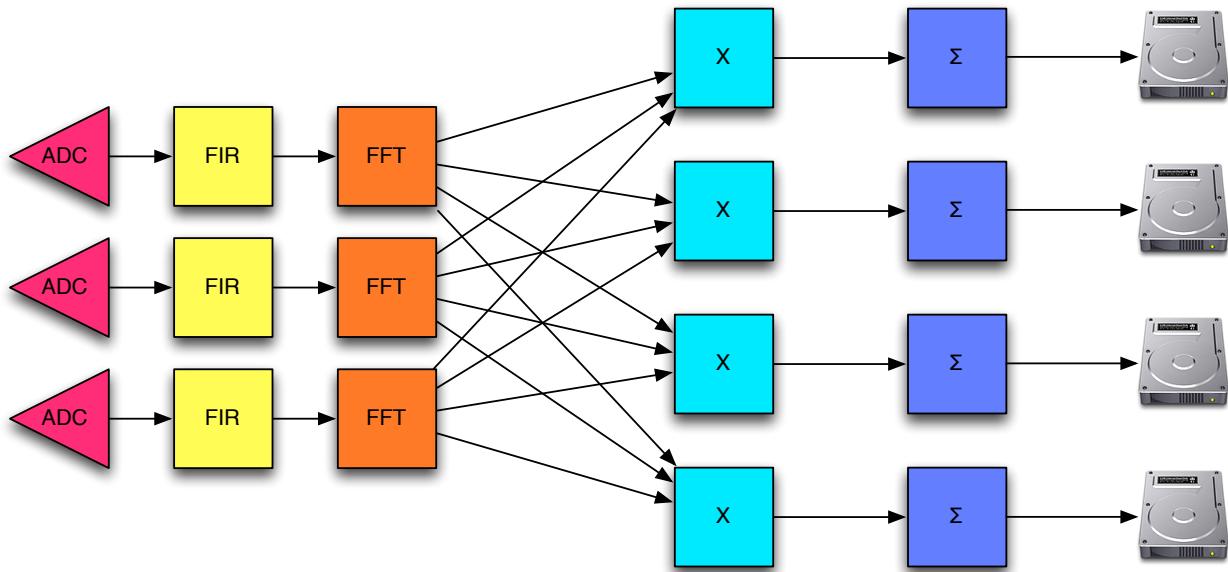


Figure 4.4: Example FX Correlator Dataflow Model Demonstrating Blocktypes

as shown in Figure 4.5. With this type of connection, the number of A blocks must be equal to the number of B blocks. The F-engines in a FX correlator are a good example of this type of connection. The correlator has an F-engine for each antenna, each containing the same blocks linked in the same way. Within an F-engine, a PFB_FIR filter must communicate with a single FFT. In general, every PFB_FIR within an F-engine, blocktype ‘A’ must communicate with exactly one FFT, blocktype ‘B’.

An ‘all-to-all’ connection occurs when every A block must send some data to every B block. Figure 4.5 shows what an all-to-all connection between three A blocks and three B blocks will look like. In this case, every A block must send some data to every B block. For example, the type of connection between the per-antenna FFTs and the per-channel X-engines in an FX correlator would be ‘all-to-all’. Each X-engine needs a small amount of data from every F-engine to compute the cross-correlations from a single channel. In the ‘all-to-all’ case, there is no reason for the number of sending nodes needs to be the same as the number of receiving nodes.

It might seem like there are two more possible types of connection that need to be defined, ‘one-to-all’ and ‘all-to-one’ . A data flow with a ‘one-to-all’ connection, shown on the left in Figure 4.6 would have exactly one A block that needs to send data to many B blocks. This is exemplified in the dataflow for a high-resolution spectrometer. The coarse channelization is done in a single FFT block, which then needs to send the data to many other FFTs to do the fine channelization. The ‘all-to-one’ connection is also shown on the right in Figure 4.6 is the reverse of the ‘one-to-all’ case. In this type of connection, there are many A blocks and they all need to send data to a single instance of a B block. An example of this arises

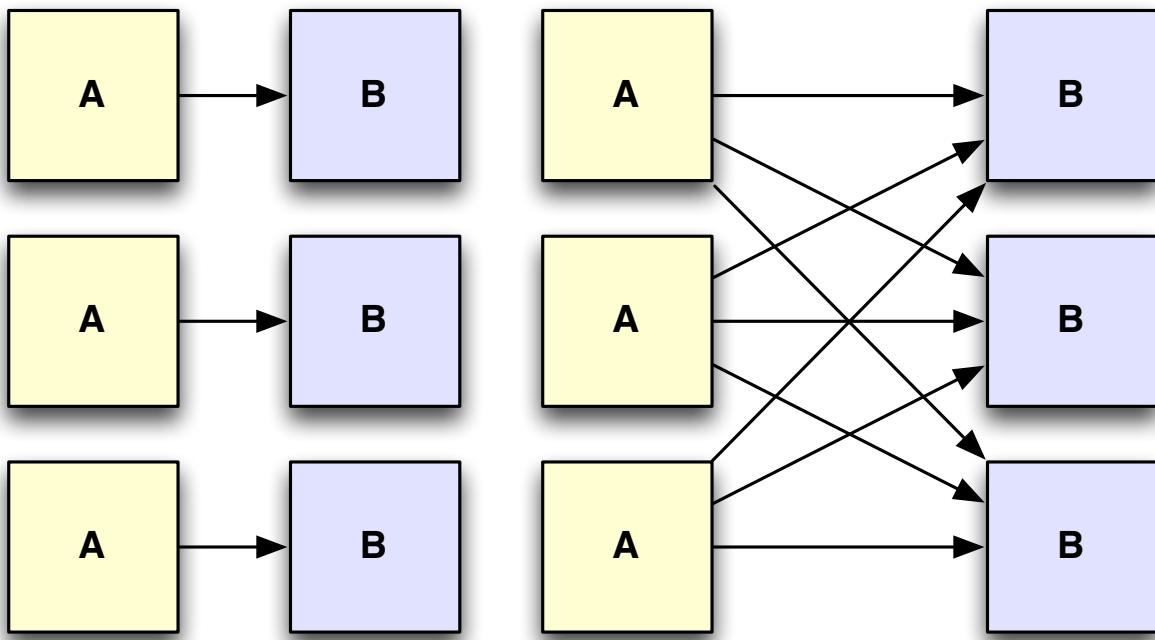


Figure 4.5: TODO

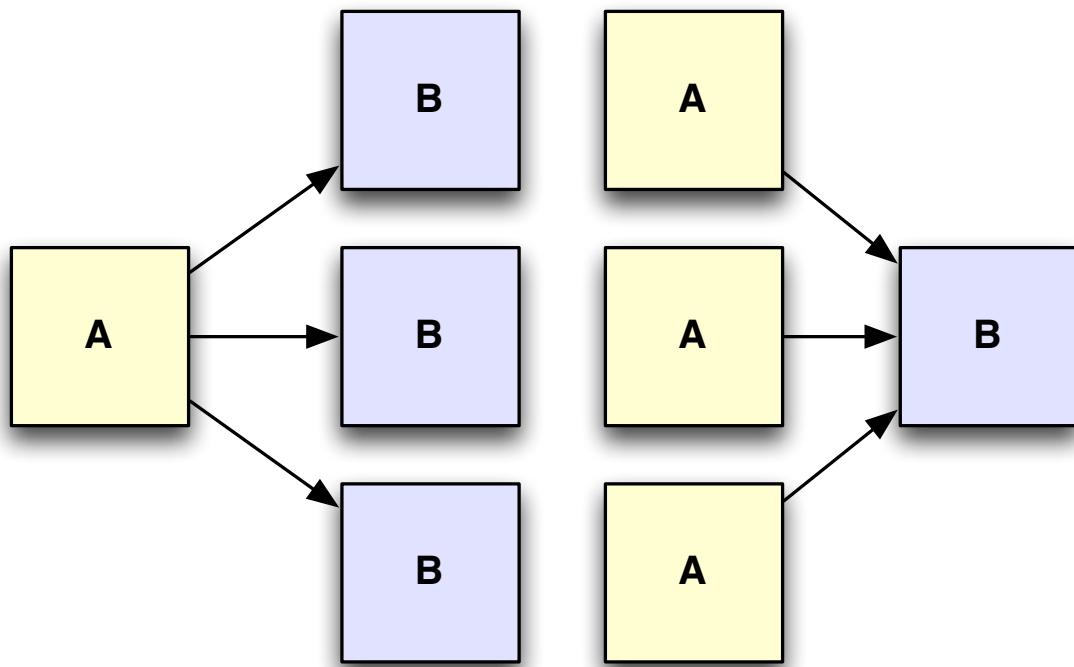


Figure 4.6: TODO

when some processing is done in a distributed manner but the instrument needs to record the final result in a central place.

It turns out, these are both special instances of the ‘all-to-all’ connection. The ‘one-to-all’ connection is simply an ‘all-to-all’ where the number of A blocks is fixed at one. Similarly, the ‘all-to-one’ connection is also an ‘all-to-all’ where the number of B blocks is fixed at one. Because of this, there is no need to include or support these cases as unique connection types.

While it may seem like additional link types exist like ‘all-to-some’ or ‘one-to-some’, this turns out to be impossible. Either a block of type A cannot send its data to only some blocks of type B because of the way blocktypes are defined. Any block of the same type should be interchangeable with another block of the same type. In an ‘all-to-some’ connection, blocks of type A would need to send data to B_1 but not send data to B_2 . But that connection patterns implies that the blocks B_1 and B_2 are *not* interchangeable and therefore cannot have the same blocktype.

Software Representation

The dataflow model is created in the initialization function for the instrument. Each block type is added to the dataflow by creating an instance of the CBlock class and adding that object to the list of blocks. The constructor for the CBlock has the user define the number of blocks in the group, and the input and output connections, indicating both the blocks that it is communicating with, the connection type, and the connection bandwidth. A CBlock also needs to be instantiated with a performance model, which is discussed further in the following section. The code below gives an example of this kind of instantiation. It creates 16 FFT blocks that must take data from FIR blocks via a 1-to-1 connection, and sends data to XEng blocks via an all-to-all connection.

```
performance_model = CBlock.getFFTModel( self.platforms , bandwidth ,
    input_bitwidth )
data_source = 'PFB'
source_connection = 0 #0 indicates a one-to-one connection
source_bandwidth = 0.4*8
data_sink = 'XEng'
sink_connection_type = 1 #1 indicates an all-to-all connection
sink_bandwidth = 0.4*8
antennas = 16 # indicates the number of blocks to create
self.blocks[ 'FFT' ] = CBlock( performance_model , data_source ,
    source_connection , source_bandwidth , data_sink ,
    sink_connection_type , fft_out_bandwidth , antennas )
```

4.4 Mapping

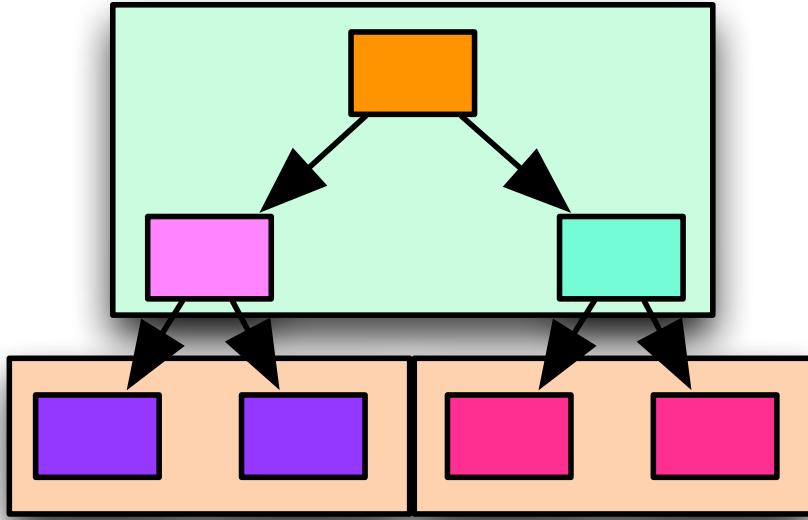


Figure 4.7: The ORCAS Toolflow: Mapped Dataflow Representation

Once the dataflow and the computational blocks are defined, ORCAS must determine how to place each computational block into hardware. In the mapping stage, ORCAS determines what type of hardware should be used for each block, while minimizing the total cost of the hardware (in dollars, watts, or another user-defined metric). The icon in Figure 4.7 represents an abstract mapped dataflow.

At this point, the computational blocks and dataflow can no longer be viewed as abstract algorithms. Each computational block must be paired with a performance model for each supported platform that shows the resource utilization and bandwidth requirements for that block. Using the performance model, the tool is able to test a number of hardware mappings and ensure that none of the available hardware or bandwidth resources are overmapped.

The optimal mapping is determined using an Integer Linear Program or ILP. The resources, such as memory, logic and CPU time, and bandwidth constraints are translated directly into ILP constraints. These constraints are used to determine a valid mapping, for example total bandwidth mapped to a link must be less than total link bandwidth. The variables represent the design decisions, determining where each block should be implemented. And finally, the cost function simply totals up the costs associated with each piece of hardware used.

A ILP was chosen because it has a number of positive features. Unlike a randomized algorithm such as simulated annealing, the results from a ILP are repeatable. While there might

be multiple solutions with the same cost, each time the same ILP is run, it is guaranteed to find one of the solutions of optimal cost. The ILP representation also makes it easy for the user to guide the algorithm. Since the design choices are represented by variables, they can also be restricted by adding additional constraints to those variables. This representation also makes it easy to build out an existing cluster, by allowing a limited amount of zero-cost hardware.

Unfortunately, the advantages of the ILP come with a high cost, namely that an ILP is NP-Hard to solve optimally. Current ILP benchmarks are able to solve problems with a hundred thousand variables in a few hours, but beyond that size the problems become infeasible. To make matters worse, the ILP runtime is very sensitive to the solver being used and the problem structure.

Because the runtime for the ORCAS mapper needs to be within a few hours, we use a number of techniques to reduce the ILP runtime. First, the easiest way to reduce the runtime without changing the ILP is to change the solver. There is a huge amount of variance in runtimes between different solvers and simply switching out the backend solver might cause a previously infeasible problem to become solvable. When that doesn't improve the runtime enough, it becomes necessary to modify the ILP. One way to do this is by reducing the number of variables it needs to solve for. Many of the radio astronomy instruments are very symmetric, so it is reasonable to assume that the optimal mapping would be symmetric as well. The symmetry can be preserved by forcing every block of a certain type to be implemented in the same kind of hardware, drastically reducing the number of decisions that the ILP needs to make. The ILP can also be modified to ensure that there is a single, unique optimal solution. When designs are very symmetric, the ILP will often find an optimal solution quickly but, because it is not unique, the ILP must spend a lot of time convincing itself that the solution is, in fact, optimal. Additional constraints can be added to ensure that only one of those solutions is valid, greatly reducing the amount of time it takes to verify optimality.

Chapter 5 goes into more detail on how the performance models are used as well as how the ILP is specified, implemented, and optimized to achieve a feasible runtime.

Mapped Dataflow Representation

Once the mapping is complete, the dataflow model is updated to describe what type of hardware will be used to implement each computational block. Figure 4.8 shows two or many possible mappings for the FX Correlator dataflow shown in Figure 6.5. In the top dataflow, each F-engine and X-engine require so many resources that they must be implemented on separate boards. Each F-engines is implemented on an FPGA-based ROACH board and the X-engines are implemented in GPUs. The bottom dataflow shows a less resource-hungry correlator, allowing all the F-engines and two X-engines to share a single board for computation.

In all designs, any link between two blocks that are not mapped to the same board must pass through a switch. In the example, all the links between FFTs and cross-correlation

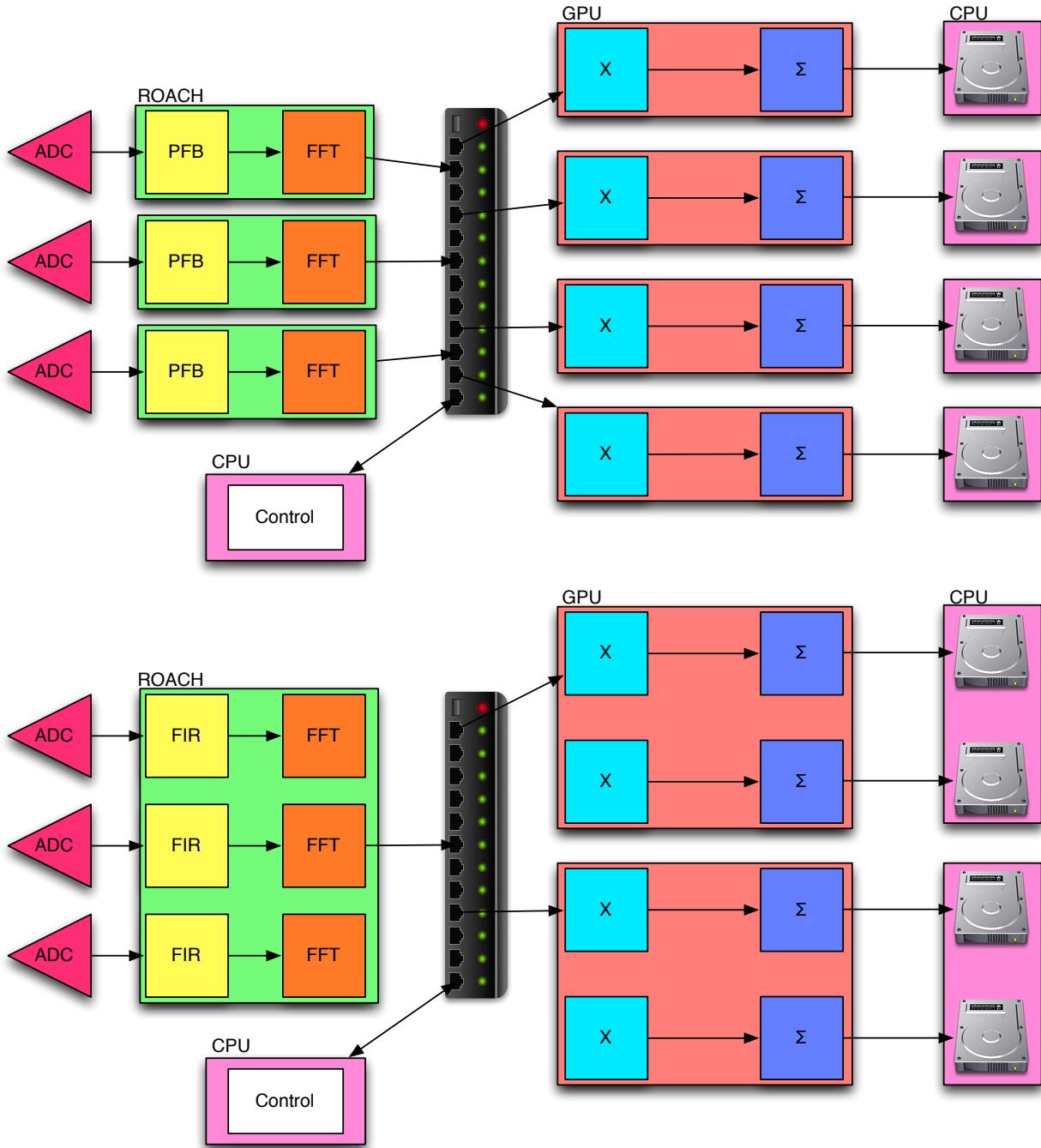


Figure 4.8: Two potential mappings for the FX Correlator

blocks must run through the switch. This is represented in the diagram by a single connection from each ROACH and GPU board to the switch.

4.5 Code Generation

The last step in the process is the implementation the instrument. This step is optional, only being used when the mapped dataflow is being used to design an instrument rather than assess cost. Since the toolflow relies on established libraries, implementation only consists of stitching together existing routines.

In addition to this, common design patterns can be stitched together into parameterized implementations. This style of instrument design greatly accelerates time to science for many projects. Separating the implementation of the instrument from the hardware specification has created a design that works well for a variety of computational resources and applications. As resources improve, the instrument can improve along with them, providing the opportunity to do new science that wasn't possible on older technology. We have shown this is possible by implementing the Packetized Astronomy Signal Processor or PASP using the CASPER toolflow. [10]

Chapter 5

Algorithm Partitioning

After coming up with an instrument description, it is necessary to determine how that instrument will be implemented in hardware. I use Integer Linear Programming (ILP) to model and solve this problem. As described in Section 3.2, ILP is a powerful technique for defining and solving optimization problems. In this chapter I explain how the ILP is defined based on the dataflow model and the techniques used to make sure the program can be solved quickly.

5.1 Variables

The variables in the ILP are represent the optimal mapping for the system. Ultimately, the ILP determines which platforms should be used, and what part of the algorithm should get implemented on each platform. This is achieved by having the ILP consider some platform, and assume it can instantiate at most n_p copies of that platform. We will call a copy of the platform a board. Each board must have some variables that determine which computation blocks get mapped to it. For some board i , the number of computation blocks of type b that get mapped to it is represented by the variable $n_{i,b}$.

A solution to the ILP, with each of the $n_{i,b}$ variables filled in, gives a complete specification of the optimal mapping for that instrument.

5.2 Constraints

The constraints serve two purposes. First, they ensure that no resource is overmapped, so that the amount of hardware the ILP generates will be sufficient to do the computation required. Second, they make sure that the correct design gets implemented.

Platform Resources

Any resource some board that gets used by mapping computation blocks to that board must be accounted for in the linear program. This is abstracted in the ILP using by adding single constraint for each resource.

Resource Limitations

For some resource, r , we use our performance model of each block to assess how much of the resource is used up by each block type. The percentage of the resource r required by some block type, or utilization, is represented by a constant (not a variable), $r_{p,b}$, where p represents the platform, and b represents the block type. Multiplying the resources required for a specific block type by the number of blocks needed on that specific board determines the total percent of resource r block type b will require on the board. Summing over all of the block types determines what percent of resource r is used in the final design, which gives us the final format of the constraint needed to ensure that some resource r is not overmapped on board i :

$$\sum_{b \in \text{Blocks}} n_{i,b} r_{p,b} \leq 1 \quad (5.1)$$

Each resource will require a separate constraint in the ILP of this form. By ensuring that the total utilization of each resource required is less than 100%, we guarantee that there are enough resources to allow all the blocks mapped to that board to complete their tasks.

Dataflow Model Constraints

After getting assurance that no resources are overused, it is important to verify that the correct design was implemented. The resource utilization on the value of $n_{i,b}$, but there are is an additional constraint on these variables. Namely, that the correct number of blocks actually get implemented. This adds a simple constraint to each block type:

$$\sum_{\substack{\text{board} \\ \text{in boards}}} n_{\text{board}, \text{blocktype}} == n_{\text{blocktype}} \quad (5.2)$$

The total number of blocks of a certain type should be equal to the number of blocks of that type we actually need in the design.

Network Resources

The ILP does not design the network topology. While it would be possible to design the network using an ILP, this would add unnecessary complexity to the program (therefore increasing the runtime) with little gain. As described in Section 3.1, most radio astronomy

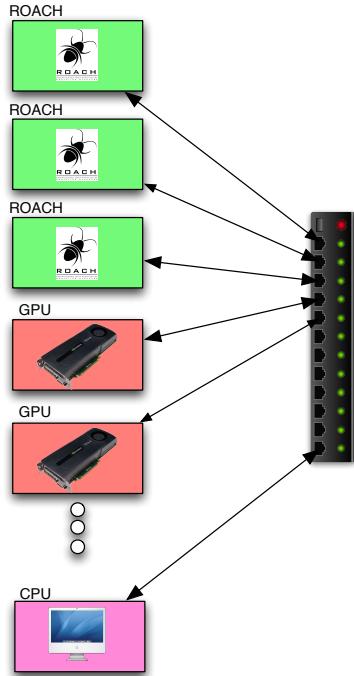


Figure 5.1: Full-Crossbar Interconnect

applications require a full-crossbar interconnect at some point, because many computational blocks have an all-to-all or one-to-all communication pattern. Rather than have the ILP redesign the same topology over and over, we simply assume this interconnect exists and every board can communicate with every other board directly over a fixed-bandwidth link.

Figure 5.1 shows an example of this topology. Each board gets connected to the same switch and can communicate with any other board on the switch, regardless of the platform type.

Bandwidth Limitations

Even a fixed network topology, there still are communication constraints that must be taken into account in the ILP. While there might be a link available between each pair of boards, the bandwidth into and out of these boards is limited. Consider Figure 5.1 again. Suppose every other board in the cluster needed to stream 10 Gbps of data to the CPU board, but it is only connected to the switch via a single 10 Gbps link. There are additional constraints to ensure that the input and output bandwidths are not exceeded.

In order to write these constraints, we must first determine how many blocks need to communicate with a block that is not located on the same board. We introduce new variables $nr_{i,b}$ to represent the number of blocks of type b on board i that need to receive data from the cluster, and, similarly, $ns_{i,b}$ to represent the number of blocks of type b on board i that

need to send data to the cluster. Given the amount of data some computational block type takes as input and the number of those blocks on the board, we can multiply them together to determine the amount of input bandwidth that computational block type will require. Summing over every block type determines the total amount of input bandwidth needed by all the computational blocks on the board, creating a constraint that the total required bandwidth must be less than or equal to the total available input bandwidth. The constraint on output bandwidth is calculated the same way, generating a pair of constraints for each board, one restricting the total amount of input bandwidth, and another restricting the total amount of output bandwidth.

$$\sum_{b \in \text{Blocks}} nr_{i,b} bw_in_b \leq bw_in_p \quad (5.3)$$

$$\sum_{b \in \text{Blocks}} ns_{i,b} bw_out_b \leq bw_out_p \quad (5.4)$$

Connection Constraints

First, we observe that these variables must be bounded by 0 and $n_{i,b}$, since there cannot be a negative number of blocks that need to communicate, and the number of blocks of type b that need to communicate can't exceed the number of blocks physically on the board. Next, we must take into account the structure of the algorithm to determine whether or not a given block needs to communicate with a separate board.

While the constraints on the total input and output bandwidth might seem simple, ensuring the values for $nr_{i,b}$ and $ns_{i,b}$ are sane requires additional constraints. Suppose we know we have 2 types of blocks: A , and B . A is a source of data, meaning it does not data from any computation block. Similarly, block B is a sink, with no data to send to another block. Regardless of how many A and B blocks get placed on platform i , none of the A blocks will need to receive data and none of the B blocks will need to send data. Knowing that A is a source and B is a sink tells us that $nr_{i,A} = 0$ and $ns_{i,B} = 0$.

In order to appropriately define the linear program, it is first important to look at the different ways the computational blocks in a design may need to communicate, and create appropriate constraints. We must revisit the connection types introduced in Section 4.3, and determine how the different types of links affect the linear program.

When two blocktypes are linked via a ‘one-to-one’ connection, communication is required when the number of A blocks is different than the number of B blocks on a single board. When there are more A blocks than B blocks, $n_{i,A} > n_{i,B}$, the number of A blocks that need to send data to the cluster is $n_{i,A} - n_{i,B}$ and none of the B blocks on that board need to receive data from the cluster. In the opposite case, $n_{i,A} < n_{i,B}$, and the none of the A blocks need to send data to the cluster, but $n_{i,B} - n_{i,A}$ blocks of type B will need to receive data from the cluster. Both of these cases are captured by the same pair of constraints:

$$ns_{i,A} \geq n_{i,A} - n_{i,B} \quad (5.5)$$

$$nr_{i,B} \geq n_{i,B} - n_{i,A} \quad (5.6)$$

When $n_{i,A} - n_{i,B}$ is non-negative, we are guaranteed that we will not underestimate $ns_{i,A}$, and when $n_{i,A} - n_{i,B}$ is negative, $ns_{i,A}$ will be forced to at least 0 because of the lower limit on the variable.

Setting the ns and nr variable is a little more complicated for the ‘all-to-all’ case, and requires the implementation of some conditional logic in the linear program. When any of the B blocks are not on the board i , then every A block must send data to the cluster, and $ns_{i,A} = n_{i,A}$. Otherwise, none of the A blocks need to send and $ns_{i,A} = 0$. Similarly on the receive side, if any of the A blocks are not of board i , $nr_{i,B} = n_{i,B}$, otherwise $nr_{i,B} = 0$. The conditional logic is easily implemented in an ILP, as described in X.

When a block must send to or receive data from multiple different block types, the constraint is the same way an ‘all-to-all’ connection gets constrained. If any of the blocks it must link to reside outside the board, it is assumed that all of the data must be sent over the link to the cluster.

Implementing the ILP this way results in an overestimation of the required bandwidth. In the case where $ns_{i,A} = n_{i,A}$, it’s true that every block of type A will need to send some data. However, they might not need to send the full bandwidth bw_out_A to the switch, since some portion of the data sent by an A block may be required by B blocks residing on the same board. A more exact version of this calculation would also take into account the smaller bandwidth, but due to the complexity and rarity of this case the approximation is sufficient. As described in Section 3.1, many architectures have this type of connection but there are very few cases where the A and B blocks connected this way reside on the same board.

5.3 ILP Implementation

In software, the ILP can be generated automatically using the dataflow model contained in an instrument object. The generic Instrument class has a single function, called `runILP()` that iterates through the dataflow model, creating variables and constraints, determines the cost model depending on what the user wants to optimize for and runs an ILP solver to generate an optimal mapping. Adding on to the instrument creation example in Section 4.2, we show how the user can create and map an instrument using only two function calls in the following code.

```
antennas = 16
channels = 128
bandwidth = 0.4 #defined in GHz
```

```
#create the instrument
```

```
myfxcorrelator = FXCorrelator( antennas , channels , bandwidth )
myfxcorrelator . runILP ()
```

5.4 Performance Modeling

The data that the ILP uses to measure resource utilization must come from a preexisting performance model. These models can take a number of forms. Benchmarks of compiled and running code provide the best information, but are also time consuming to obtain if they don't already exist and they require a real implementation of the block. Estimates are faster to obtain but won't be as accurate. Nevertheless, many DSP blocks have predictable performance and a performance estimate based on similar benchmarks can be very reliable.

Benchmarks of the mapped design or running code serve as a very accurate way to assess performance. Figure 5.2 shows the type of benchmarks that would be useful for an FPGA block, measuring utilization of available FPGA resources. The graphs show the utilization of registers, LUTs, BRAMs and DSPs used on the Virtex 5 SX95T. The top graph is the data for an 800 MHz FIR filter with 4 taps, and the bottom shows utilization for an 800 MHz FFT. Getting the performance model for a specific block simply requires a looking up the FFT size in the graph.

Similarly, Figure 5.3 has benchmarks for the same blocks, FIR filter above and FFT below, but these are tested on a GTX 580 GPU. These benchmarks measure the runtime of each function.

Many papers also provide these kinds of benchmarks, making it easy to get accurate numbers without installing or running the code. The xGPU paper has a number of graphs demonstrating kernel performance that can be used directly as an ORCAS benchmark, which will be shown in Chapter 6.

Performance data can also be represented using formulas. Figure 5.2 clearly shows some predictable trends in the FIR and FFT utilization. Primiami, et al turned this predictability into a set of equations that determine the requisite resources with a set of formulas [21].

We can also use existing benchmarks to project how a block might perform on newer technology. In FPGAs, we notice the number of resources required for a block remains nearly constant between different chips. Table 5.1 shows the resource utilization of an 800MHz 32k channel FFT on three chips from three different generations, the Virtex 5 SX95T, Virtex 6 SX475T, and Virtex 7 VX980T. Aside from the number of LUTs, which slightly dropped between the Virtex 5 and Virtex 6, the number of resources required remains very stable across chips. Although the utilization will change, because different chips will provide different amounts of resources, recording the number of resources required on one chip makes it possible to predict the utilization on another chip.

Similarly projections can be made with new processor technology. Even if a new technology is not yet available to buy, a conservative and optimistic estimate of the speedup can be used to generate a conservative and optimistic estimate of the instrument cost using that new technology.

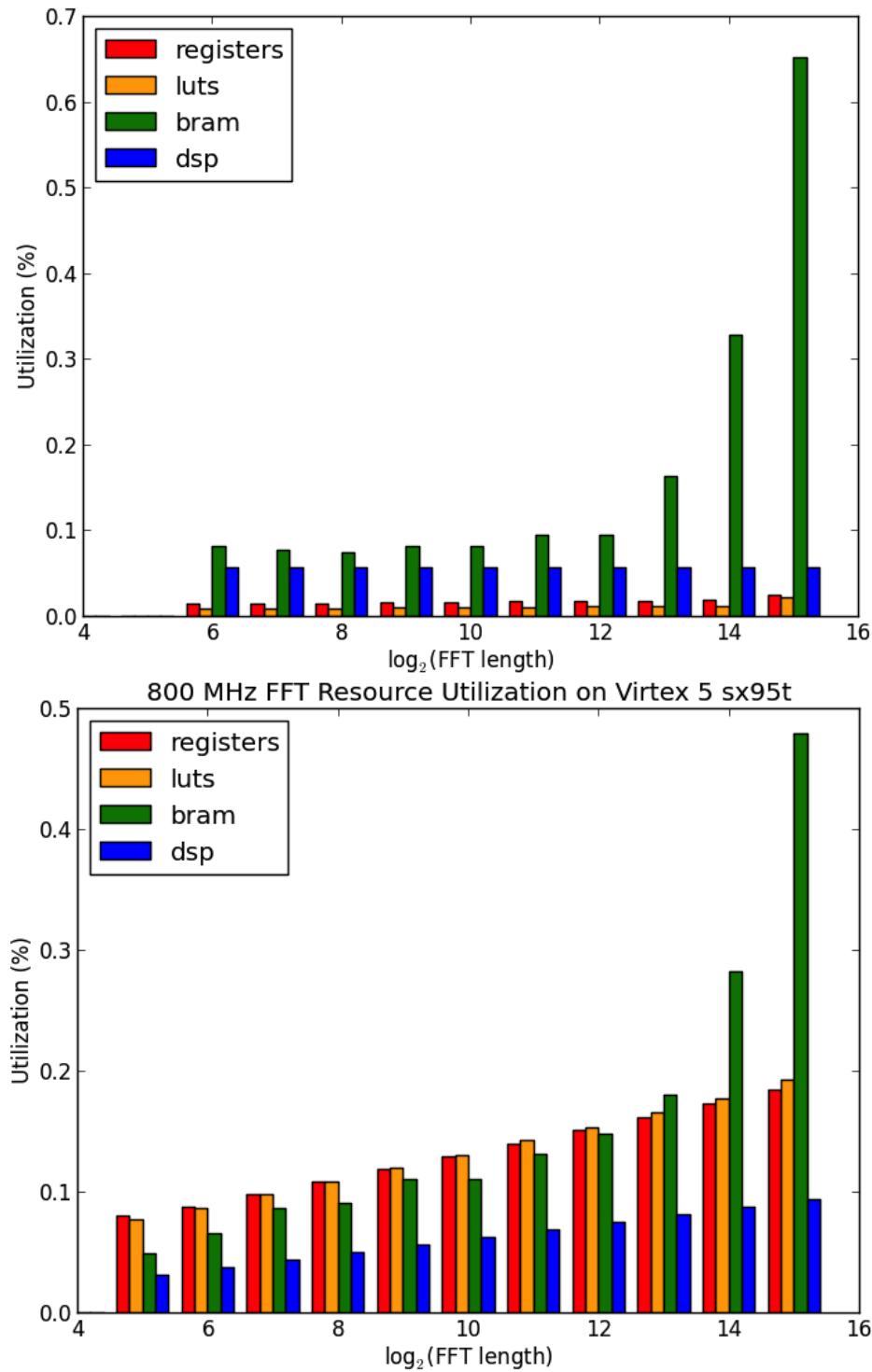


Figure 5.2: PFB FIR and FFT benchmark data on the Virtex V SX95T

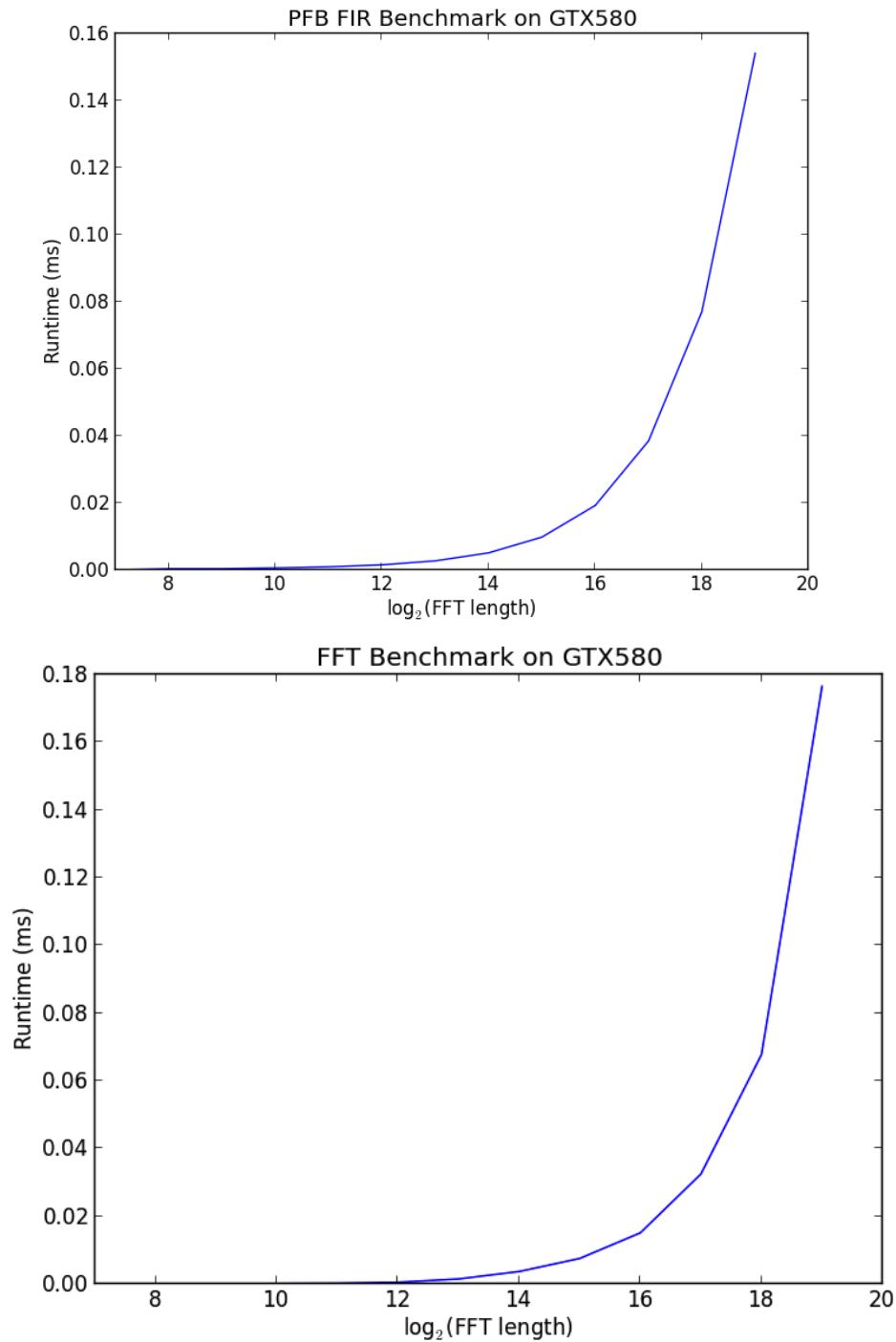


Figure 5.3: PFB FIR and FFT benchmark data on the GTX 580

| Resource | Virtex 5 SX95T | Virtex 6 SX475T | Virtex 7 VX980T |
|--------------|----------------|-----------------|-----------------|
| Registers | 10881 | 10789 | 10788 |
| LUTs | 11358 | 9632 | 9773 |
| 36k BlockRAM | 100 | 100 | 100 |
| 18k BlockRAM | 30 | 30 | 30 |
| DSPs | 60 | 60 | 60 |

Table 5.1: Comparative Resource Utilization of a 32k Channel 800 MHz FFT

| Platform | Specification | Monetary Cost | Idle power | Average power | Maximum power |
|-------------|-----------------|---------------|------------|---------------|---------------|
| ROACH | Virtex 5 SX95T | \$6,700 | 55 W | 65 W | 75 W |
| ROACH 2 | Virtex 6 SX475T | \$10,500 | 60 W | 70 W | 80 W |
| ROACH 3 | Virtex 7 VX980T | | 60 W | 70 W | 80 W |
| NRAO Server | GTX 580 | \$3,500 | 225 W | 400 W | 475 W |

Table 5.2: Monetary and Power Costs for Common CASPER Platforms

| GPU | Cost | Idle power | Average power | Maximum power |
|---------|---------|------------|---------------|---------------|
| GTX 580 | \$500 | 125 W | 150 W | 175 W |
| GTX 680 | \$500 | 100 W | | |
| GTX 690 | \$1,000 | | | |

Table 5.3: GPU Board Costs

5.5 Cost Modeling

The cost function in the linear program can represent a number of properties like monetary cost, power, development time or rack space. In this work I focus on monetary costs and power. Both can be calculated simply by iterating through the boards used and adding the cost of the board to the total cost. Table 5.2

The model can also take into account donated or existing hardware that will not contribute to the monetary cost of an instrument. Adding another platform, like a 'Free ROACH', with same specifications as a ROACH but a monetary cost of \$0 will allow the model to use the hardware without incurring any cost. This feature is a useful way to determine if it is worthwhile to replace existing hardware for an instrument upgrade.

The final entry in Table 5.2 is a high performance server, and the reported costs include a GTX 580 GPU, but that might not be the best GPU to use. The data in Table 5.3 is useful to determine how a different GPU will affect the total cost of the server.

5.6 Optimization

While Integer Linear Programming has a number of desirable properties, the lack of an efficient algorithm to solve it constitutes a significant obstacle in designing an ILP with reasonable performance. This section describes how the solver selection, design of this ILP, and the introduction of a few extra constraints serve to improve the performance and scalability of the program defined in this chapter.

ILP Solver Selection

The ILP is described using an open source Python package called PuLP, available at <http://www.coin-or.org/PuLP/>. PuLP does not include an integer linear program solver. Instead, it supports a number of existing solvers, allowing the user to choose which one to use.

ORCAS was originally tested using the GNU Linear Programming Kit or GLPK [] a free open source package that is supported by PuLP. Unfortunately, as the instrument models became more complex, GLPK often failed to converge on an optimal solution after running for 6 hours on my personal laptop, a 2011 Macbook Air and an attempt to get better performance by using a high powered server was futile.

Because PuLP makes it simple to change the solver, only requiring a change to the single line of code that calls the solver, I decided to test other solvers before editing the ILP. Another solver was chosen by referring to a set of integer linear programming benchmarks publish by Koch et al.[16]. These benchmarks show the feasibility of an ILP is highly sensitive to the solver used. Recent results from those benchmarks, available online [18], found that the Gurobi solver [13] has a relatively high number of successes. Gurobi was able to solve the existing models within minutes and is the solver used to provide all the results in this work.

Guided Optimization

Another solution relies on user aid to guide the mapping. The ILP may spend time going over solutions that are obviously wrong to a human user. In this case, the user could intervene by setting some of the ILP variables manually and letting the ILP find a solution for the remaining variables. While this may be a feasible solution for a computer expert who might have some idea of what the optimal mapping should be, this not a useful technique for the domain specific experts who are as familiar with the hardware and computational block implementations. This violates one of the basic goals of this tool described in Section 4.1. The tool needs to be accessible and usable by domain specific experts as well as computer experts, and dealing with symmetry in this way will require a computer expert in the loop to generate a mapping and get a cost estimate. To maintain usability for domain experts, guided optimization will not be the sole solution to this issue.

Breaking Symmetry

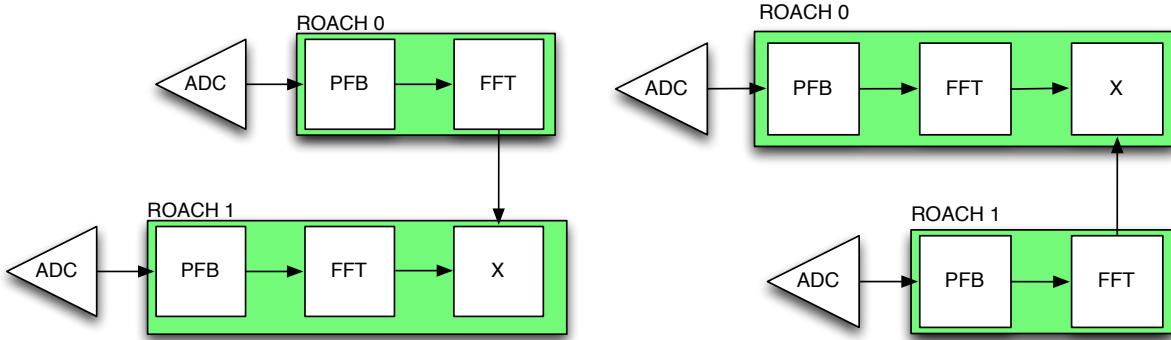


Figure 5.4: TODO

In this type of program, symmetry significantly increases the amount of time required to confirm an optimal solution. The boards that have the same platform type are interchangeable, so if board i implements some design and board j implements a different design in the optimal solution, there is another optimal solution where their designs are swapped. For example, suppose 2 ROACH boards are available to implement a FIR filter and an FFT. The ILP might observe that both blocks cannot fit on a single board and assigns the FIR to the first ROACH and the FFT to the second ROACH. This obviously seems like an optimal solution, but the ILP may also need to check the case where the FFT is placed on ROACH_0 and the FIR is on ROACH_1, only to find that it has the same cost as the previous result. In this simple example there was only one other solution to search, but as the ILP and the search space grows the number of solutions that are symmetric to the optimal case will also grow.

Searching symmetric solutions can become a major time sink, because the ILP solver may find an optimal solution early on, but will require a long time to confirm that it is actually the optimal result, spending time going over many other solutions that are isomorphic to the first one.

This returns the current best result the solver knows of, but it cannot guarantee that the solution is globally optimal or, in the case where it is not the a globally optimal mapping, determine if it is close to the optimal solution, (determining that is analogous to solving the ILP). Early stopping works well when it's clear that symmetry is the cause of the long runtime and the amount of time it would take to get one of the isomorphic optimal solutions is short. Even so, the lack of predictable and repeatable results makes early stopping an unappealing solution.

The solutions above outline ways to cope with the existing symmetry. Another way to reduce the runtime is to remove the symmetry altogether. In order to do this, the ILP must be modified so that only 1 of the isometric optimal solutions is a valid solution to the ILP. First, a variable lex_order_i is added for each board. This variable is meant to uniquely

identifies the design running on the board; it is simply the concatenation of all the $n_{i,b}$ variables for that board. Note that the ordering of $n_{i,b}$ variables when the concatenation is done is irrelevant, the only thing that matters is that the order is consistent for every board. When $\text{lex_order}_i = \text{lex_order}_j$ we can infer that for all blocks b , $n_{i,b} = n_{j,b}$. Otherwise, there must be some block b where $n_{i,b} \neq n_{j,b}$. Now that the designs can be identified, they can be ordered. They are simply ordered lexicographically, by adding the constraints in Equation 5.7 for every $i \geq 1$

$$\text{lex_order}_{i-1} \geq \text{lex_order}_i \quad (5.7)$$

While this make the ILP, adding both constraints and variables, it reduces the amount of time the solver takes to find a solution. This lexicographic ordering makes it impossible to swap designs between different blocks, resulting in unique and valid mappings.

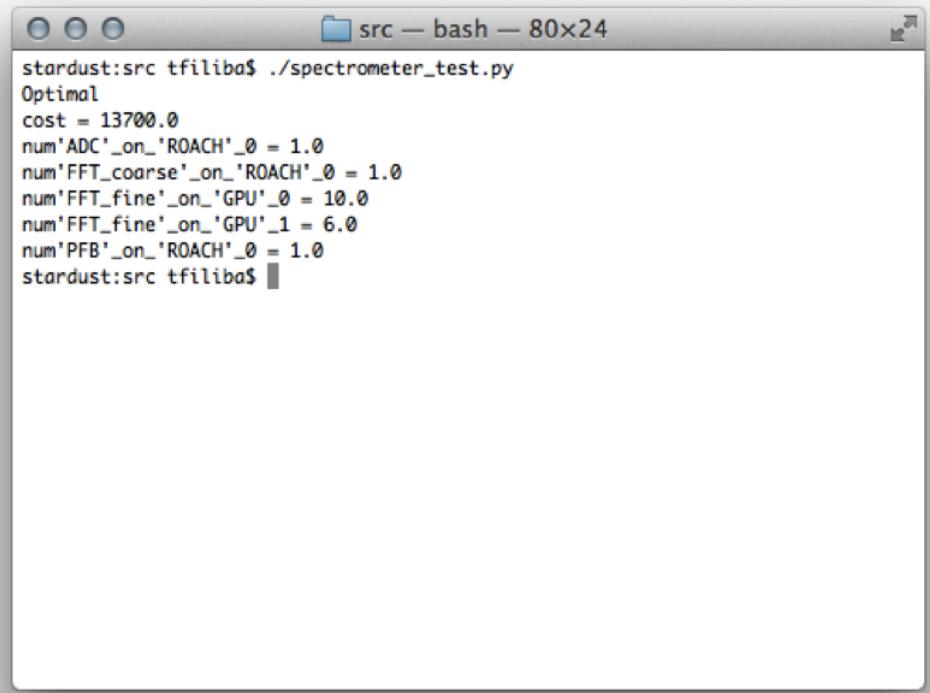
Revisiting the symmetry example at the beginning of this section, the design with 1 FIR and no FFTs would be encoded with a $\text{lex_order} = 10$, and the design the no FIRs and a single FFT would get the encoding $\text{lex_order} = 01$. When the FIR is placed on ROACH_0, then $\text{lex_order}_0 = 10 \geq \text{lex_order}_0 = 01$, satisfying the new constraint. The solution where the blocks are swapped and FIR is on ROACH_1 violates the new constraint $\text{lex_order}_0 = 01 \not\geq \text{lex_order}_0 = 10$, and will not be considered by the ILP solver.

Generalizing this, it is impossible to take a valid solution (with the lexicographic constraint) and get another valid solution by swapping distinct designs between boards. Suppose, without loss of generality, board i has a design with $\text{lex_order}_i = x$ and board j has a distinct design with $\text{lex_order}_j = y$ and $i < j$. Knowing that the design is valid implies $x \geq y$. Another optimal mapping exists where the designs are swapped and $\text{lex_order}_i = y$, $\text{lex_order}_j = x$, but we are guaranteed that this is not a valid solution to the ILP because it violates the lexicographic ordering constraint.

These constraints have been implemented in the final ILP and drastically reduces the amount of time it takes to solve the ILP. The additional constraints do not change the cost of the optimal solution, instead they just reduce the number of valid optimal solutions. By modifying the ILP, the performance is greatly improved without sacrificing optimality or usability.

5.7 Final Mapping

The final mapping produced by the ILP is just a list of variable indicating which blocks go on which boards and the cost of the design. Figure 5.5 shows an example of the output produced by ORCAS. The design is a simple wideband spectrometer model. The mapping indicates that the design will use one ROACH board and two GTX 580 servers, placing the FIR and coarse FFT on the FPGA and the remaining fine FFTs on the GPU.

A screenshot of a terminal window titled "src — bash — 80x24". The window contains the following text:

```
stardust:src tfiliba$ ./spectrometer_test.py
Optimal
cost = 13700.0
num'ADC'_on_-'ROACH'_0 = 1.0
num'FFT_coarse'_on_-'ROACH'_0 = 1.0
num'FFT_fine'_on_-'GPU'_0 = 10.0
num'FFT_fine'_on_-'GPU'_1 = 6.0
num'PFB'_on_-'ROACH'_0 = 1.0
stardust:src tfiliba$
```

Figure 5.5: ORCAS Output

Chapter 6

Analysis

The ORCAS tool is analyzed by looking at three case studies. We created three instrument types: Spectrometer, High Resolution Spectrometer, and FX Correlator, as defined in Sections 2.1, 2.2, and 2.4. Each case study was chosen to illustrate a different aspect of the toolflow. The spectrometer gives an example of the end to end toolflow using a simple dataflow, the high resolution spectrometer allows us to explore tradeoffs in the design space and the FX correlator shows how the tool behaves when designing very large instruments.

6.1 Spectrometer Case Study

The spectrometer is a simple instrument, making it easy to follow the end to end toolflow. We design a 50 MHz and an 800 MHz spectrometer that breaks the band into 1024 channels.

Spectrometer Definition

Defining a simple spectrometer requires very few parameters. First, as with most instruments, the astronomer must specify the sky bandwidth the instrument must process, defined in MHz and the number of bits in each ADC sample. Then, the desired spectral resolution is defined in MHz per channel, or analogously, the number of channels that should be used to break up the bandwidth. Finally, the integration time needs to be defined.

One optional parameter, number of antennas, can also be defined. This describes the number of independent spectrometers that need to be created. While this parameter does not affect the end to end processing for each antenna, knowing how many spectrometers are needed allows for more efficient use of the hardware.

The 800 MHz spectrometer is created simply by defining the parameters and instantiating a Spectrometer object as follows:

```
# 800MHz spectrometer
numchannels = 1024
accumulation_length = 10
```

```

bandwidth = 0.8 #defined in GHz
input_bitwidth = 8
fft_out_bitwidth = 4

#create the instrument
myspectrometer = Spectrometer(numchannels, accumulation_length,
    bandwidth, input_bitwidth, fft_out_bitwidth)

```

Spectrometer Dataflow

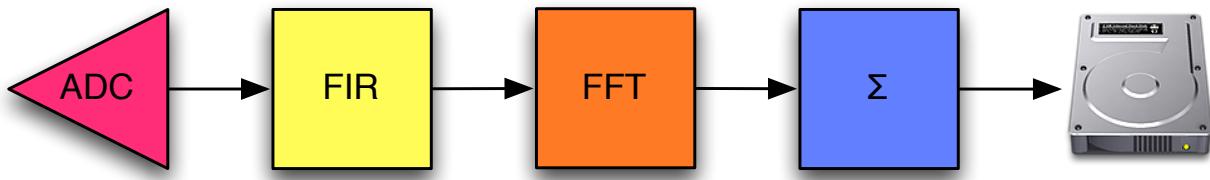


Figure 6.1: General Spectrometer Dataflow Model

The spectrometer instrument definition generates a very simple dataflow. Figure 6.1 shows the general dataflow model for a single antenna spectrometer. This model can be applied to any spectrometer, as the spectrometer parameters do not affect how many computational blocks are required or the interconnect layout. The ADC feeds data into a FIR filter. Then the filtered signal is transformed into channels in the FFT and those channels are accumulated and saved to disk. Regardless of the parameters the astronomer chooses, the dataflow will be the same.

The parameters for each block come directly from the instrument definition. The FIR parameters come from the number of FIR taps and window shape, the FFT is simply defined by the FFT length parameter and the accumulator also is parameterized by the FFT length as well as the integration time.

```

# add the ADC
adc_bw = bandwidth*input_bitwidth
self.blocks[ 'ADC' ] = CBlock( CBlock.getADCModel( self.platforms ,
    bandwidth, input_bitwidth ), -1, 0, 0, 'PFB' , 0, adc_bw, antennas )
self.totalblocks += antennas

# add the PFB
self.blocks[ 'PFB' ] = CBlock( CBlock.getPFBModel( self.platforms ,
    input_bitwidth, numchannels ), 'ADC' , 0, adc_bw, 'FFT' , 0, adc_bw ,
    antennas )

```

```

self.totalblocks += antennas

# add the FFT
fft_out_bandwidth = bandwidth* fft_out_bitwidth
self.blocks[ 'FFT' ] = CBlock( CBlock.getFFTModel( self.platforms ,
    numchannels ), 'PFB' ,0 ,adc_bw , 'VAcc' ,0 ,fft_out_bandwidth ,antennas
)
self.totalblocks += antennas

#add the Vacc
self.blocks[ 'VAcc' ] = CBlock( CBlock.getVAccModel( self.platforms ,
    fft_out_bitwidth , accumulation_length ), 'FFT' ,0 ,
    fft_out_bandwidth ,-1 ,0 ,0 ,antennas )
self.totalblocks += antennas

```

Spectrometer Mapping

Optimal
cost = 6700.0
is_used_ 'ROACH' _0 = 1.0
num'ADC' _on_ 'ROACH' _0 = 1.0
num'FFT' _on_ 'ROACH' _0 = 1.0
num'PFB' _on_ 'ROACH' _0 = 1.0
num'VAcc' _on_ 'ROACH' _0 = 1.0

Optimal
cost = 5500.0
is_used_ 'GPU' _0 = 1.0
is_used_ 'IBOB' _0 = 1.0
num'ADC' _on_ 'IBOB' _0 = 1.0
num'FFT' _on_ 'GPU' _0 = 1.0
num'PFB' _on_ 'GPU' _0 = 1.0
num'VAcc' _on_ 'GPU' _0 = 1.0

6.2 High Resolution Spectrometer Case Study

High Resolution Spectrometer Definition

The main difference between a spectrometer and a high resolution spectrometer is the need for two stages of channelization rather than just one. The sky bandwidth, integration time, and number of antennas are defined in the same way as the previous spectrometer type.

The spectral resolution is defined differently, because both the coarse and fine resolutions need to be defined. The coarse resolution defines how many channels the whole sky bandwidth should be broken up into initially. The fine resolution defines how many channels each coarse channel is broken into. Both can be described in MHz per channel.

High Resolution Spectrometer Dataflow

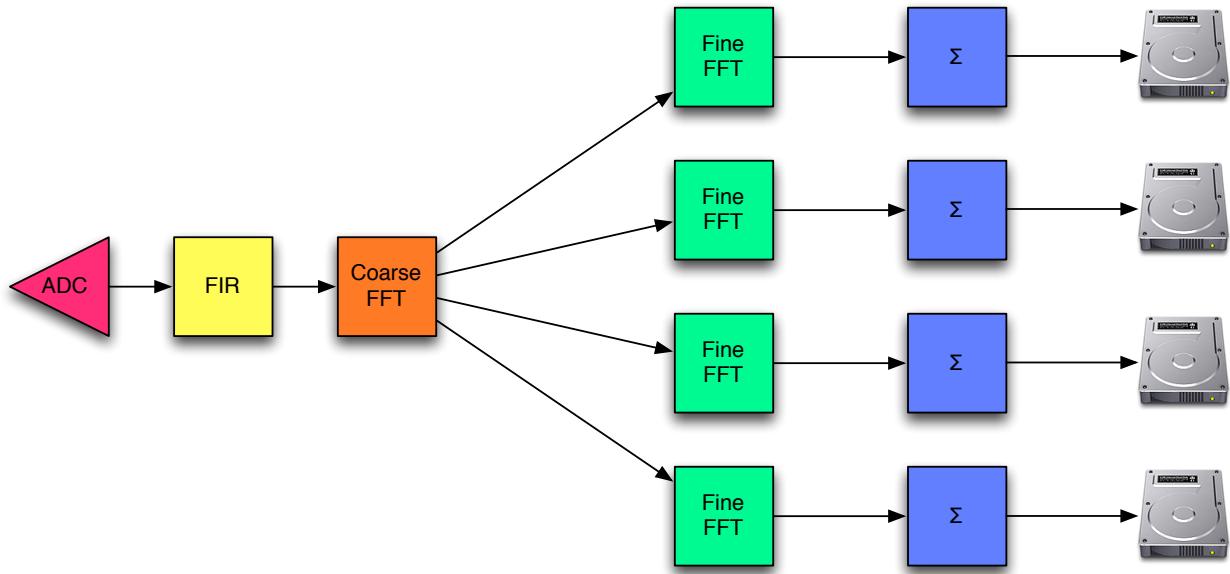


Figure 6.2: Example High Resolution Spectrometer Dataflow Model *This figure shows an example dataflow for a high resolution spectrometer with 4 coarse FFT channels. Like the previous spectrometer, the data comes in through an ADC, is filtered and then channelized using a coarse FFT. Then, to achieve a higher resolution, each coarse channel is divided into sub-channels using the 4 fine FFT blocks in the dataflow.*

The high resolution spectrometer dataflow does depend on the parameters specified in the instrument description. An example dataflow is shown in Figure 6.2. The first three blocks in the dataflow are exactly the same as the Spectrometer dataflow described in the previous section. An ADC feeds data into a FIR filter followed by an FFT. After the FFT, the algorithm is modified to accommodate the higher resolution required. The first FFT divides the band into a number of coarse channels and then each coarse channel must be further divided into a number of fine channels. The coarse FFT must feed its data to a separate fine FFT for each coarse channel, so the number of fine FFTs in the dataflow diagram will vary based on the number of coarse channels. At this point, each coarse channel is processed in an independent pipeline the finely channelized data is accumulated and recorded to disk. The example in Figure 6.2 shows a spectrometer that divides the data into 4 coarse channels.

results for hi res spectrometer (gbt and seti)
 Serendip 6 300MHz 7 ant 2 pol
 1Hz resolution (256 million channels)
 Greenbank 2.5GHz 1 beam 2 pols
 1 Hz resolution
 Same benchmarks as before, just discuss large bw, multi stage fft

Cost

Power

Algorithmic Exploration

6.3 FX Correlator Case Study

FX Correlator Definition

An FX Correlator is also defined by the amount of bandwidth it processes, number of channels, and integration time, but now the number of antennas is a necessary parameter.

FX Correlator Dataflow

The FX dataflow model is based on the algorithm used by the CASPER correlator described in Section 3.1. The processing model is described by replicating 2 basic pipelines, called an F-Engine and an X-Engine. The number of times each pipeline needs to be replicated depends on the number of antennas and number of channels this correlator requires.

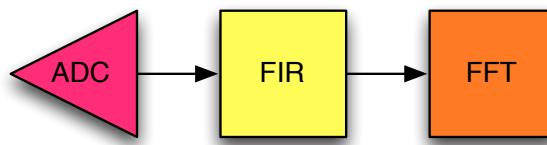


Figure 6.3: FX Correlator F-Engine Model

An F-Engine, pictured in Figure 6.3 is responsible for channelizing the data from a single antenna. It takes in data from an ADC, and channelizes the data using a FIR and FFT to create a polyphase filter bank. The number of F-Engines in the correlator dataflow will vary with the number of antennas.

The second pipeline, the X-Engine processes the channelized data. Each X-Engine takes a single channel of data from every antenna in the array, cross-correlates the data, accumulates

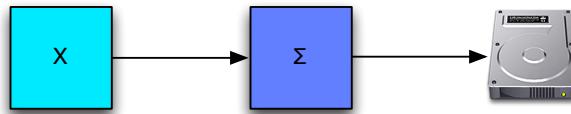


Figure 6.4: FX Correlator X-Engine Model

it and stores it to disk. Figure 6.4 shows the pipeline for a single X-Engine. Since each X-Engine only operates on a single channel, the total number of X-Engines in the correlator must be the same as the number of channels in the FFT.

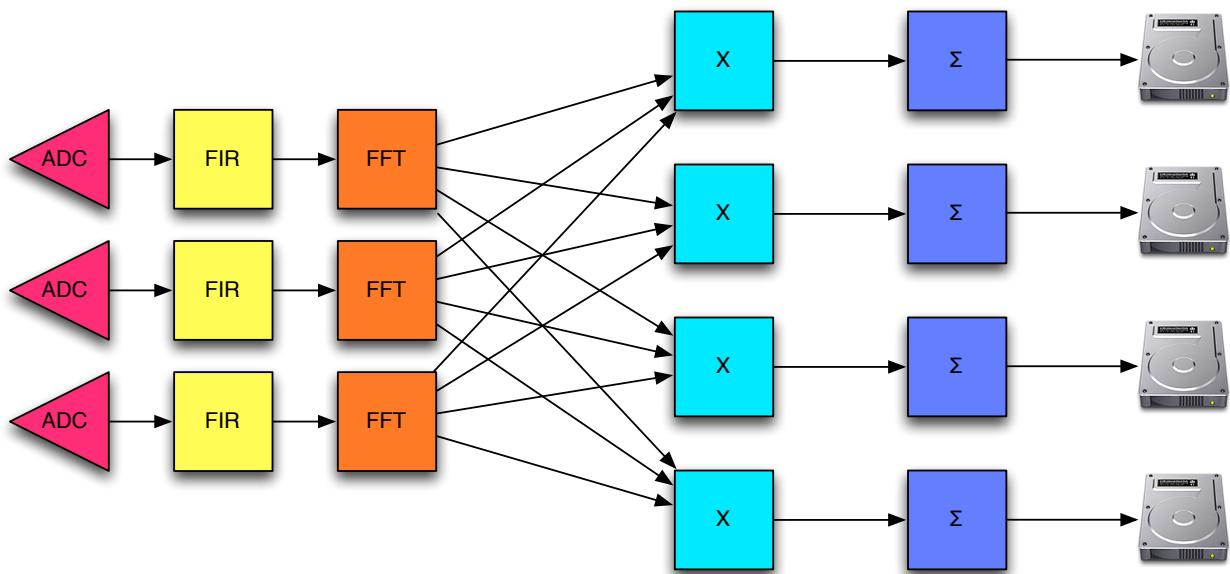


Figure 6.5: Example FX Correlator Dataflow Model

The dataflow for an FX correlator will vary quite a bit based on the input parameters. Figure 6.5 shows an example there antenna four channel FX correlator. The left half of the figure has three F-engines, for each of the three antennas. The right half has four X-engines, one for each channel. And in the center, since each X-engine requires data from every F-engine, the cross-correlation blocks, represented by an X and the FFT blocks are connected in an all-to-all configuration.

xgpu benchmarks, correlator placement

Need benchmarks for xengine

HERA Correlator

572 ant dual pol (currently can run in under 5 minutes)

100 MHz

1024 channels

FPGA F, xGPU X
LWA @ OVRO
10k ant dual pol
120k signals
60MHz
4096 channels
1MWatt est w/GPU x-engine
SKA
4000 ant (8k signals)
10 GHz
16k channels

Chapter 7

Conclusions

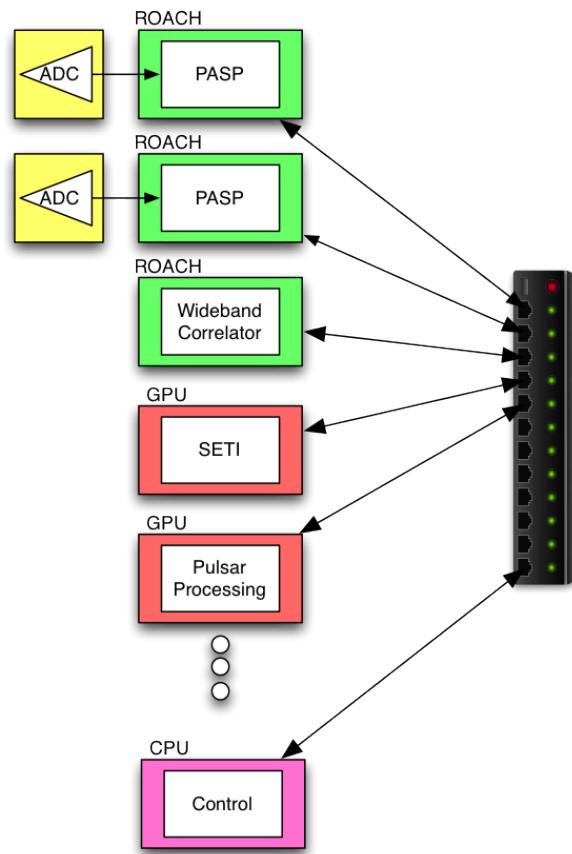


Figure 7.1: A potential architecture for multiple scientific instruments simultaneously processing data from the same telescope

7.1 Future Work

Bibliography

- [1] M M Bailes et al. “Transformation of a star into a planet in a millisecond pulsar binary.” In: *Science* 333.6050 (Sept. 2011), pp. 1717–1720.
- [2] F. Balarin et al. “Metropolis: an integrated electronic system design environment”. In: *Computer* 36.4 (2003), pp. 45–52.
- [3] J Bunton. “ALMA Memo 342-An Improved FX Correlator”. In: *CSIRO, Telecommunications and Industrial Physics, Australia* (2000).
- [4] M A Clark, Pc La Plante, and L J Greenhill. “Accelerating radio astronomy cross-correlation with graphics processing units”. In: *International Journal of High Performance Computing Applications* 27.2 (May 2013).
- [5] Abhijit Davare. “Automated Mapping for Heterogeneous Multiprocessor Embedded Systems ”. PhD thesis. University of California, Berkeley, 2007.
- [6] A T Deller et al. “DiFX: A Software Correlator for Very Long Baseline Interferometry Using Multiprocessor Computing Environments”. In: *Publ Astron Soc Pac* 119.853 (Mar. 2007), pp. 318–336.
- [7] P Demorest et al. “Gravitational Wave Astronomy Using Pulsars: Massive Black Hole Mergers & the Early Universe”. In: *Astro2010: The Astronomy and Astrophysics Decadal Survey 2010* (2009), p. 64.
- [8] Douglas Michael Author Densmore et al. “Metro II Execution Semantics for Mapping” . In: () .
- [9] Sheperd S Doeleman et al. “Jet-Launching Structure Resolved Near the Supermassive Black Hole in M87”. In: *Science* 338.6 (Oct. 2012), p. 355.
- [10] T Filiba and D Werthimer. “Automatic generation of heterogeneous spectrometers for radio astronomy”. In: *General Assembly and Scientific Symposium, 2011 XXXth URSI* (2011), pp. 1–4.
- [11] Greg Gibeling. “RDLC2: The RAMP Model, Compiler & Description Language”. MA thesis. University of California, Berkeley, 2008.
- [12] N.K. Govindaraju et al. “High performance discrete Fourier transforms on graphics processors”. In: *High Performance Computing, Networking, Storage and Analysis, 2008. SC 2008. International Conference for* (2008), pp. 1–12.

- [13] Gurobi Optimization Inc. *Gurobi Optimizer Reference Manual*.
- [14] M Jun, S Yoo, and EY Chung. “Mixed integer linear programming-based optimal topology synthesis of cascaded crossbar switches”. In: *Proceedings of the 2008 Asia and South Pacific Design Automation Conference* (2008), pp. 583–588.
- [15] S Kestur, JD Davis, and O Williams. “BLAS Comparison on FPGA, CPU and GPU”. In: *2010 IEEE Annual Symposium on VLSI* (2010), pp. 288–293.
- [16] Thorsten Koch et al. “MIPLIB 2010”. In: *Math. Prog. Comp.* 3.2 (June 2011), pp. 103–163.
- [17] Hirofumi Kondo et al. “A Multi-GPU Spectrometer System for Real-Time Wide Bandwidth Radio Signal Analysis”. In: *ISPA ’10: Proceedings of the International Symposium on Parallel and Distributed Processing with Applications* (Sept. 2010).
- [18] H Mittelmann. *Mixed Integer Linear Programming Benchmark*. URL: <http://plato.asu.edu/ftp/milpc.html>.
- [19] A Parsons et al. “PetaOp/Second FPGA Signal Processing for SETI and Radio Astronomy”. In: *Signals, Systems and Computers, 2006. ACSSC ’06. Fortieth Asilomar Conference on* (2006), pp. 2031–2035.
- [20] Aaron Parsons et al. “A Scalable Correlator Architecture Based on Modular FPGA Hardware, Reuseable Gateware, and Data Packetization”. In: *Publ Astron Soc Pac* 120.873 (Sept. 2008), pp. 1207–1221.
- [21] Rurik A Primiani, Jonathan Weintraub, and Jonathan deWerd. “SMA Wideband Correlator: Memo 1 Polyphase Filter-bank Utilization”. In: (June 2011), pp. 1–4.
- [22] A Rakhshanfar and JH Anderson. “An integer programming placement approach to FPGA clock power reduction”. In: *Proceedings of the 16th Asia and South Pacific Design Automation Conference* (2011), pp. 831–836.
- [23] Scott M Ransom et al. “GUPPI: Green Bank Ultimate Pulsar Processing Instrument”. In: *American Astronomical Society Meeting Abstracts #214* 214 (Dec. 2009).
- [24] A Siemion et al. “Current and Nascent SETI Instruments in the Radio and Optical: SERENDIP V.v, OSPOSH and HRSS”. In: *Astrobiology Science Conference 2010: Evolution and Life: Surviving Catastrophes and Extremes on Earth and Beyond* 1538 (Apr. 2010), p. 5378.
- [25] Hayden Kwok-Hay So. “BORPH: An Operating System for FPGA-based Reconfigurable Computers”. PhD thesis. ProQuest, 2007.
- [26] G Theodoridis, N Vassiliadis, and S Nikolaidis. “An integer linear programming model for mapping applications on hybrid systems”. In: *IET Computers & Digital Techniques* 3.1 (2009), p. 33.
- [27] A Thompson, J Moran, and G Swenson. “Interferometry and synthesis in radio astronomy”. In: *New York* (1986).

- [28] KH Tsoi and W Luk. “Axel: A heterogeneous cluster with FPGAs and GPUs”. In: *Proceedings of the 18th annual ACM/SIGDA international symposium on Field programmable gate arrays* (2010), pp. 115–124.
- [29] R Van Nieuwpoort and John W Romein. “Correlating Radio Astronomy Signals with Many-Core Hardware”. In: *International Journal of Parallel Programming* (2009).
- [30] Randall B Wayth, Lincoln J Greenhill, and Frank H Briggs. “A GPU-based Real-time Software Correlation System for the Murchison Widefield Array Prototype”. In: *Publ Astron Soc Pac* 121 (Aug. 2009), pp. 857–865.