

Internal Technical Report

CPU-GPU based DIGITAL Backend

S. Harshavardhan Reddy & Irappa M. Halagali
Ver. 2.0 , 11/06/2014.

Index

1. INTRODUCTION
2. BLOCK DIAGRAM
3. SPECIFICATIONS
 - a. uGMRT
 - b. GWB-II
 - c. iADC
 - d. ROACH
 - e. GPU Board
 - f. HOST Machines
 - g. Infiniband Switch
4. DESIGN DETAILS
 - a. Digitisation and Packetising in FPGA
 - b. Acquisition and Shared Memory buffering
 - c. Time slicing and data sharing using Infiniband
 - d. Correlation in GPU
 - e. Beamformer
 - f. Software Model
 - g. Online control, visibility and beam data recording
 - h. Computation and IO requirements
5. TESTS and RESULTS
 - a. LAB :
 - b. Antenna :
6. CONCLUSIONS
 - Appendix - 1 SOP
 - Appendix - 2 SOFTWARE PROGRAMS

Appendix - 3 POWER ON/OFF PROCEDURE

Appendix - 4 GWB NETWORK DIAGRAM

1. Introduction :

The GMRT consists of an array of 30 antennas, each of 45 m diameter, spread over a region of 25 km diameter, and operating at 5 different wave bands from 150 MHz to 1450 MHz. The maximum instantaneous operating bandwidth at any frequency band is 32 MHz. Each antenna provides signals in two orthogonal polarizations, which are processed through a heterodyne receiver chain and brought to the central receiver building, where they are converted to baseband signals and fed to the digital back-end consisting of correlator and pulsar receiver. The existing GMRT Software back-end (GSB) is built on software based approach designed from off-the-shelf components, PCI based ADC cards and a Linux cluster of 48 nodes with gigabit inter-node connectivity for real-time data transfer requirements.

GMRT is upgrading to uGMRT and the back-end systems are undergoing major changes to achieve the upgrade system specifications like increased bandwidth of 400MHz, direct processing of RF signals, increased dynamic range, improved channel resolution. The digital backend part of this upgrade was named GWB(GMRT WIDEBAND BACKEND). As part of this upgrade, a version of GWB(GWB-II) has been built and released for users. This is a 4-antenna dual polarization or 8-antenna single polarization correlator developed on CPU-GPUs and FPGA boards.

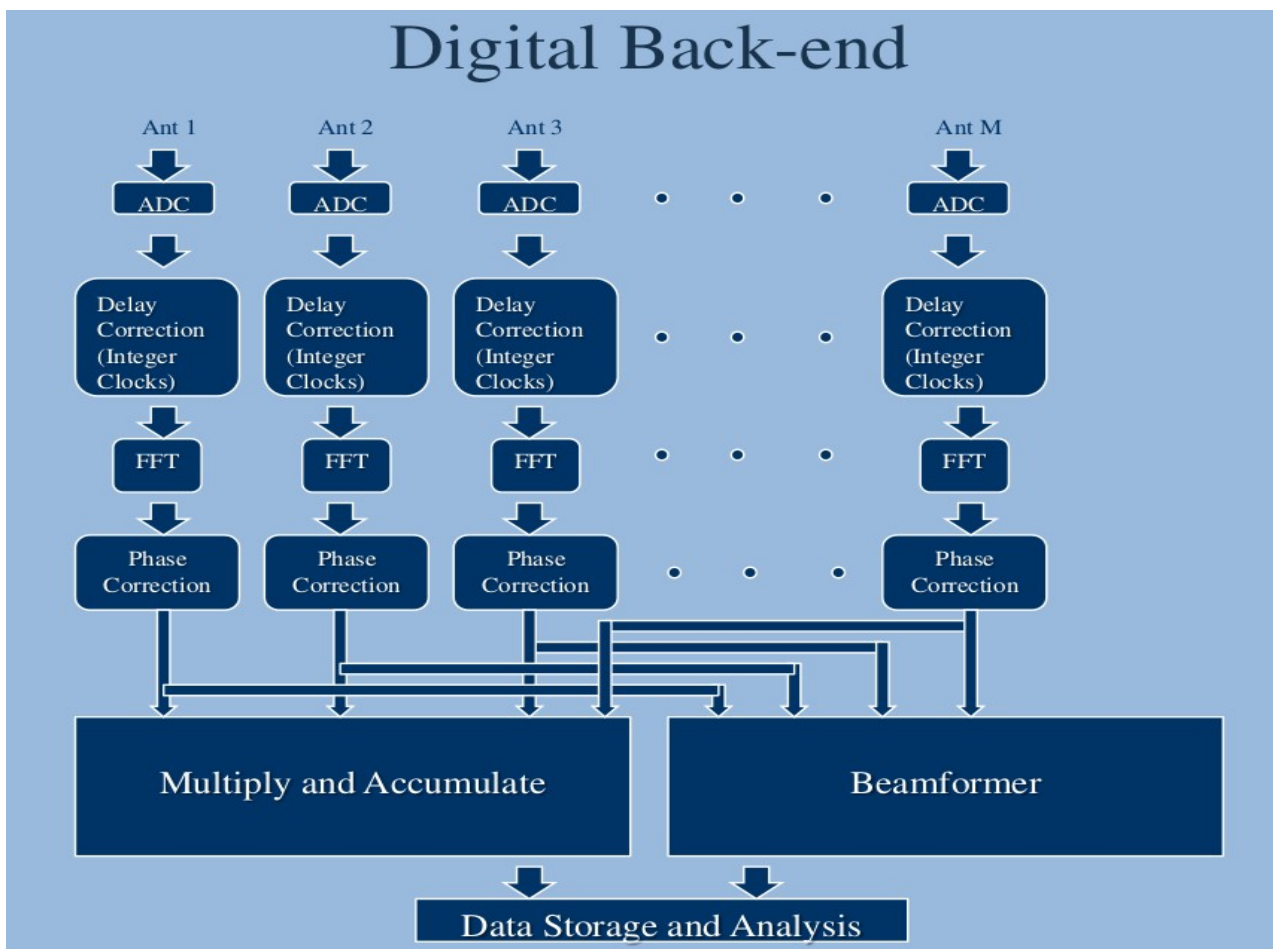


Fig : Digital Back-end

The design is a hybrid one using FPGAs and CPU-GPUs for various processes in the digital back-end chain. FPGAs connected with ADCs perform the digitisation and packetising the data while CPU-GPUs acquire the data, perform correlation and record the visibilities onto a disk for post-processing and analysis.

2. BLOCK DIAGRAM :

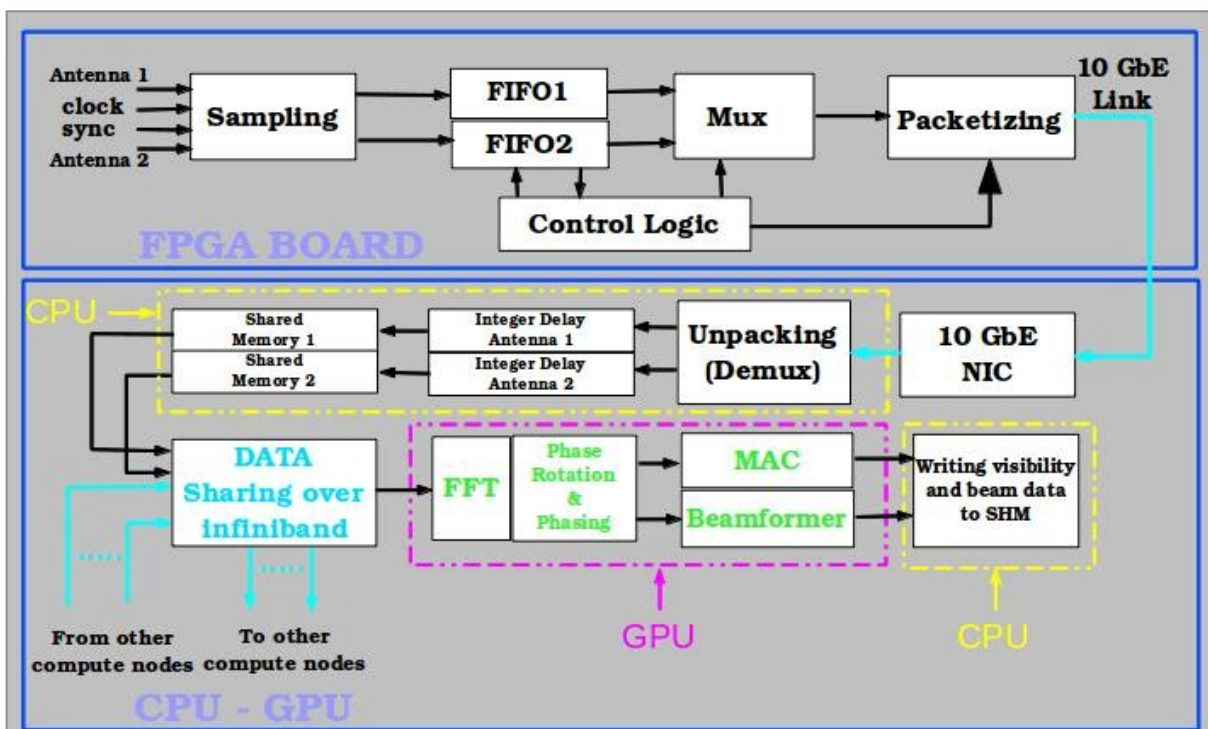
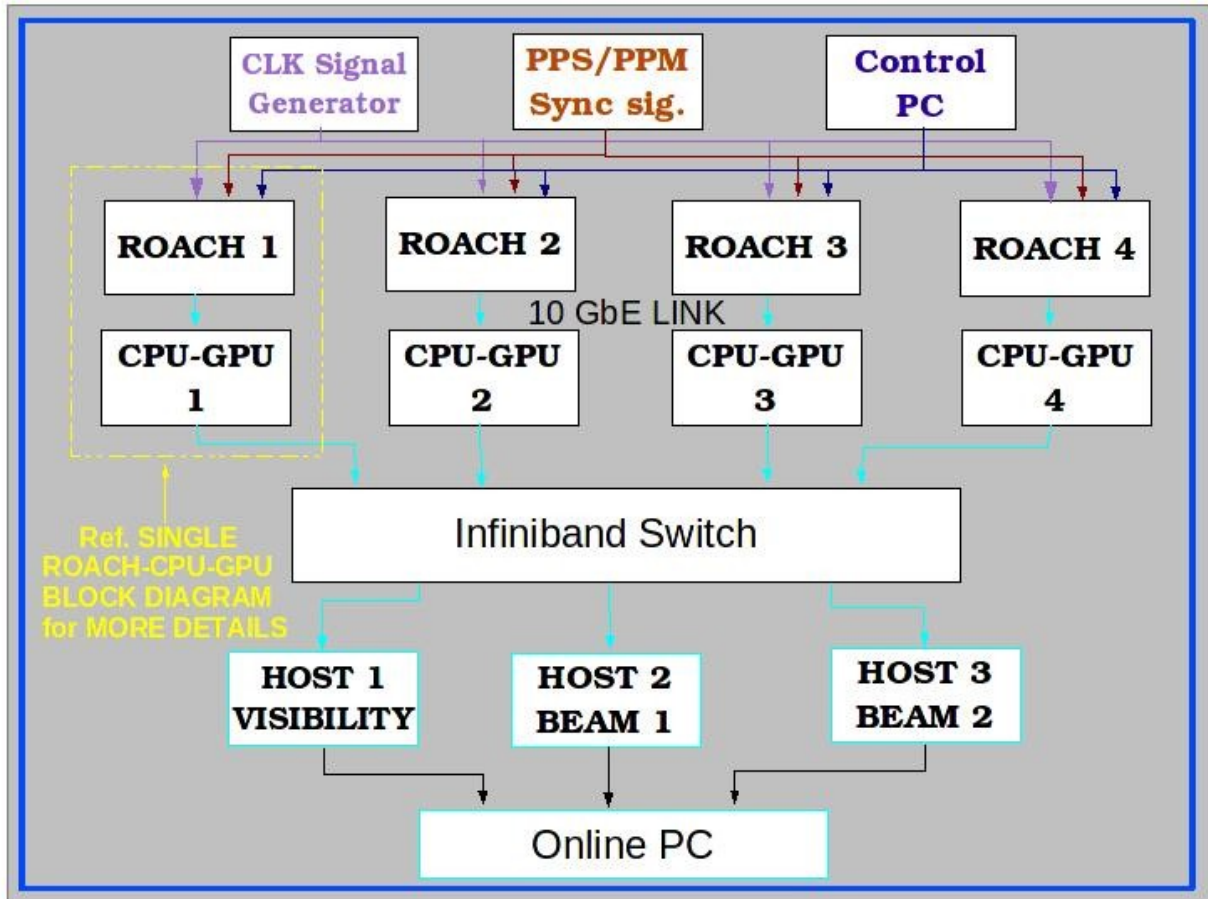


Fig : Data flow on single ROACH and CPU-GPU

Hardware :

1. iADC boards :- 1x Atmel/e2V AT84AD001B 8-bit dual 1Gsp/s, with clock 0dbm 50Ω 10MHz – 1GHz.
2. ROACH boards:- ROACH stands for Reconfigurable Open Architecture Computing Hardware. The board is a standalone FPGA processing board built around Xilinx Virtex-5.
3. CPU-GPUs with Nvidia GPU (Tesla C2050 or above) with Myricom 10GbE NIC card and Infiniband NIC.
4. Mellanox QDR 18-port Infiniband switch
5. Control PC for ROACH boards configuring and programming
6. PCs with Infiniband NIC as host, IA recording, PA recording nodes
7. Signal generator as clock source for the ADCs.

Software :

1. CASPER MSSGE Toolflow - Matlab-Simulink + System Generator +Xilinx 11.5 EDK
2. Centos Linux 5.6 (Kernel 2.6.18-194.32.1.el5 or higher on X86_64)
3. Python 2.6 on Control PC for ROACH boards
4. Nvidia GPU device drivers + CUDA toolkit on CPU-GPUs
5. OpenMPI 1.4.5 or higher.

Implementation :

The 4-antenna dual polarization correlator is implemented using two ROACH boards each with two ADCs and four CPU-GPUs and an 8-port infiniband switch(40 gbps).

3. SPECIFICATIONS :

a. uGMRT

Number of Stations	:	32
Maximum instantaneous Bandwidth	:	400MHz
Number of spectral channels	:	2048 – 8192
Number of input polarizations	:	2
Full Stokes capability	:	Yes
Dump time	:	min. 128ms
Coarse and fine delay tracking	:	+/- 128 us
Fringe rotation	:	up to 5 Hz
Subarray support	:	Yes
Incoherent and Phased array beams(for pulsar work)		

b. GWB : The specifications achieved in the current GWB system

Number of Stations	:	4/8
Maximum instantaneous Bandwidth	:	200MHz
Number of spectral channels	:	2048

Number of input polarizations	:	1 for 8 antennas, 2 for 4 antennas
Full Stokes capability	:	Yes
Dump time	:	min. 671ms
Coarse and fine delay tracking	:	+/- 128 us
Fringe rotation	:	up to 5 Hz
Subarray support	:	No
Incoherent and Phased array beams	:	Yes

c. iADC :

ADC with 1x [Atmel/e2V AT84AD001B](#) 8-bit Dual 1Gsp/s. It needs

- Inputs
 - Clock: 10MHz-1GHz 50Ohm 0dBm
 - Signal: 2 analog inputs.
 - Sync: LVTTTL (5V tolerant, ([SN65LVDS1](#)))
- Outputs
 - 1x [Tyco Z-DOK+](#) 40 differential pair connector
 - 8x (1x8 or 2x4) 8-bit offset binary data
 - 2x digital clocks to CASPER standard Z-DOK clock pins

d. ROACH (FPGA) Boards :

ROACH (Reconfigurable Open Architecture Computing Hardware) is Vertex 5 FPGA based reconfigurable hardware. The centrepiece of ROACH is a Xilinx Virtex 5 FPGA (either LX110T for logic-intensive applications, or SX95T for DSP-slice-intensive applications). A separate PowerPC runs Linux and is used to control the board (program the FPGA and allow interfacing between the FPGA "software registers/BRAMs/FIFOs" and external devices using Ethernet).

Two quad data rate (QDR) SRAMs provide high-speed, medium-capacity memory (specifically for doing corner-turns), and one DDR2 DIMM provides slower-speed, high-capacity buffer memory for the FPGA. The PowerPC has an independent DDR2 DIMM in order to boot Linux/BORPH.

The two Z-DOK connectors allow ADC, DAC and other interface cards to be attached to the FPGA, in the same manner as the [IBOB](#) allowed (with backwards compatibility for the ADC boards used with the IBOB).

Four CX4 connectors provide a total of 40Gbits/sec bandwidth for connecting ROACH boards together, or connecting them to other XAUI/10GbE-capable devices.

100Mbps ethernet port is used for both control and acquisition of the data.

e. GPU Boards : C2050

Architecture	:	Tesla 20-series GPU
Number of cores	:	448
Caches	:	64 KB L1 cache + Shared Memory / 32 cores, 768 KB L2 cache
PCI-e DMA Engines	:	2
Floating Point Peak Performance	:	1030 Gflops (single) 515 Gflops (double)
GPU Memory	:	3 GB
Memory Bandwidth	:	144 GB/s (GDDR5)
System I/O	:	PCIe x16 Gen2
Power	:	247 W (max)

f. Host Machines : Dell T7500

Processor	:	Intel Xeon X5550
cpu MHz	:	2.67GHz
cache size	:	8192 KB
No. of cores	:	4 per CPU
No. of CPUs	:	2
Memory	:	12 GB (scalable upto 192 GB)
No. of PCI-e slots	:	5 (one X4, two X8, two X16)
Operating System	:	Redhat Linux Centos 5.6
SMPS	:	1100 watt 85PLUS Power Factor Correcting(PFC) power supply

g. Infiniband switch : Mellanox 18-port switch

- 18-port QSFP non-blocking ports
- Aggregate data throughput up to 1.44 Tb/s
- Port-to-port latency < 100ns
- 9 virtual lanes: 8 data + 1 management
- Adaptive routing
- Congestion control
- Port mirroring
- 48K entry linear forwarding data base

-Supports Passive/Active copper or fiber cables

4. Design details :

a. Digitisation and Packetising in FPGA :

The design is made of Matlab Simulink blocks, CASPER BEE-XPS block set. The ADC is an 8-bit sampler and is connected to the ROACH board through a 40-pin Z-DOK connector. Each ADC can sample two input signals. The sampled data from each channel is stored in the FIFOs. Each FIFO is 8192 bytes in size. The data from FIFOs is then packetised in the 10GbE block using control logic and is sent over 10GbE network. The packet size is 8242 bytes. The packet structure is 42 bytes of UDP packet header, 8 bytes of packet counter, 4096 bytes from channel 1 and 4096 bytes from channel 2. The sync signal to the ADC is a PPS/PPM signal to synchronize the data from many ROACH boards. For 200MHz BW, 8 bits per sample and two input channels, the data is 6.4 gbps.

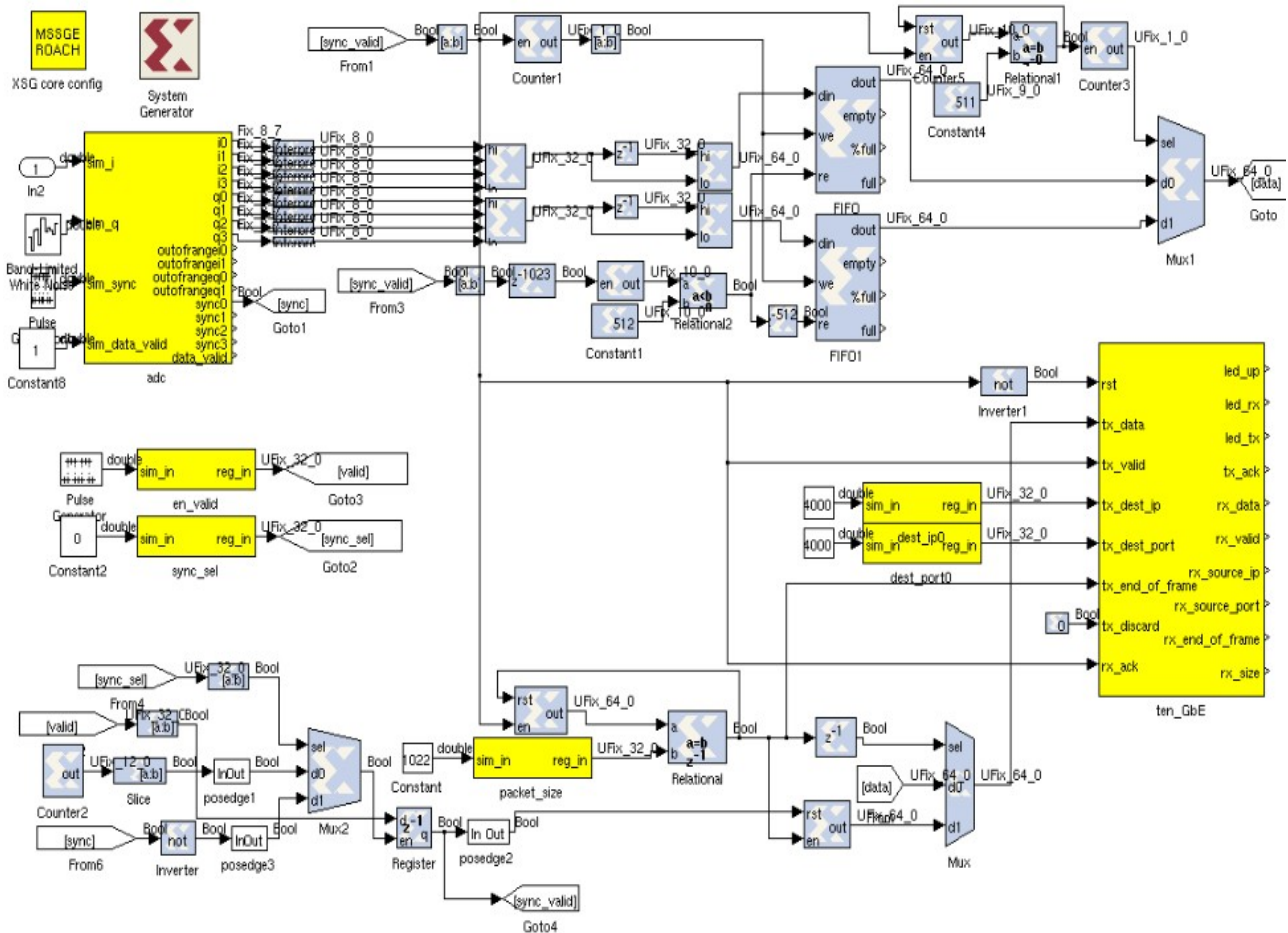


Fig : FPGA design

b. Acquisition and Shared Memory buffering :

The packets are received on the CPU side through the Myricom 10GbE NIC. The data thus received is written into shared memories. Two shared memories are created on each node for each input channel. Each shared memory is of four buffers, each buffer size is 256 MB. The data before writing into the shared memory is corrected for coarse delay. Also the time when the first packet is received in the first acquisition node is noted and shared

with other nodes for time synchronization. The nodes are NTP synchronized with the GPS time server. The coarse delay is corrected for every buffer of 256 MB.

c. Time slicing and data sharing using Infiniband :

The data from the shared memories is read and is shared over infiniband network with the other nodes. The shared memory buffer is read and sliced into number of nodes times slices. Node 0 keeps the first slice and sends the second slice to node1, third slice to node2 and so on. At the same time, node 0 receives slice 0 from node1, node 2 and so on. Similarly, node 1 keeps slice 1, sends the remaining slices to corresponding nodes and receives slice 1 from all the other nodes. This happens over all the nodes at the same time. To keep the nodes busy all the time, non-blocking MPI calls are used for this communication.

Algorithm :

- step 1 : initialise variable $i = 0$
- step 2 : Repeat steps 3 to 4 till $i = \log_2(\text{number of compute nodes}) - 1$
- step 3 : variables $wr = \exp_2(i)$, $cr = \exp_2(i+1)$, $\text{new_rank} = \text{rank} \% cr$, $\text{idx} = \text{rank}/cr$, $j = 0$
- step 4 : Repeat steps 5 to 15 till $j = wr - 1$
- step 5 : if $\text{new_rank} < wr$ perform steps 6 to 9 else perform steps 10 to 13
- step 6 : $\text{dest} = \text{new_rank} + wr + j$
- step 7 : if $\text{dest} > cr - 1$ perform step 8 else perform step 9
- step 8 : $\text{dest} = \text{dest} - wr + (cr * \text{idx})$
- step 9 : $\text{dest} = \text{dest} + (cr * \text{idx})$
- step 10 : $\text{dest} = \text{new_rank} - (wr + j)$
- step 11 : if $\text{dest} < 0$ perform step 12 else perform step 13
- step 12 : $\text{dest} = \text{dest} + wr + (cr * \text{idx})$
- step 13: $\text{dest} = \text{dest} + (cr * \text{idx})$
- step 14 : send the data to node ranked dest
- step 15 : receive the data from node ranked dest

Before MPI transfers				After MPI transfers			
Node0		Node1		Node0		Node1	
1	Ant1	1	Ant3	1	Ant1	2	Ant1
2		2		1	Ant2	2	Ant2
3		3		1	Ant3	2	Ant3
4		4		1	Ant4	2	Ant4
1	Ant2	1	Ant4	1	Ant5	2	Ant5
2		2		1	Ant6	2	Ant6
3		3		1	Ant7	2	Ant7
4		4		1	Ant8	2	Ant8
Node2		Node3		Node2		Node3	
1	Ant5	1	Ant7	3	Ant1	4	Ant1
2		2		3	Ant2	4	Ant2
3		3		3	Ant3	4	Ant3
4		4		3	Ant4	4	Ant4
1	Ant6	1	Ant8	3	Ant5	4	Ant5
2		2		3	Ant6	4	Ant6
3		3		3	Ant7	4	Ant7
4		4		3	Ant8	4	Ant8

Fig : Data sharing

d. Correlation in GPU :

A slice of data of all antennas in each node is copied to GPU for correlation process. FFT is performed using CUFFT library. The fourier transformed data is fine delay corrected and compensated for fringe rotation. The delay cal function calculates the fstc and fringe values from the starting time which was noted when the first packet was received.

FFT : Real to Complex FFT is performed efficiently by performing Complex to Complex CUFFT and reconstructing the output taking advantage of the symmetry.

For N-point FFT form a complex N/2 sample sequence(even samples are real and odd samples are imaginary).

$$z(n') = fe(n') + jfo(n') \quad n' = 0..N/2 - 1$$

Take the FFT of z using CUFFT to get

$$Z(k) = FFT_{N/2}(k, z) \quad k = 0..N/2 - 1$$

Launch a GPU kernel to perform

$$F(k) = \frac{1}{2} (\{Z(k) + Z(N/2 - k)^*\} - je^{-j2\pi k/N} \{Z(k) - Z(N/2 - k)^*\})$$

for k = 0, 1, ----, N/2 - 1

$$F(N/2) = \frac{1}{2} (\{Z(0) + Z(0)^*\} + j\{Z(0) - Z(0)^*\}) \quad \text{for } k = N/2$$

Parallelism is achieved over number of channels at the thread level and over the no. of

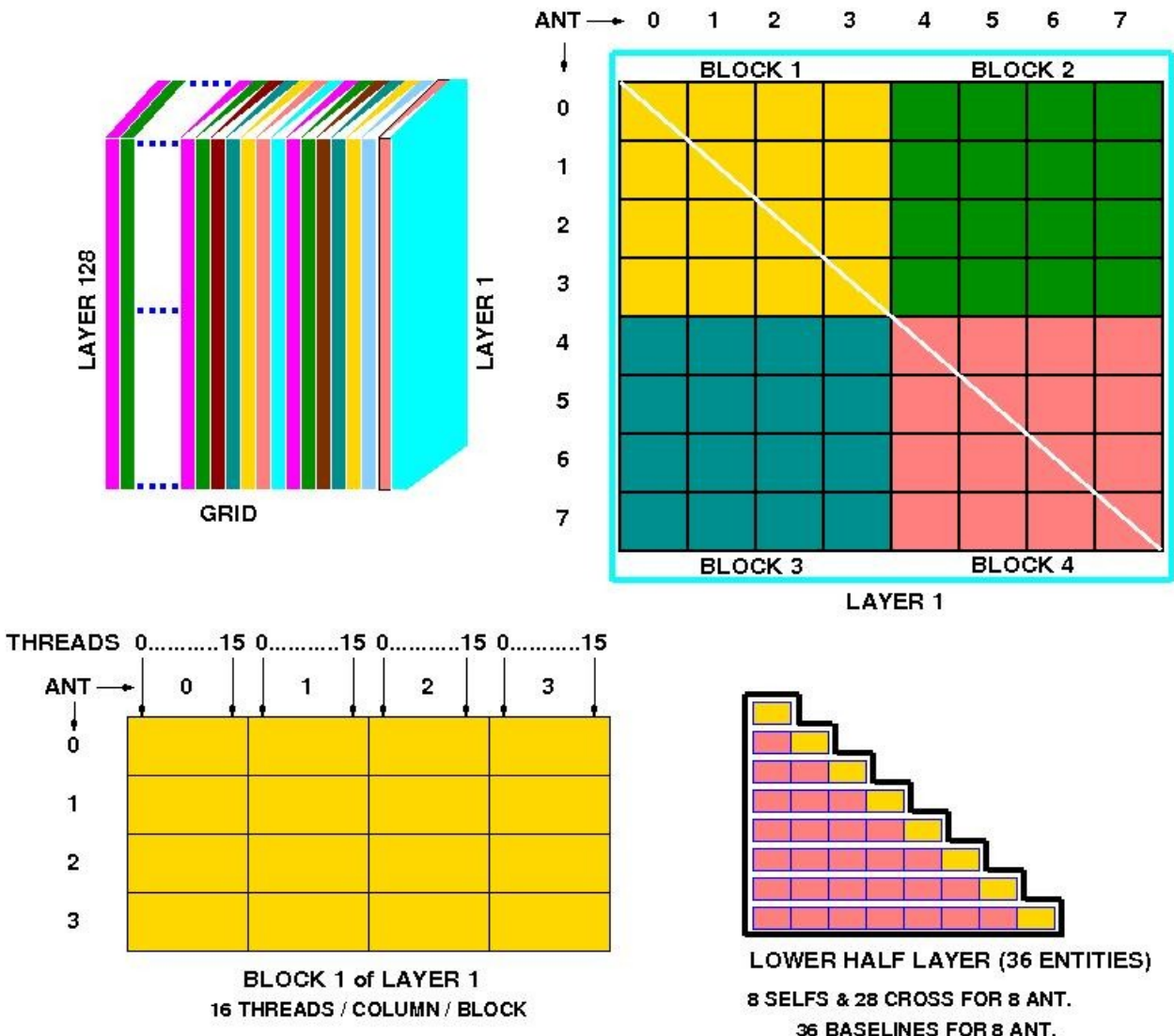
FFTs at the block level.

Phase rotation : This is done to compensate the fractional delay and fringe rotation. The delay cal function provides the values of delay in seconds, rate of change of delay, initial phase and fringe. From the values, total amount of phase that is to be applied is calculated and applied to every channel for all the antennas. Parallelism is achieved over the channels at the thread level and over the no. of FFTs per integration at the block level.

MAC : Multiply and Accumulation part is the most compute intensive one of all the other processes (see section 4.h Computations & IO requirements). For an array of N antennas, there will be $N*(N+1)/2$ baselines and correspondingly $N*(N+1)/2$ visibilities (including selfs of all N antennas). In the current system, the number of antennas should be in the integer powers of 2. While computing, the parallelism is achieved at both the thread level and block level in no. of frequency channels in X-dimension and no. of baselines in Y-dimension.

Let us say there are N antennas and K frequency channels. The block dimensions are selected as (bd.x,bd.y). Then grid dimensions (gd.x,gd.y) must be $(K/bd.x, (N/bd.y)^2)$. The grid dimensions are selected in such a way that the product of (square of bd.y) and gd.y must be equal to (square of N). Each block does multiplication and accumulation of $(bd.y)^2$ baselines for bd.x channels. Here is an example for 8 antennas and 2048 channels. For 8 antennas, 36 baselines are formed (including selfs).

Multiplication and Accumulation (MAC)



are needed to be calculated from the baseline matrix(see fig). When we select block dimensions as (16,4) then grid dimensions become (128,4), which means each block calculates products of one 4x4 matrix for 16 channels. If any block is completely out of the lower triangle it simply quits without performing any operations. To reduce accessing the global memory for every baseline, shared memory is used.

e. Beamformer :

The beamformer specification is to provide two modes of operations : (i) an incoherent array mode, where the voltage samples from the selected antennae are added after converting to intensities, (ii) a coherent or phased array mode, where the voltage signals from the selected antennae are first added, and then converted to intensity samples. Both these modes require antenna based time delays to be corrected prior to the addition of the signals. Furthermore, for the coherent beamformer mode, it is also required to calibrate out antenna based phase offsets.

Currently, the sampling period is set to 1.3ms which is equal to averaging of 128 FFTs at 2048 spectral channels. The design provides two beams at the same time in two different beam host nodes. In each beam, either Incoherent array beam(IA) or Coherent beam beam(PA) can be selected. Further, if 4 stokes of PA are needed, then Interferometry needs to be run in 4 stokes mode. In the design, the parallelism is achieved over no. of spectral channels at thread level and no. of sampling period blocks in each integration of visibility data(which is 671ms) at the block level. Each thread hence calculates average of one spectral channel for all the 128 FFTs in a sampling period block. The design was optimised to read the global memory only one for the processing of both beams. Also, array masking has been implemented at the lowest level to improve the performance and remove computations for antennas masked out of the array. Phasing for Coherent array(PA) mode has been implemented in the phase rotation stage.

f. Software Model :

The various processes mentioned above in the CPU-GPU are performed parallelly using nested openMP sections. The main thread creates two threads. One for data acquisition and shared memory buffering. The other thread creates four more threads for (1)data reading from shared memories, (2)MPI sharing of data, (3)correlation in GPU and (4)writing the visibilities to shared memory to be written by a host node. Double buffering ping-pong scheme has been implemented to enable the processing and data transfer in parallel.

g. Online control and visibility recording :

PROGRAMING Used :

Control and Monitoring programs for ONLINE system and GWB Systems are implemented using various programing and scripting languages, Such as C, C++, QT3, Bash, Perl, Perl-Tk, Tcl-Tk etc.

ONLINE System :

GMRT Wide-band Back-end (GWB) can be controlled from Online Machines (both '*lenyadri*' as well as '*shivneri*'). For this, a socket is created, through which online sends the commands to GWB and receives their acknowledgments. This Socket program has been written in C, C++. In this way, data control and monitoring information is shared between the two systems (ONLINE and GWB).

GWB System :

GWB data acquisition, processing and recording programs are written in C and C++. Where as Graphical user interface such as GWB DasConsole, GWB configuration, GWB Pulsar configuration, etc. are designed using C, C++ and QT3. Power Equalization, GAC-Config are written using C, Perl-Tk. Also, various other scripting languages like Bash, Perl, Perl-Tk, Tcl-Tk etc. are used for most online as well as offline applications.

USER INTERFACE :

After starting the data acquisition through GWB, we can control and Monitor it using both GUI and Command line tools as per the user's/observer's requirement. For this, DasConsole, GWB-Config, DasMon, Pulsar-DasConsole, Pulsar-Config, PowerEq, Phasing, GAC-Config etc. are Interface tools available.

OUTPUT DATA :

Data acquired from GWB is stored in Standard GMRT LTA format. And Pulsar Data is stored in RAW formats. This data can be monitored using online/offline tools viz. DasMon. Itahdr, tax, xtract, rantsol, pmon, etc.

h. Computation and IO requirements :

For uGMRT, FFT(16k point) computation requires 2.9 teraflops, phase rotation takes 0.1 teraflops and MAC takes 6.6 teraflops. The IO requirement at 400 MHz bandwidth and 4-bits per sample is nearly 25 GB/s.

The same computation requirements for GWB-II are : 153 gigaflops for FFT(4k point), 6.4 gigaflops for phase rotation and 230 gigaflops for MAC. IO requirement is 3.2 GB/s. The performance achieved on C2050s is 345 gigaflops for FFT and 490 gigaflops for MAC. For host to device transfer the bandwidth achieved is nearly 4.3 GB/s. On infiniband network the bandwidth achieved is 5 GB/s.

5. TESTS and RESULTS :

Interferometer :

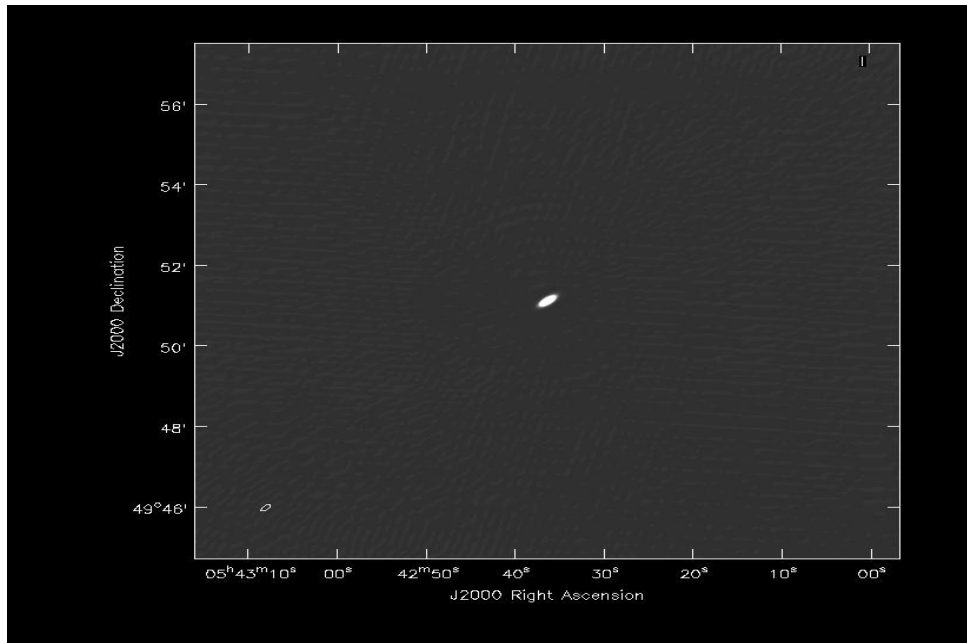


Image of 3C147 : First image of GWB correlator taken with 8 antennas at 1280MHz RF and 32 MHz signals. Data recorded on 07 June 2011

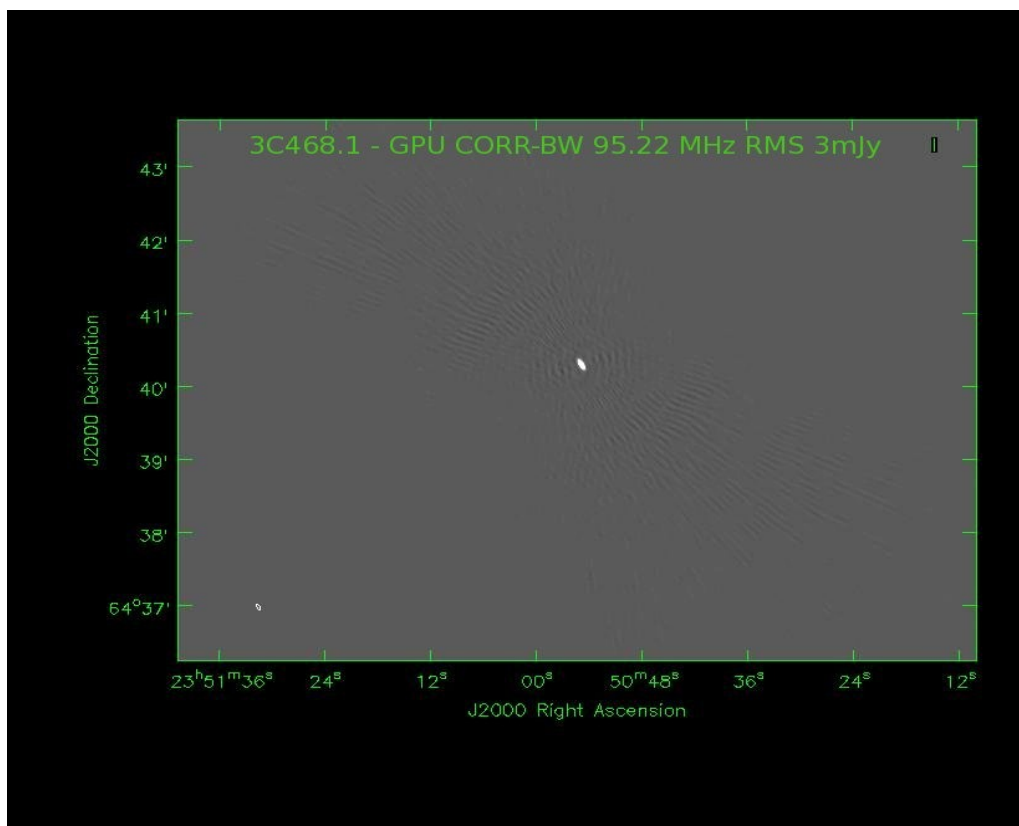


Image of 3C468.1 : First image using broadband signals

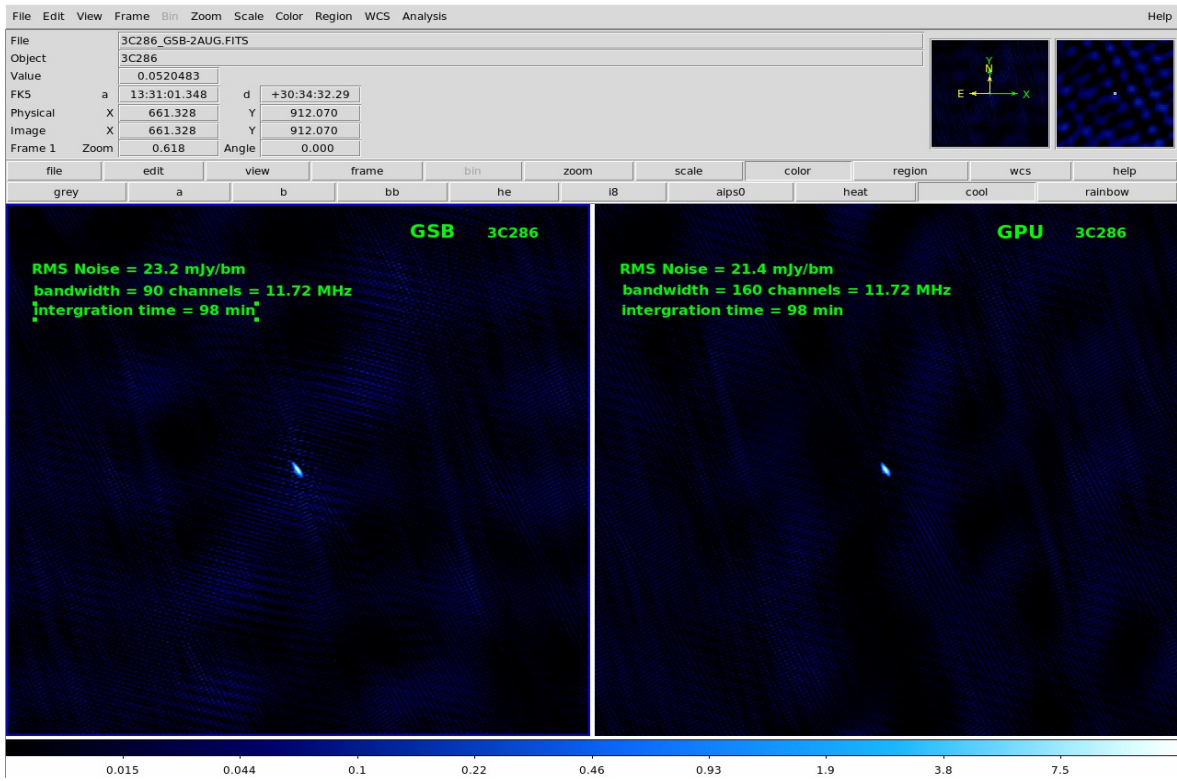
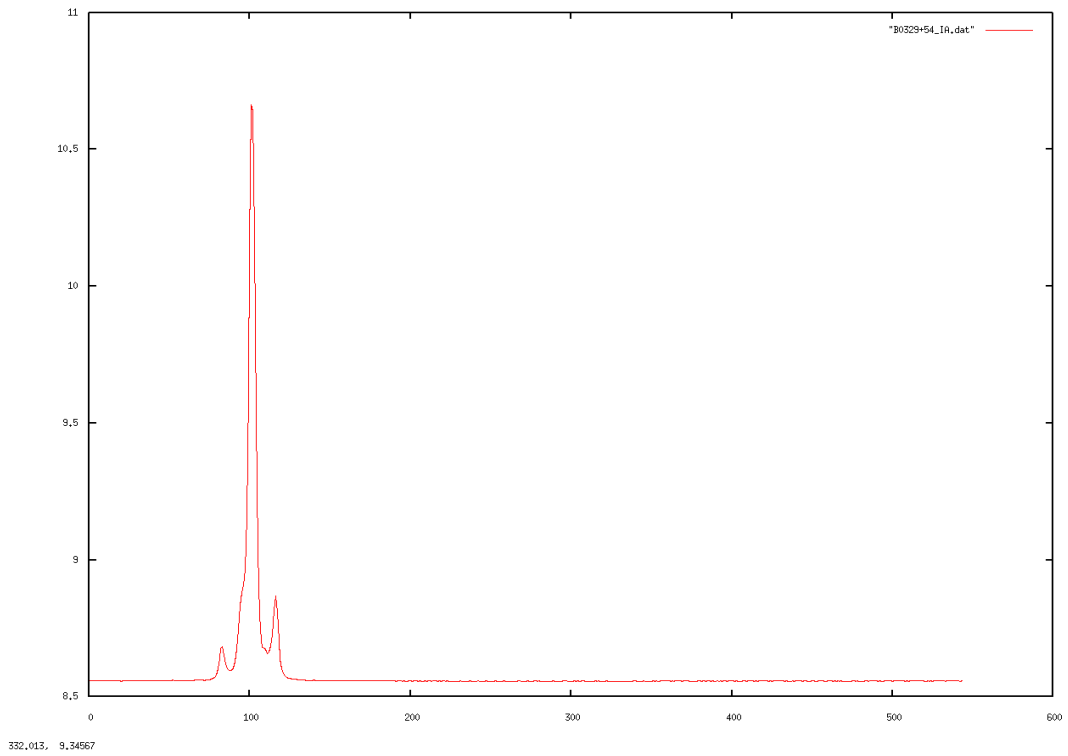
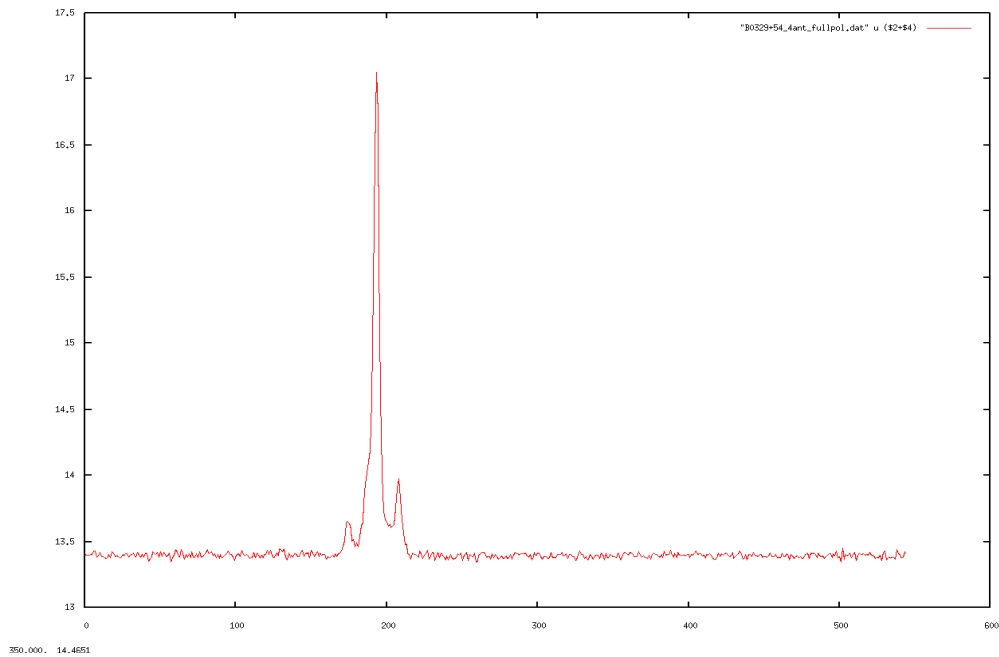


Image of 3C286 : A comparison between GSB and GWB

Beamformer :



Profile of B0329+54 : Incoherent array beamformer at 610MHz RF and sampling period of 1.3ms



Profile of B0329+54 : Phased array beamformer at 610Mhz RF and sampling period of 1.3ms

9. Conclusions : An 8-antenn single pol/4-antenna dual pol backed has been released for testing by users. Scaling this prototype to 32-antennas completes the digital backend for uGMRT.

Appendix - 1 SOP

Please visit the page : http://corrbk/corr_web/IRU/SOPs/GWB/ and download the file SOP_GWB2_Ver3.1.pdf

Appendix - 2 SOFTWARE PROGRAMS

Code is located at [/home/gpuuser/GWB/ver3/code/psrdada/gmrt_online_dual_pol_coherent_beam](http://home/gpuuser/GWB/ver3/code/psrdada/gmrt_online_dual_pol_coherent_beam) on gpunode6 machine(192.168.4.75).

Username : gpuuser

Password : gmrt123

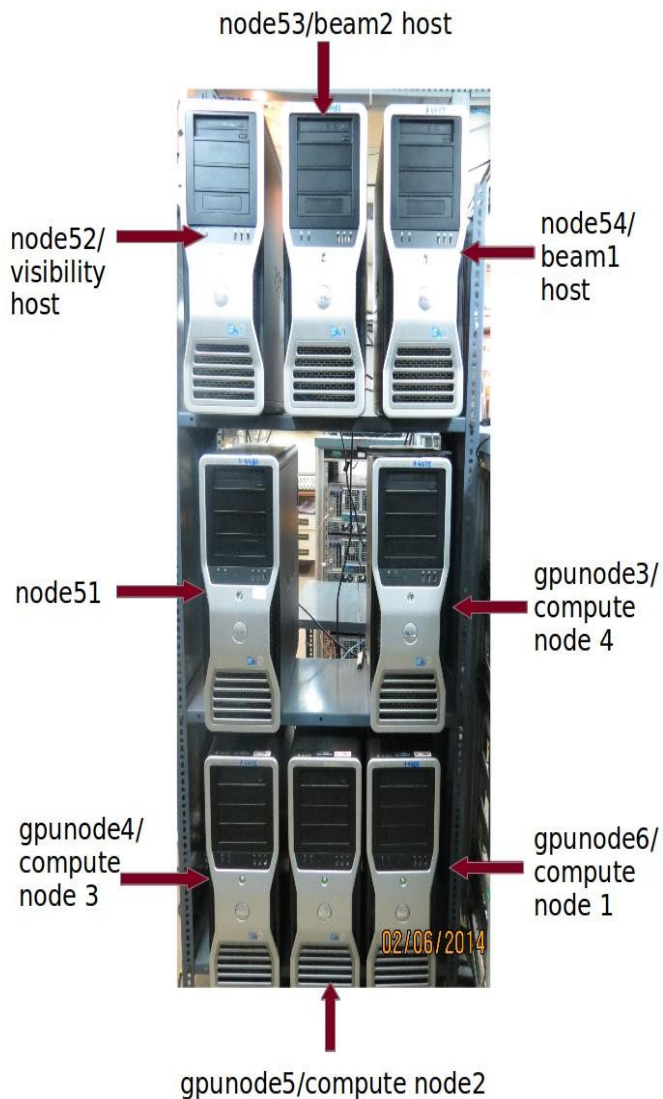
Appendix - 3 POWER ON/OFF PROCEDURE

GWB racks



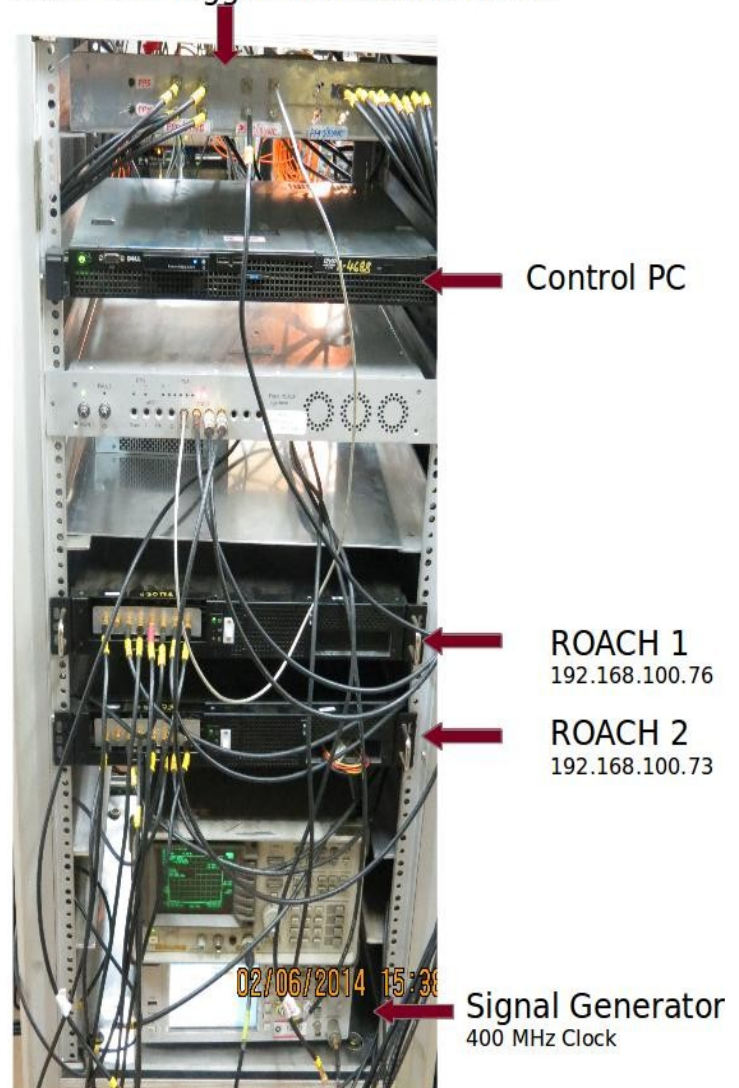
Image courtesy : Irappa M. Halagali

Compute and Host nodes rack



Sampler and Packetizer rack :

Clock and Trigger distribution board



Images courtesy : Irappa M. Halagali

1. Switch OFF procedure.

a. Switch off the ROACH UNITS 192.168.100.76 & 192.168.100.73 in the GWB rack by holding down the Black switch on the front panel for ~5 sec.

b. Switch off the Signal Generator, kept at bottom of the GWB rack. This feeds clock signal of 400 MHz, +17dBm to the ROACH boards.

c. No need to switch off the SMPS used for PPS unit. This will get switched off directly from mains.

d. No need to switch off the infiniband switch. This will get switched off directly from mains.

e. Halt the control PC (192.168.4.68) which is a 1U pc in the GWB rack.

NOTE : a. ssh -X root@192.168.4.68 (gmrttifr) b. halt -p

f. Halt the compute nodes 192.168.4.75(gpunode6), 192.168.4.76(gpunode5), 192.168.4.77(gpunode4) & 192.168.4.78(gpunode3).

NOTE1 : a. ssh -X root@192.168.4.xx (123\$abc) b. halt -p

NOTE2 : Switch off the 4 host nodes, node 192.168.16.51(node51), 192.168.16.52(node52), 192.168.16.53(node53) & 192.168.16.54(node54) with GSB system only.

NOTE1 : a. ssh -X root@192.168.16.xx (123\$abc) b. halt -p

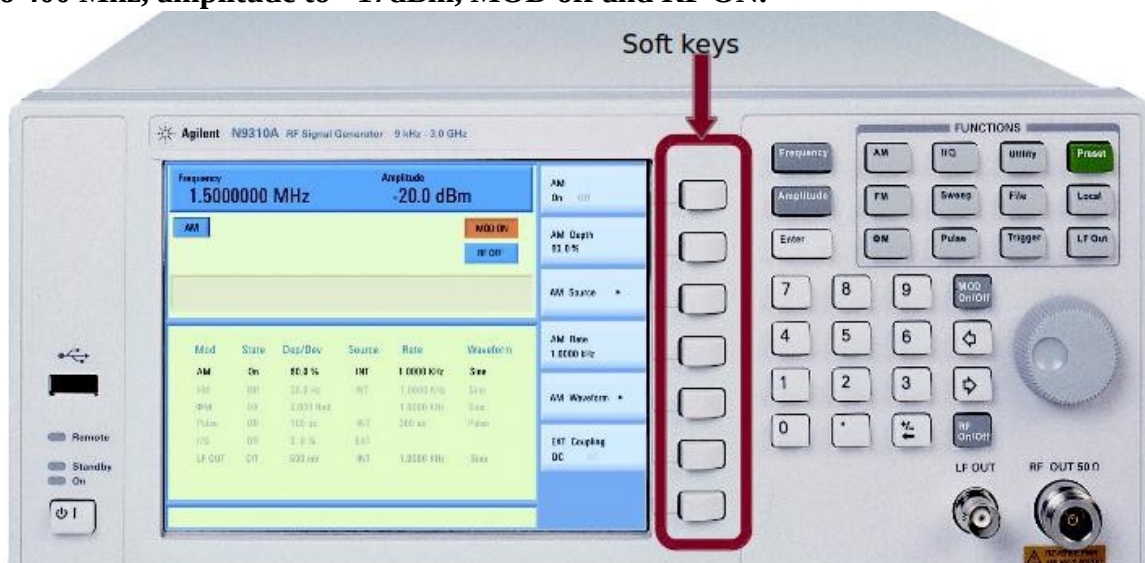
2. Switch ON procedure.

a. Switch ON the control PC (192.168.4.68).

b. Make sure the infiniband switch is ON.

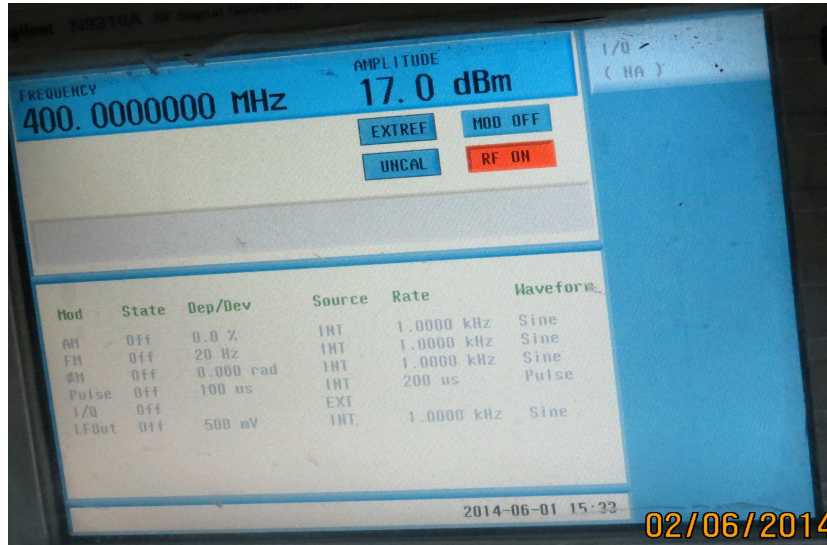
c. Make sure the PPS unit is Switched ON.

d. Switch ON the Signal Generator, kept at bottom of the GWB rack. Set the frequency to 400 Mhz, amplitude to +17dBm, MOD off and RF ON.



1. Press “Frequency” key
2. Type “400” from the numbers
3. Select “MHz” from the “soft keys”
4. Press “Amplitude” key
5. Type “17” from the numbers
6. Select “dbm” from the soft keys
7. Press “MOD On/Off” key to turn it off.
8. Press “RF On/Off” key to turn it on.
9. Press “Utility” key
10. Select “Reference” from the soft keys
11. Select “EXT REF” from the soft keys

Final screen on the Signal generator should look like as given in the picture below :



e. Switch ON the ROACH UNITS 192.168.100.76 & 192.168.100.73 in the GWB rack by holding down the Black switch on the front panel for ~2 sec.

f. Switch ON the compute nodes 192.168.4.75(gpunode6), 192.168.4.76(gpunode5), 192.168.4.77(gpunode4) & 192.168.4.78(gpunode3). And host nodes, node 192.168.16.51(node51), 192.168.16.52(node52), 192.168.16.53(node53) & 192.168.16.54(node54) with GSB system.

Appendix - 4 GWB NETWORK DIAGRAM

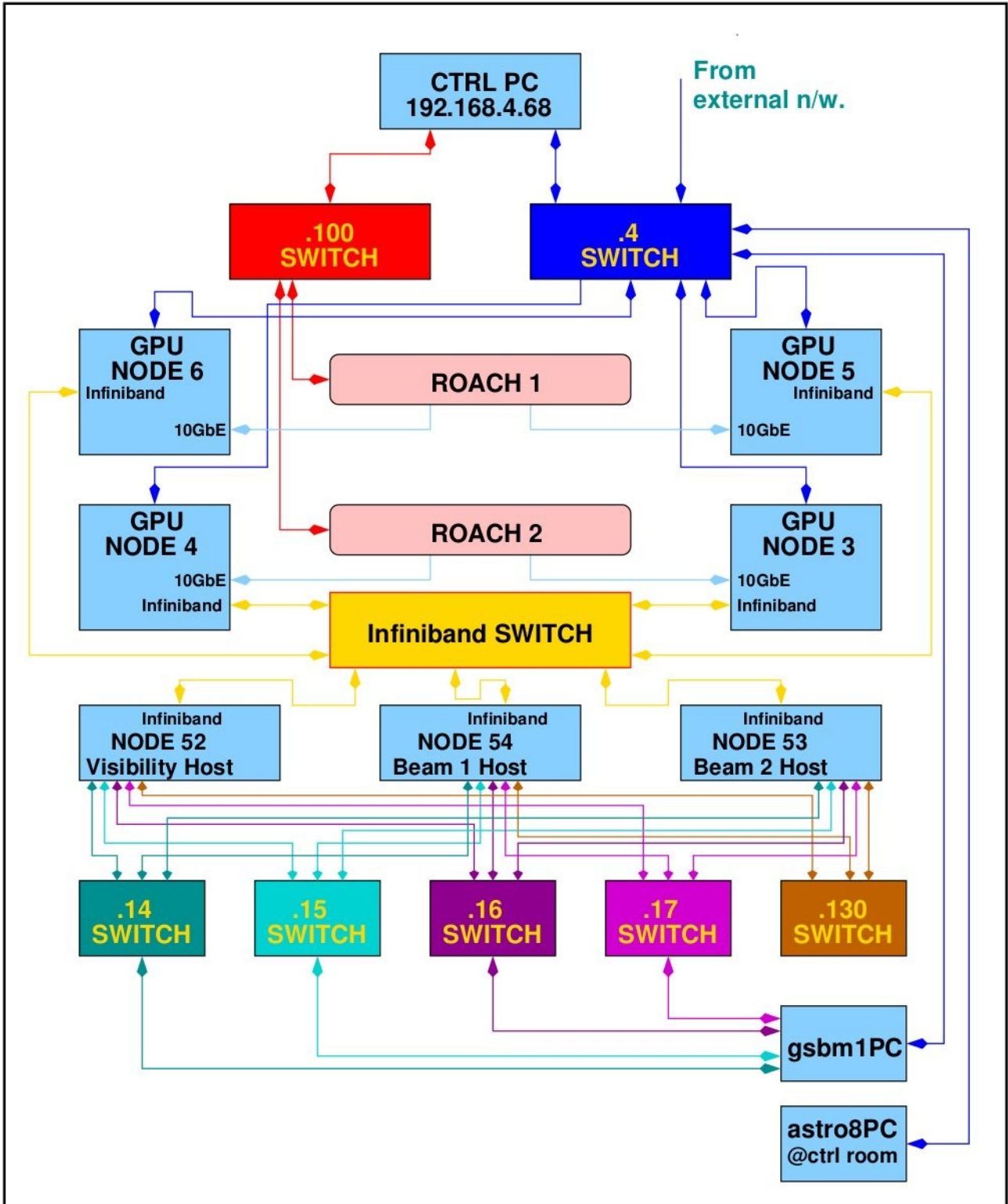


Image courtesy : Irappa M. Halagali