



GPUS IN GRAVITATIONAL-WAVE DATA ANALYSIS

DREW KEPPEL^{1,2}

FOR THE LIGO SCIENTIFIC COLLABORATION
& VIRGO COLLABORATION

3RD ASPERA COMPUTING AND ASTROPARTICLE PHYSICS WORKSHOP
HANNOVER, GERMANY
2012-05-03

LSC AND VIRGO GPU CONTRIBUTORS

- ABILENE CHRISTIAN UNIVERSITY

- JOSH WILLIS

- AEI

- CARSTEN AULBERT, OLIVER BOCK, TITO DAL CANTON, HEINZ-BERND EGGENSTEIN, DREW KEPPEL, BADRI KRISHNAN, BERND MACHENSCHALK, KARSTEN WIESNER

- CANADIAN INSTITUTE OF THEORETICAL ASTROPHYSICS

- KIPP CANNON

- EÖTVÖS LORÁND UNIVERSITY

- MÁTÉ NAGY

- SYRACUSE UNIVERSITY

- DUNCAN BROWN, ALEXANDER NITZ

- TSHINGHUA UNIVERSITY

- ZHIHUI DU, YUAN LIU

- UNIVERSITÀ DEGLI STUDI “CARLO BO” DI URBINO

- RICCARDO STURANI

- THE UNIVERSITY OF BIRMINGHAM

- BEN AYLOTT

- UNIVERSITY OF WESTERN AUSTRALIA

- DAVID BLAIR, SHIN KEE CHUNG, AMITAVA DATTA, SHAUN HOOPER, LINQING WEN

- WIGNER RESEARCH CENTRE FOR PHYSICS

- DEBRECZENI GERGELY

GW INTERFEROMETERS

LIGO HANFORD
OBSERVATORY



GEO600



LIGO LIVINGSTON
OBSERVATORY

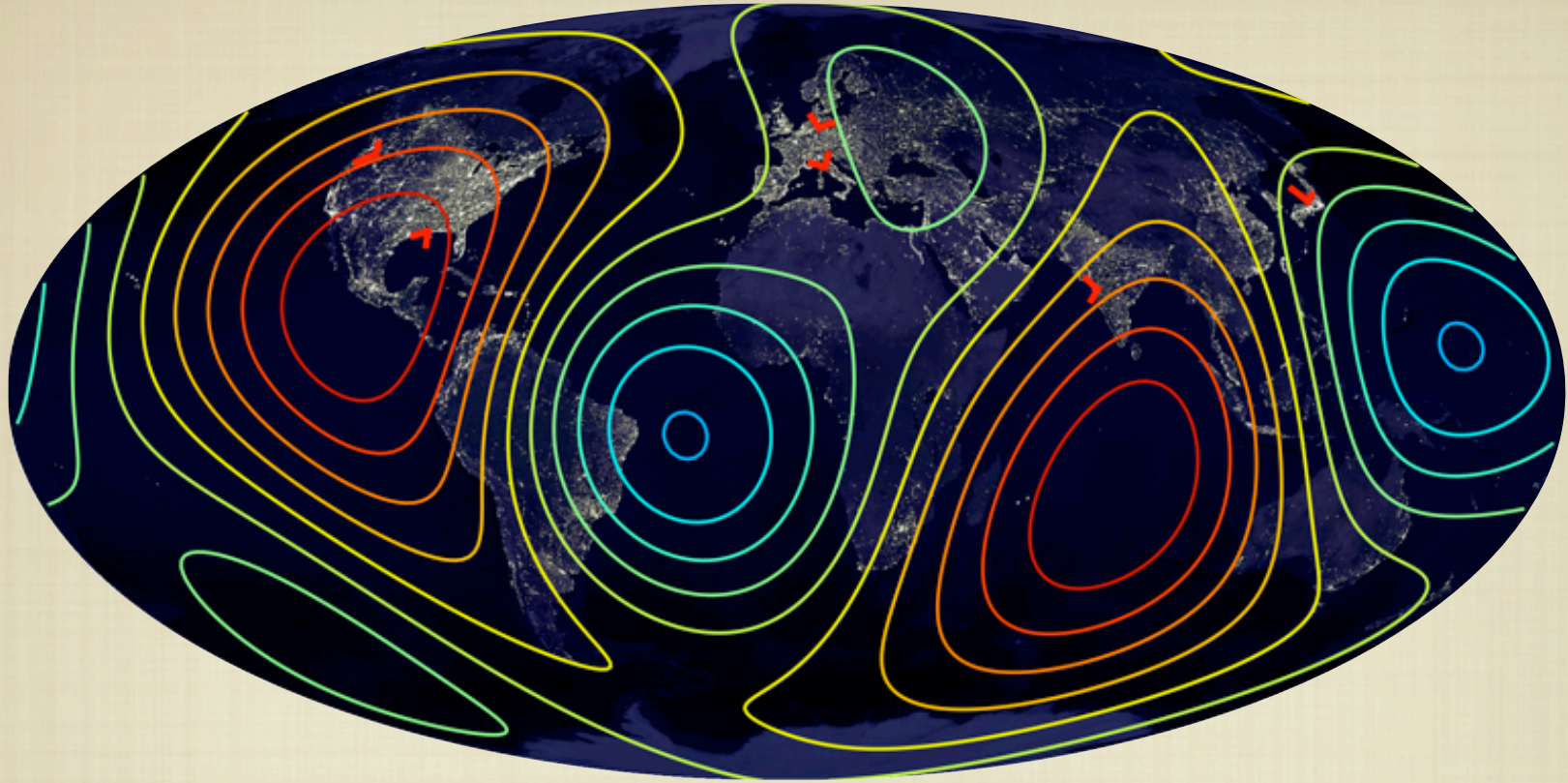


VIRGO



- PUMPED UP MICHELSON INTERFEROMETERS AT THEIR CORE
- SENSITIVE TO DIFFERENTIAL DISPLACEMENTS

GW DETECTOR NETWORK



■ ERA OF THE FIRST GENERATION OF INTERFEROMETRIC DETECTORS HAS ENDED

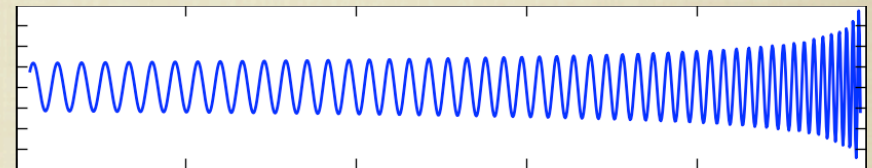
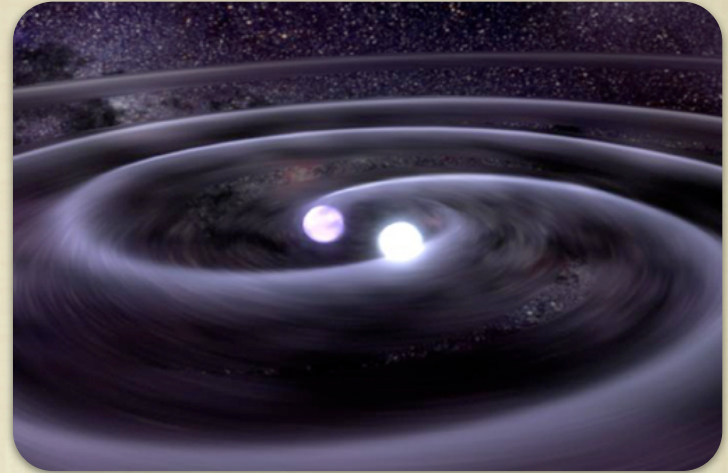
■ UPGRADES UNDERWAY TO SECOND GENERATION DESIGNS

■ EXCEPTION: GEO600 OPERATING IN “ASTROWATCH” MODE

Photo Credit: Satellite data courtesy Marc Imhoff of NASA GSFC and Christopher Elvidge of NOAA NGDC. Earth image by Craig Mayhew and Robert Simmon, NASA GSFC.

GW SIGNALS FROM INSPIRALING COMPACT BINARIES

- NEUTRON STAR AND/OR BLACK HOLE BINARIES
- MOST PROMISING SOURCE FOR SECOND GENERATION DETECTORS
- 10S OF DETECTIONS EXPECTED EACH YEAR
- CHARACTERIZED BY A “CHIRPING” SIGNAL

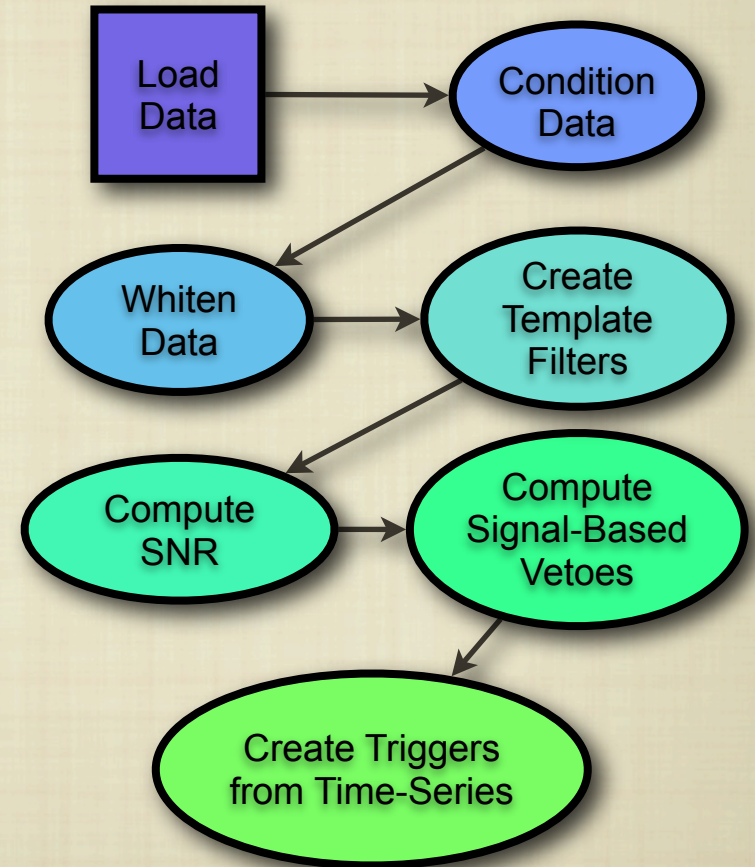


$$\Delta t = (5\mathcal{M})^{-5/3} \left(\frac{5}{8\pi f} \right)^{8/3}$$

FREQUENCY BAND (HZ)	DURATION (S)
40 - 2048	~25
10 - 2048	~975

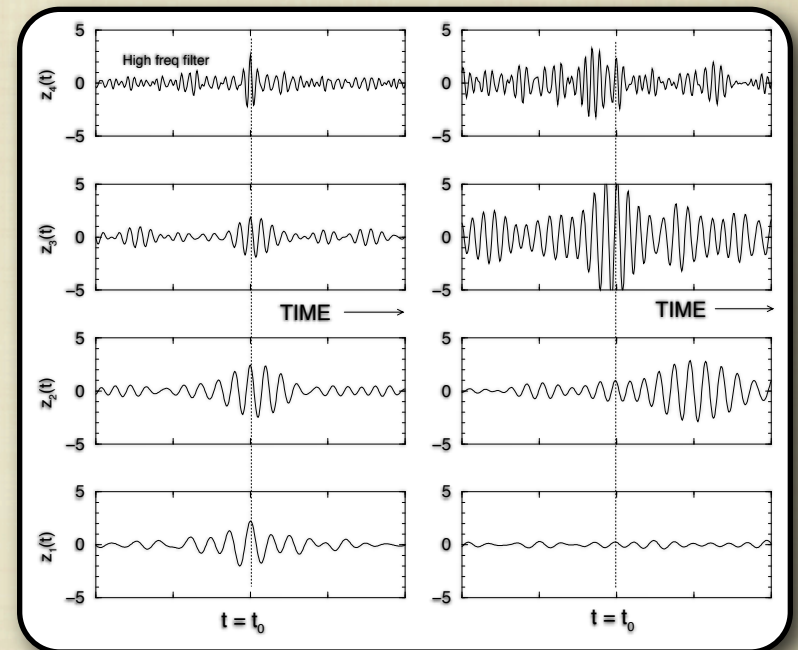
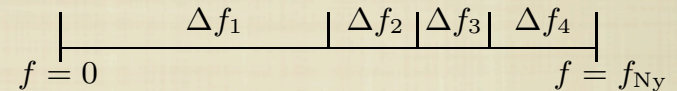
SEARCHING FOR INSPIRAL SIGNALS: SINGLE DETECTOR METHOD

- **MATCHED FILTER USED TO SEARCH FOR KNOWN WAVEFORMS**
- **OVERLAP-SAVE ALGORITHM USED TO REDUCE COMPUTATIONAL COST AT EXPENSE OF LATENCY**
- **SNR = 1 IFFT**
- **NON-GAUSSIAN DATA REQUIRES SOMETHING MORE**
- **SIGNAL-BASED VETOES REDUCE EFFECTS OF NON-GAUSSIAN DETECTOR GLITCHES**
- **CHISQ = 16 IFFTs**



SEARCHING FOR INSPIRAL SIGNALS: SINGLE DETECTOR METHOD

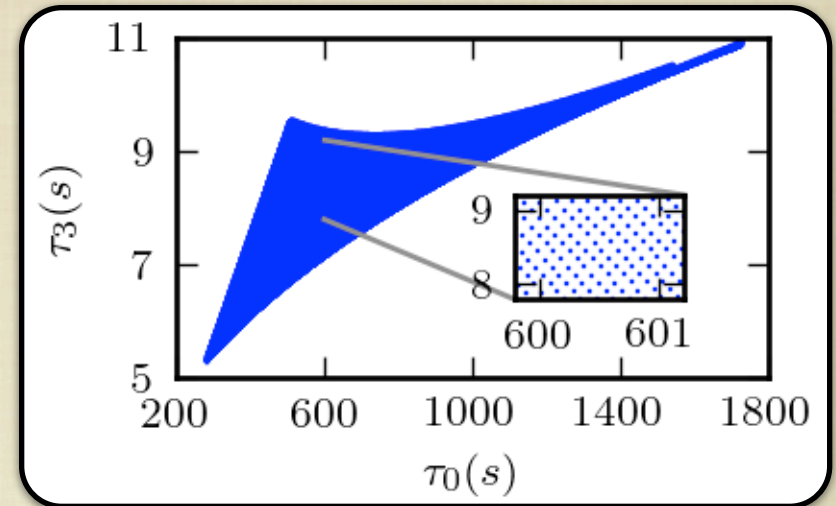
- **MATCHED FILTER USED TO SEARCH FOR KNOWN WAVEFORMS**
- **OVERLAP-SAVE ALGORITHM USED TO REDUCE COMPUTATIONAL COST AT EXPENSE OF LATENCY**
- **SNR = 1 IFFT**
- **NON-GAUSSIAN DATA REQUIRES SOMETHING MORE**
- **SIGNAL-BASED VETOES REDUCE EFFECTS OF NON-GAUSSIAN DETECTOR GLITCHES**
- **CHISQ = 16 IFFTS**



$$\chi^2 = \sum_{i=1}^{16} |\rho_i - \rho/16|^2$$

SEARCHING FOR INSPIRAL SIGNALS: SINGLE DETECTOR COST

- FULL BNS SEARCH IN SECOND GENERATION DETECTORS
- TEMPLATE BANKS MADE UP OF FEW $\times 10^5$ WAVEFORMS
- WAVEFORMS HAVE FEW $\times 10^6$ SAMPLES
- ~ 100 GFLOPS NEEDED TO PRODUCE SNR DATA IN REAL TIME (WITH A LATENCY OF ~ 15 MINUTES)
- $\sim 10\times$ MORE INCLUDING CHISQ VETO CALCULATION
- SMART ALGORITHM DESIGNS CAN REDUCE COMPUTATIONAL COST AND MEMORY FOOTPRINT



SPEED OF 2^{21} -LENGTH FFTS	# OF UNITS FOR REAL TIME PROCESSING
ATLAS CPU CORE @3GFLOPS	~ 10
NVIDIA C2050S @200GFLOPS	$\sim 1^*$

*NEGLECTING MEMORY RESTRICTIONS

GPU ACCELERATION EFFORTS: THE FFT & CHISQ CALCULATION

- PORTING JUST THE FFT ROUTINES OF “lalapps_inspiral”
- LIMITED BY THE TRANSFERRING DATA TO AND FROM GPU FOR EACH FFT
- PORTING ALSO THE CHISQ CALCULATION OF “lalapps_inspiral”
- REDUCES NUMBER OF TRANSFERS TO AND FROM GPU BY 8X

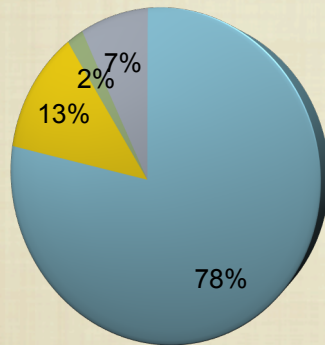
LALAPPS_INSPIRAL W/	EXECUTION TIME (SECONDS)
CPU FFT	480
CUDA FFT	164
CPU FFT AND CPU CHISQ	1530
CUDA FFT AND CPU CHISQ	895
CUDA FFT AND CUDA CHISQ	188

GPU ACCELERATION EFFORTS:

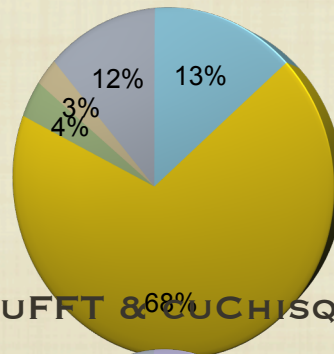
“lalapps_inspiral” Limitations

- LIMITED BY:
 - OTHER CPU OPERATIONS
 - MEMORY TRANSFERS BETWEEN TO AND FROM GPU
 - RESTRUCTURING CODE NECESSARY FOR OBTAINING FULL GPU BENEFITS

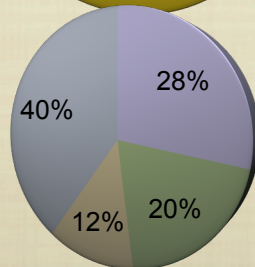
FFTW & CHISQ



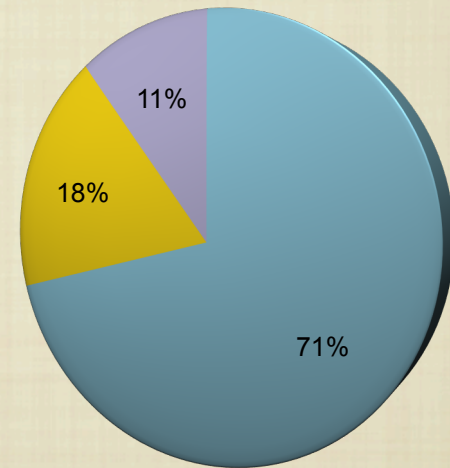
CUFFT & CHISQ



CUFFT & cuCHISQ



GPU OPERATIONS



- FFTW / CUFFT
- CHISQ
- CUFFT & cuCHISQ
- XLALBANKVETOCMAT
- KERNEL
- OTHERS

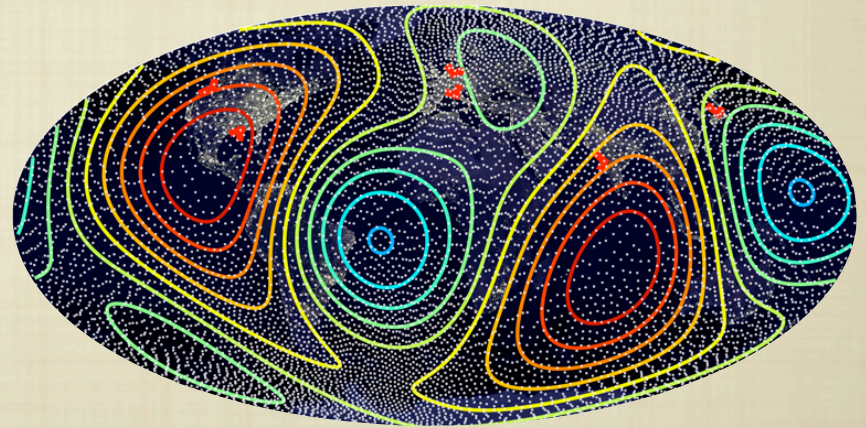
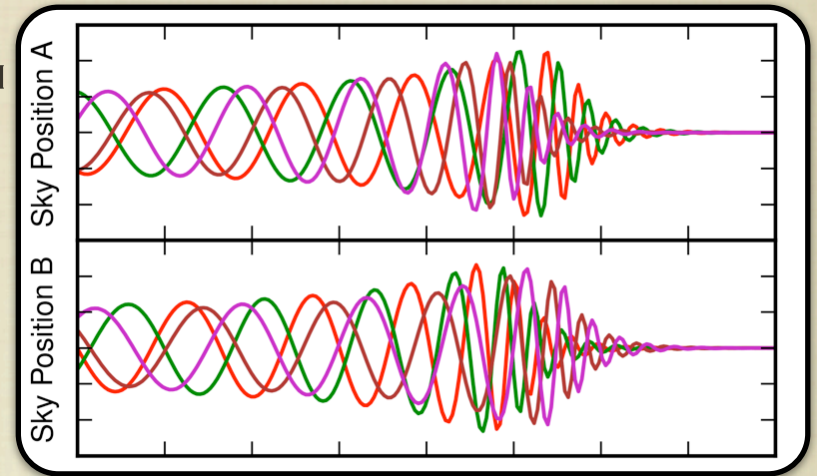
- MEMORY TRANSFERS
- CUFFT
- cuCHISQ

GPGPU ACCELERATION EFFORTS: GWTOOLS

- BASED ON OPENCL FOR PORTABILITY
- ACHIEVED THEORETICAL SPEEDUP FOR SNR USING OVERLAP-SAVE ALGORITHM
- PROTOTYPING CODE FOR NEW ALGORITHMS
 - WAVEFORM GENERATION
 - SNR CALCULATION
 - SIGNAL-BASED VETOES
 - MAXIMIZATION AND TRIGGER PRODUCTION
 - TRIGGER CLUSTERING
- KERNELS DEVELOPED WILL BE INCORPORATED INTO OTHER TOOLKITS

SEARCHING FOR INSPIRAL SIGNALS: COHERENT SEARCH

- COHERENTLY COMBINE DATA FROM DETECTOR NETWORK
- EACH SKY LOCATION REQUIRES:
 - N TIME-SHIFTS AND IFFTS
 - RECOMBINATION
- $\sim 10^3$ - 10^4 SKY LOCATIONS FOR NETWORK OF N ADVANCED DETECTORS
- HUNDREDS OF GFLOPS BECOME FEW PFLOPS
- A HIERARCHICAL APPROACH COULD REDUCE COST



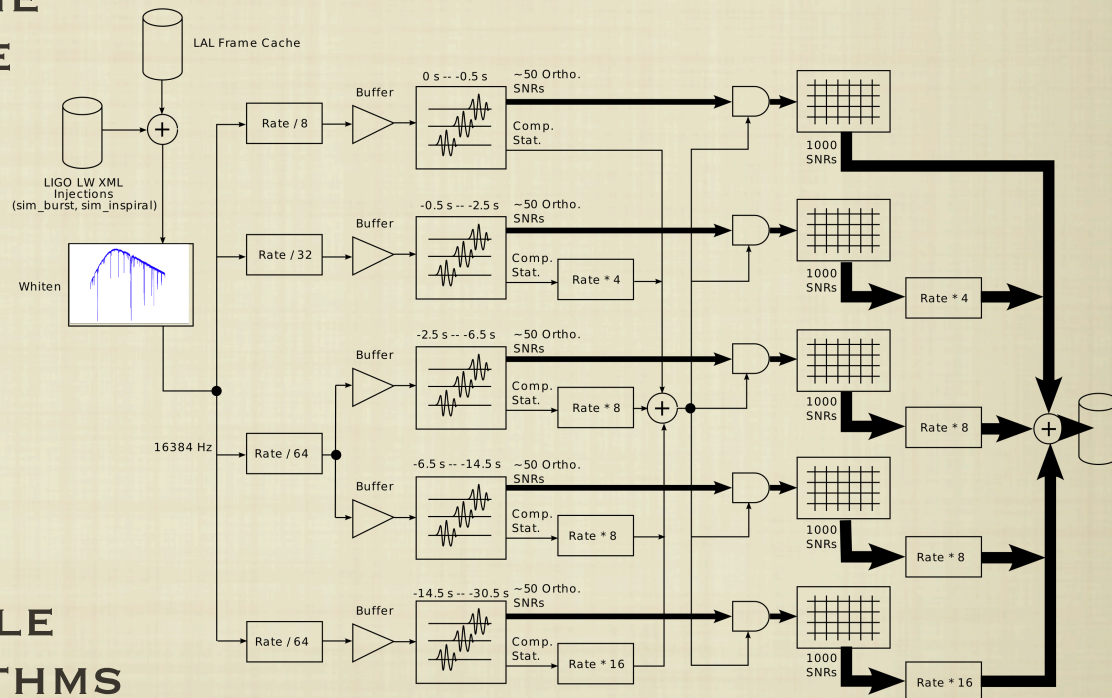
LOW-LATENCY PIPELINES

- HIGH PAYOFF COULD COME FROM PROMPT TELESCOPE POINTING

- REDUCE LATENCY AS MUCH AS POSSIBLE

- EX: GSTLAL BASED ON GSTREAMER MULTIMEDIA FRAMEWORK AND LSC'S ALGORITHMS LIBRARY

- FILTER ENGINE SWAPPABLE WITH DIFFERENT ALGORITHMS



GPU ACCELERATION EFFORTS: LOW LATENCY WITH LLOID

- DECREASED LATENCY MEANS INCREASED COMPUTATIONAL COST
- COMPUTATIONAL COST REDUCTIONS:
 - MULTIRATE FILTERING
 - ALSO IMPLEMENTED WITHIN MBTA
 - SIGNIFICANCE-BASED FILTERING
 - INVESTIGATING BENEFIT OF PORTING INDIVIDUAL TOOLS TO GPU

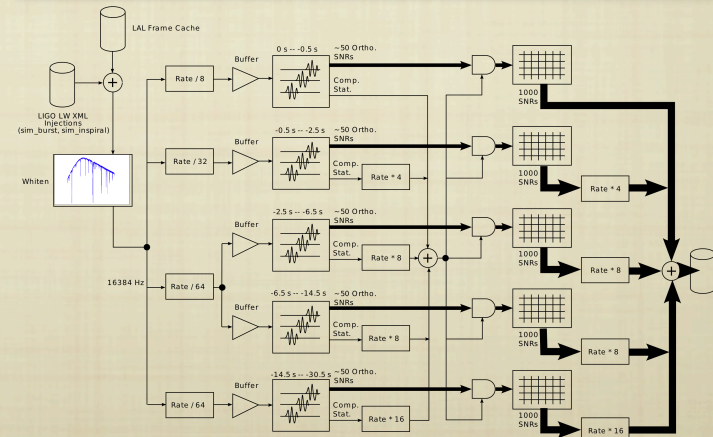
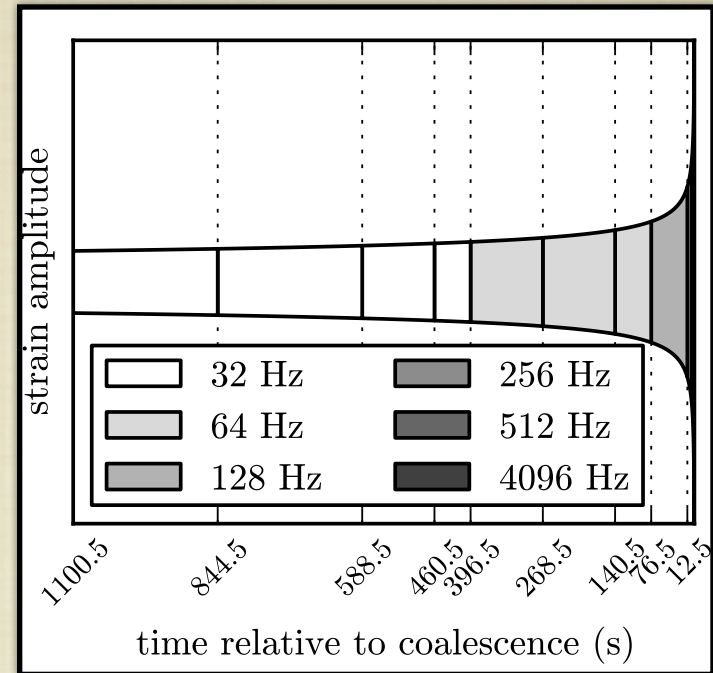


Figure Credit: K. Cannon, ..., DK, et al, ApJ 748, 136 (2012)

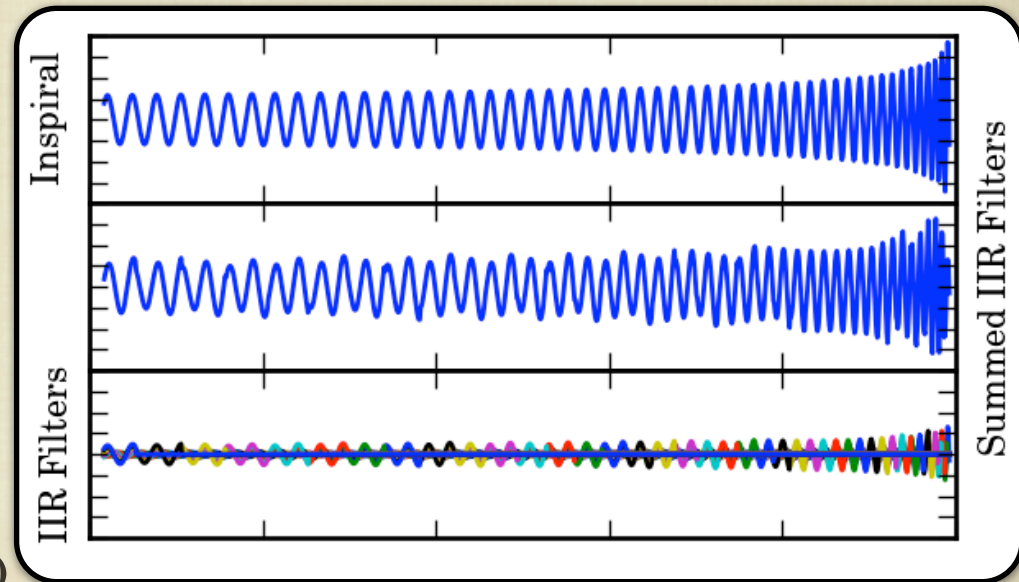
GPU ACCELERATION EFFORTS: LOW LATENCY WITH SPIIR

- WAVEFORMS APPROXIMATED BY SUMMED PARALLEL IIR FILTERS

- COMPUTATIONALLY LESS EXPENSIVE THAN FIR REPRESENTATION

- 1 (COMPLEX) ADD + 1 (COMPLEX) MUL PER SAMPLE PER IIR FILTER (EACH IS A SINGLE POLE)

- HUNDREDS OF IIR FILTERS PER WAVEFORM



GPU ACCELERATION EFFORTS: LOW LATENCY WITH SPIIR

- GOOD TARGET FOR ACCELERATION

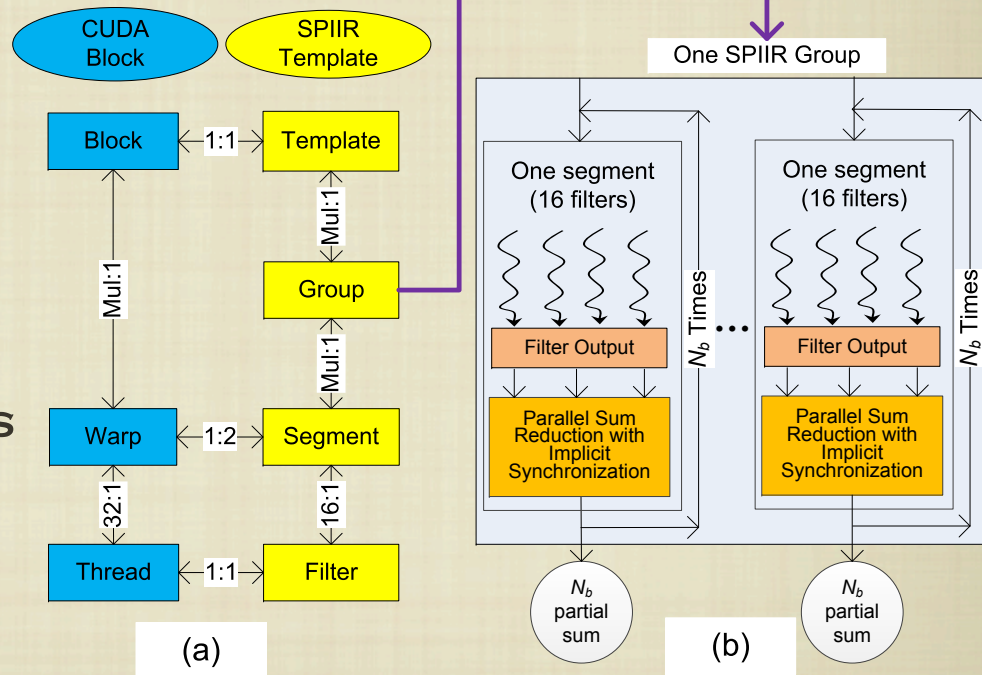
- PRO: ALGORITHM HIGHLY PARALLELIZED

- CON: REQUIRES SYNCHRONIZED OUTPUTS

- SUB-CALCULATIONS REGROUPED FOR PERFORMANCE ENHANCEMENTS

- SYNCH EVENTS REDUCED BY NUMBER OF SAMPLES PER BATCH

- ALIGN NUMBER OF FILTERS PER SEGMENT WITH SIZE OF WARP



GPU ACCELERATION EFFORTS: LOW LATENCY LIMITATIONS

- CHOICE OF IMPLEMENTATION
 - MAKE GPU KERNELS FOR SMALL TOOLS OR LARGE BLOCKS?
 - WHERE IS THE LINE BETWEEN CPU AND GPU COMPUTATIONS?
- OBSERVED LIMITS IN SPIIR
 - FEW OPERATIONS PER DATA SAMPLE
 - MORE ADDS THAN MULS
 - COMPARABLE TO INEFFICIENCIES IN FFT ALGORITHMS

GPU ACCELERATION EFFORTS: USING PYFFT, PYCUDA, PYOPENCL

- BUILD A FRAMEWORK WITH TRANSPARENT ACCELERATION FOR SIMPLE ALGORITHMS
- SOME INVESTIGATIONS COMPUTE FILTER OUTPUT 10^3 - 10^6 TIMES WITH VARIOUS PARAMETERS
 - TEMPLATE BANK COVERING STUDIES
 - PARAMETER ESTIMATION SEARCHES
- AMOUNT OF ACCELERATION ALGORITHM DEPENDENT
 - MEMORY TRANSFERS
 - GENERATION OF WAVEFORMS
- COULD BE USED TO INVESTIGATE AND IMPLEMENT NEW PIPELINES

THE END



QUESTIONS?