

FPGA-based Acceleration of SKA Pulsar Search Modules Using OpenCL

Haomiao Wang and Tyrone Sherwin

12/07/2018

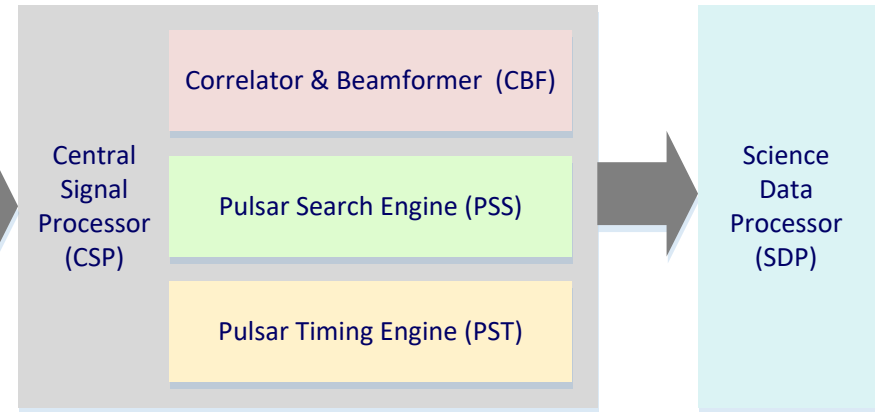


**PARALLEL AND
RECONFIGURABLE
COMPUTING LAB**

Pulsar Search



Pre-processed
Digital Data



[SKA-Mid Africa Widefield](#) Image credit: SKA Organisation

- Test of general relativity and gravitational waves => **Pulsar searches**
- Pulsar Search Engine (PSS)
 - Single Pulse Search
 - Time-domain acceleration search
 - Frequency-domain acceleration search

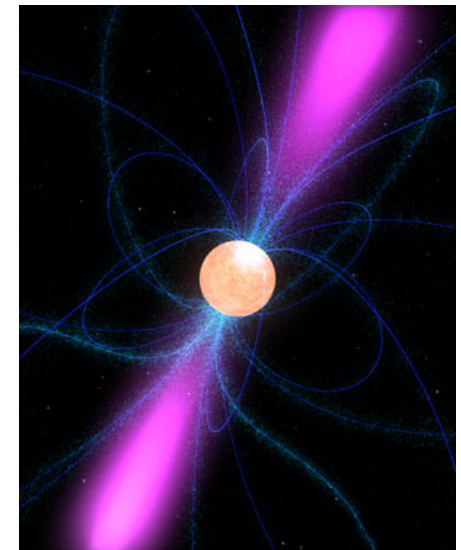
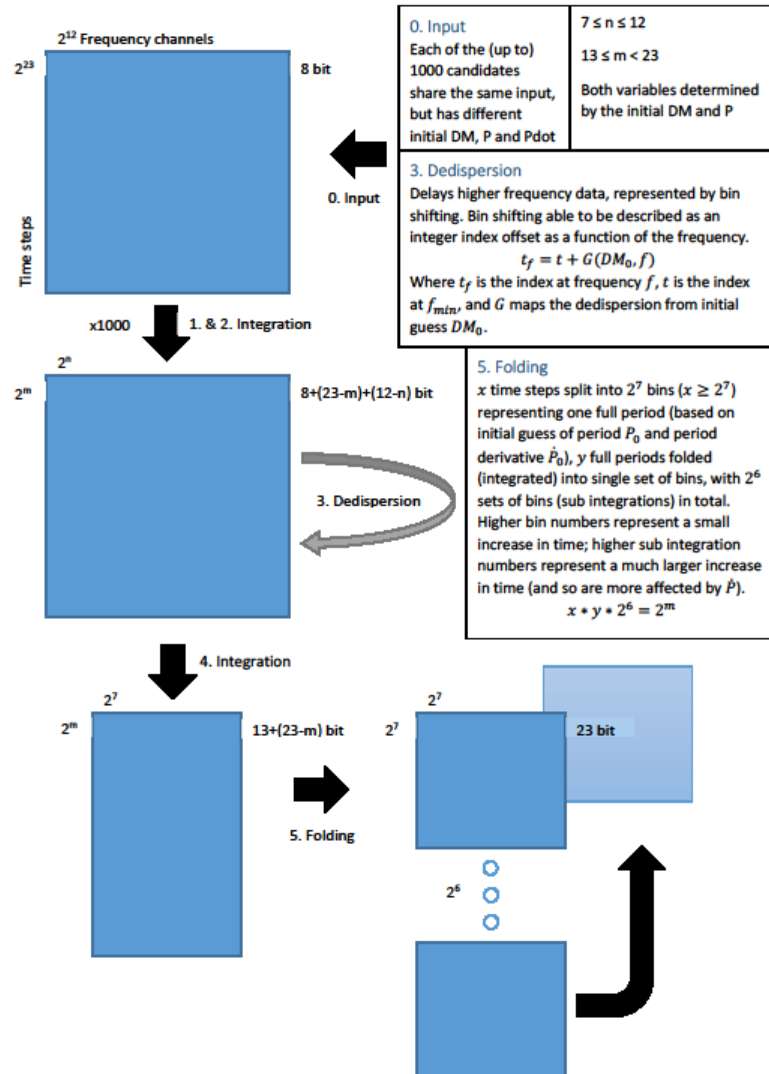


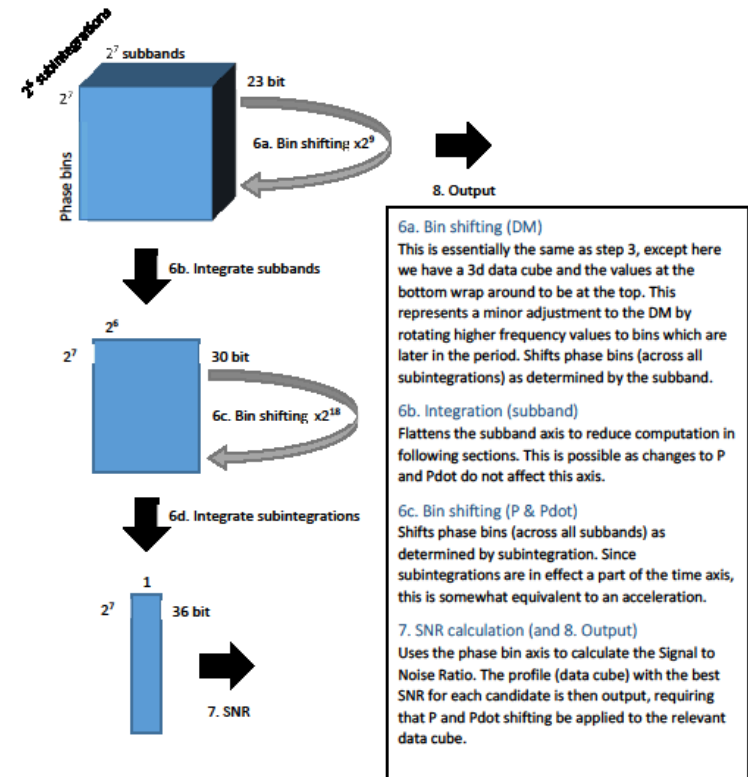
Image credit: NASA

Example-FLDO

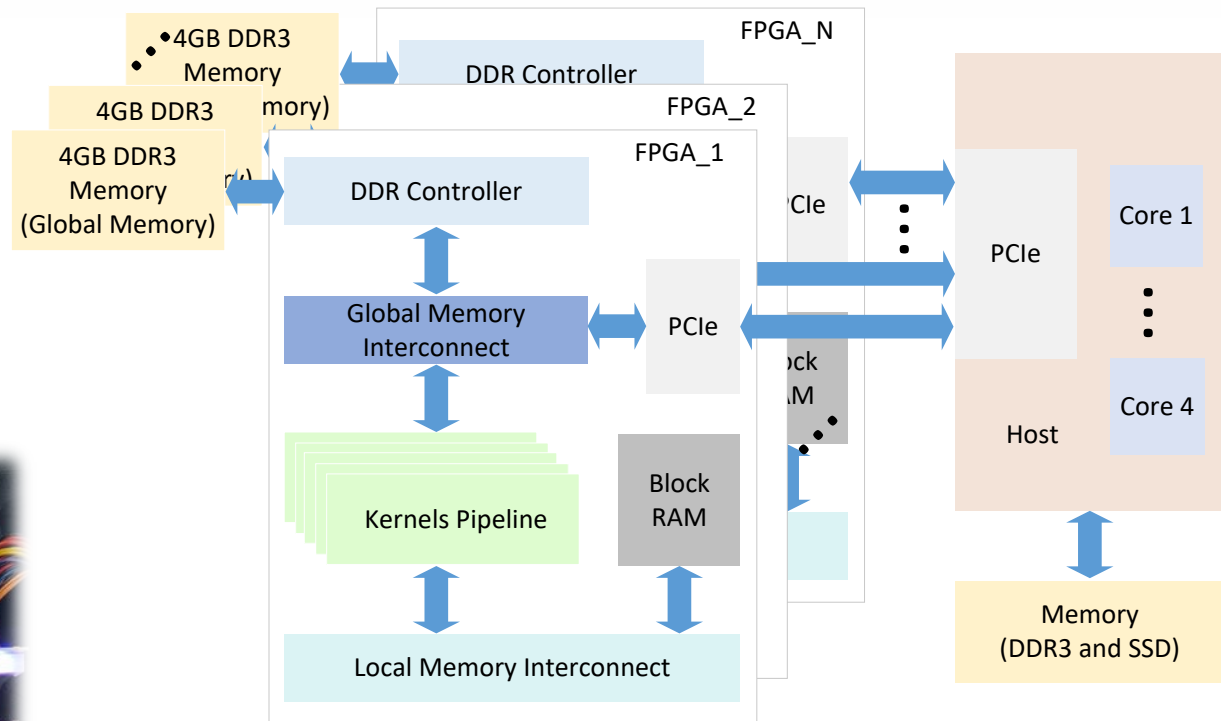
1-5 Preprocessing (Dedispersion & Folding)



6-7 Optimisation



OpenCL for FPGA



GeForce GTX TITAN X



520N

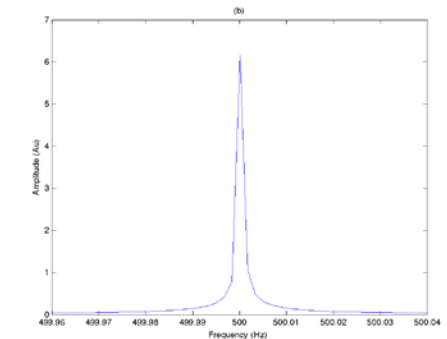
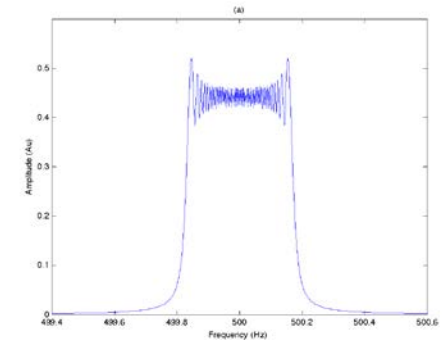
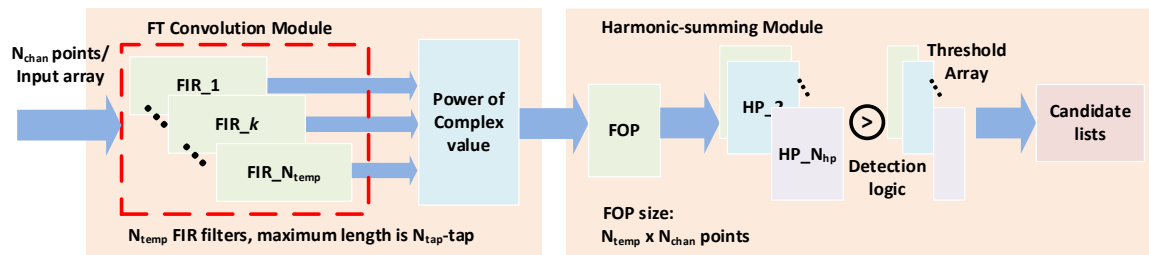
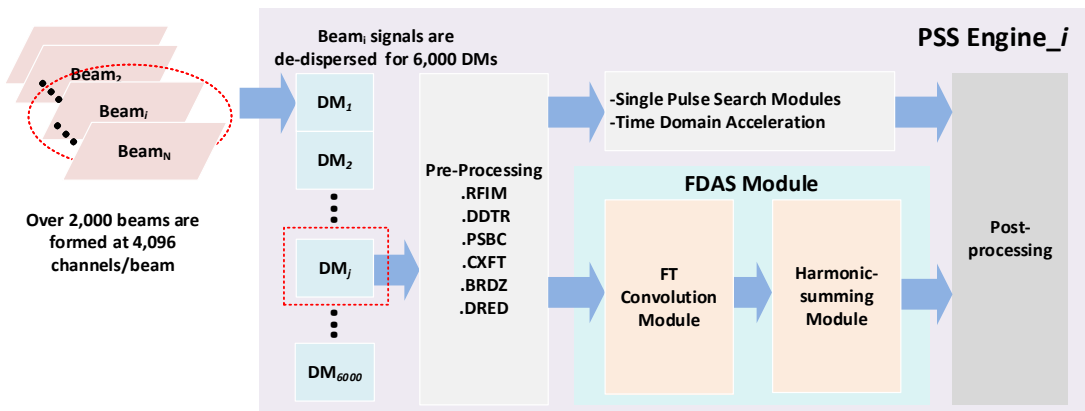
Intel Stratix 10 FPGA

Designing pipeline stages of Pulsar Search Pipeline

- Fourier Domain Acceleration Search (FDAS)
 - FT convolution
 - Harmonic summing
- RFI Mitigation (RFIM)
 - Median filtering

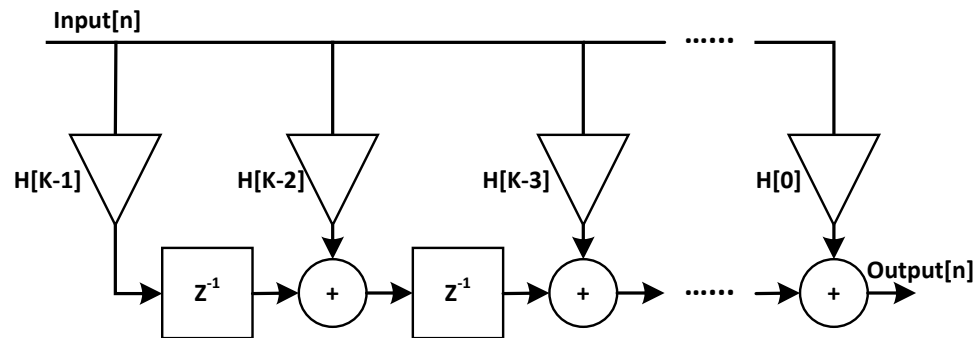
- Search for (binary) pulsars with constant frequency derivatives
- Matched filtering (a group of FIR filters with different lengths, complex floating-point)

$$A_{r_0} = \sum_{k=[r_0]-\frac{m}{2}}^{[r_0]+\frac{m}{2}} A_k A_{r_0-k}^*$$



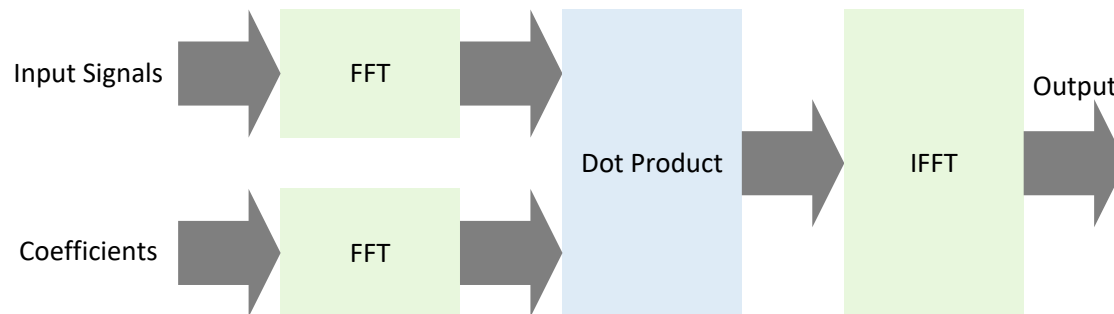
- Time-domain FIR Filter (**TDFIR**)

$$y_m[i] = \sum_{k=0}^{K-1} x_m[i-k]h_m[k], \quad \text{for } i = 0, 1, \dots, N-1$$

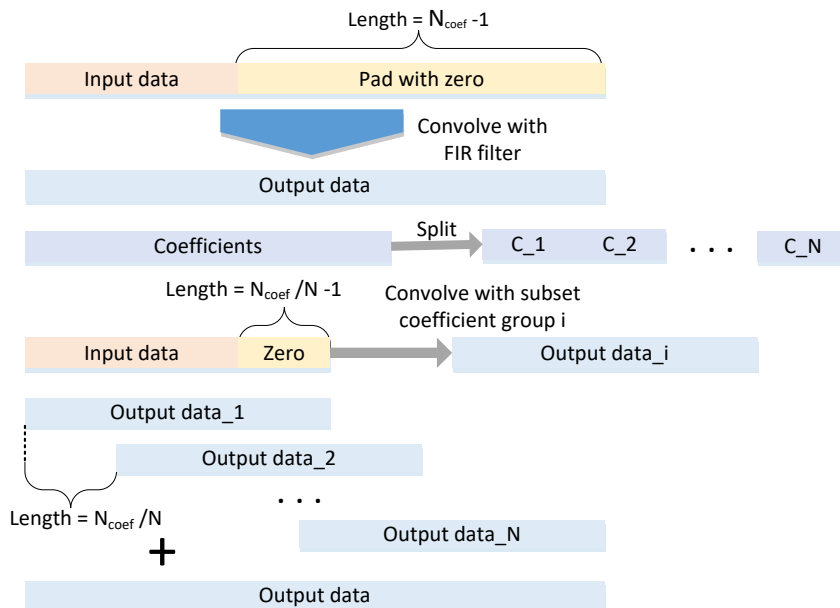


- Frequency-domain FIR Filter (**FDFIR**)

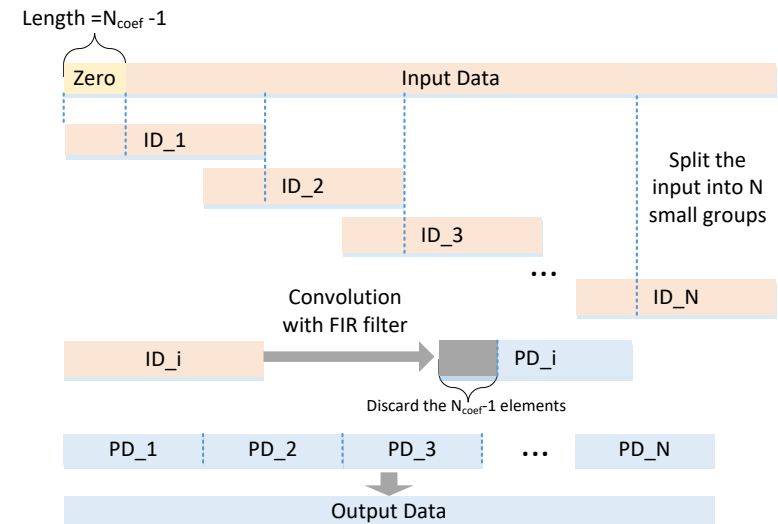
$$x * h = F^{-1}\{F\{x\} \cdot F\{h\}\}$$



- **Overlap-add algorithm (OLA)**
 - split the **coefficient array** (TDFIR)
- **Overlap-save algorithm (OLS)**
 - split the **input array** (FDFIR)

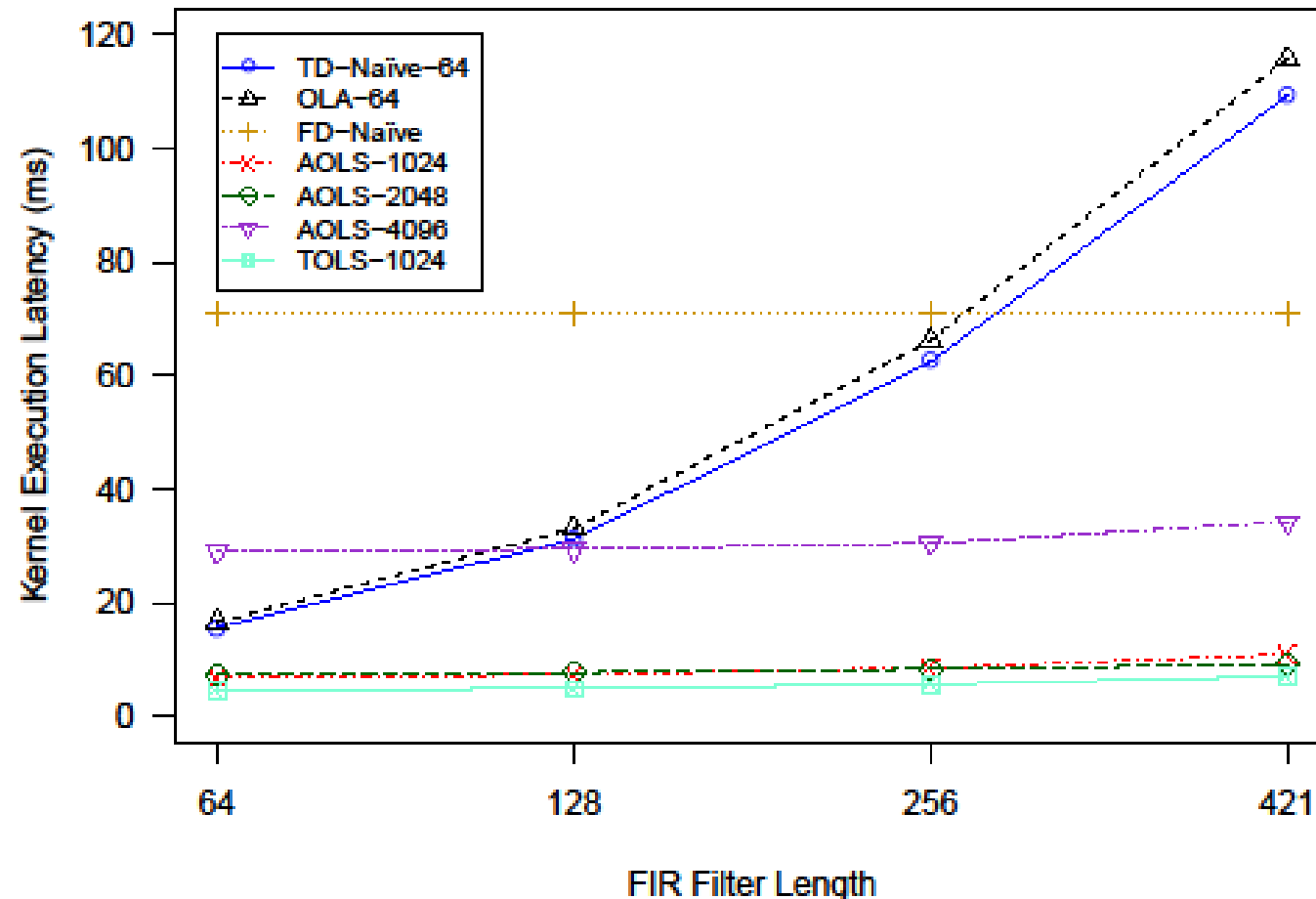


Overlap-add Algorithm



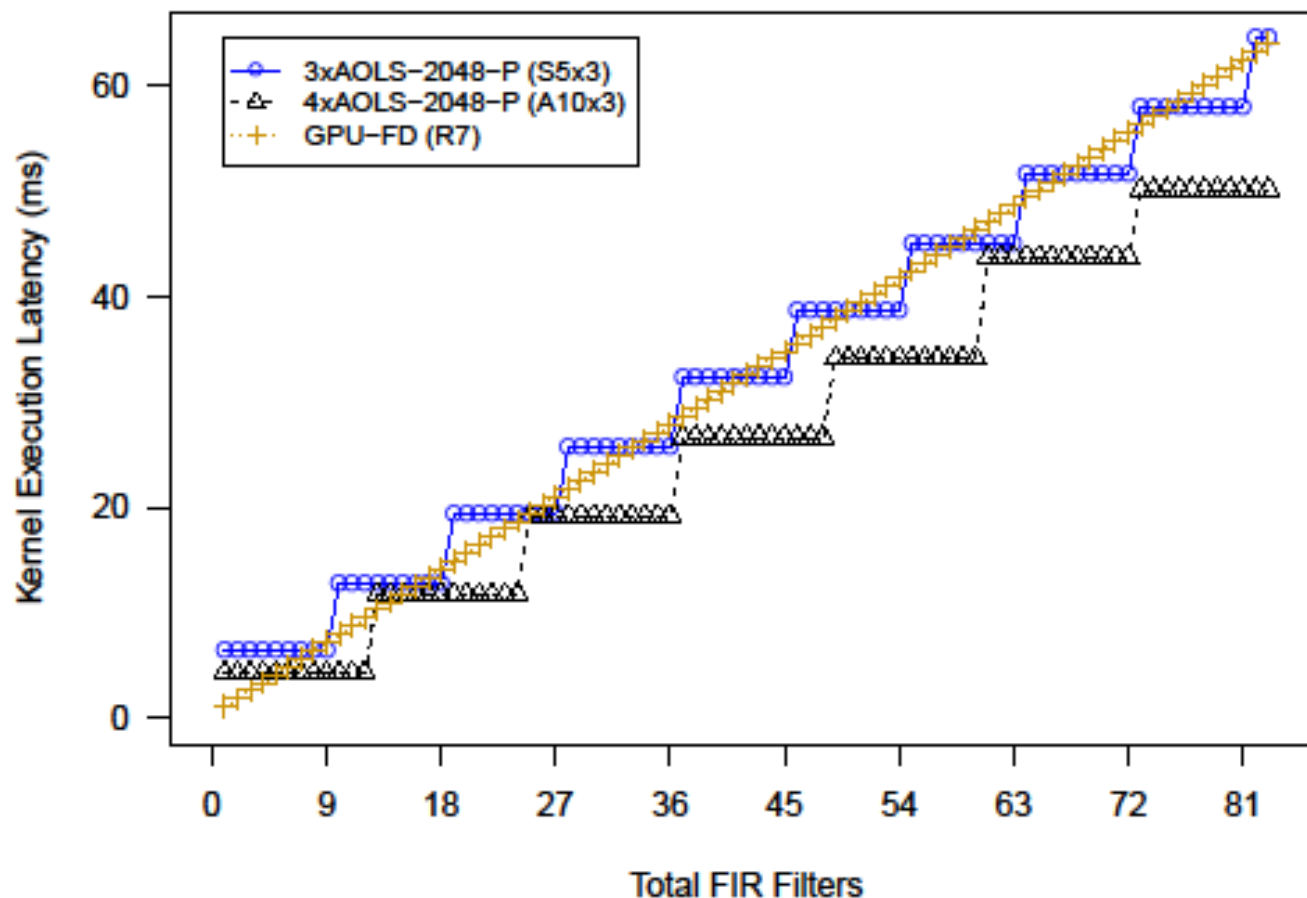
Overlap-save Algorithm

Latencies of a single FPGA (Intel Stratix V A7) in processing same input array using 7 different methods

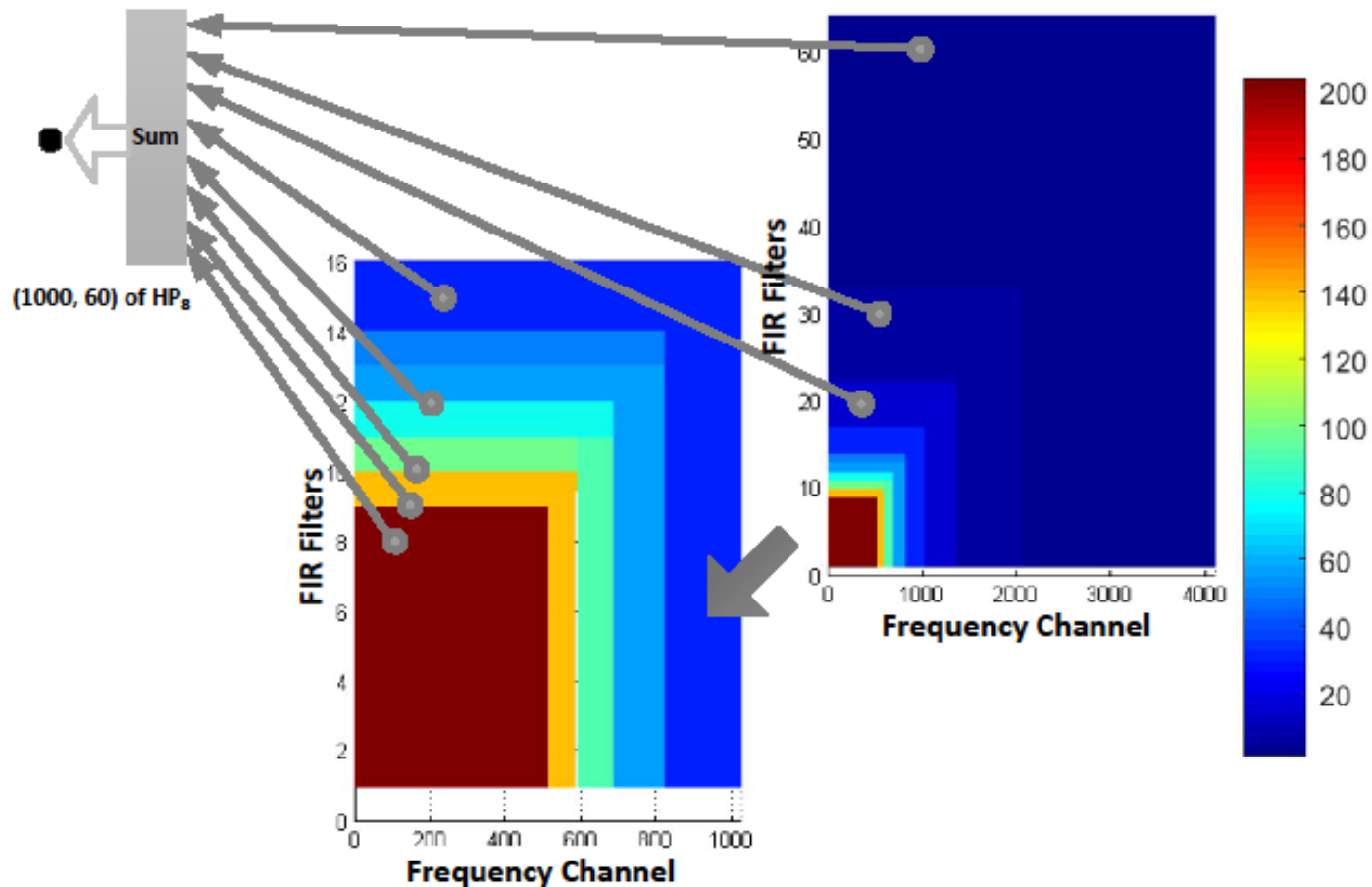


Latencies of a single GPU (AMD Radeon R7) and 3 FPGAs (Intel Stratix V and Arria 10)

3xFIRs on one Stratix V image and 4xFIRs on one Arria 10 image



Optimizing irregular accesses to off-chip memory



- **SingleHP method**

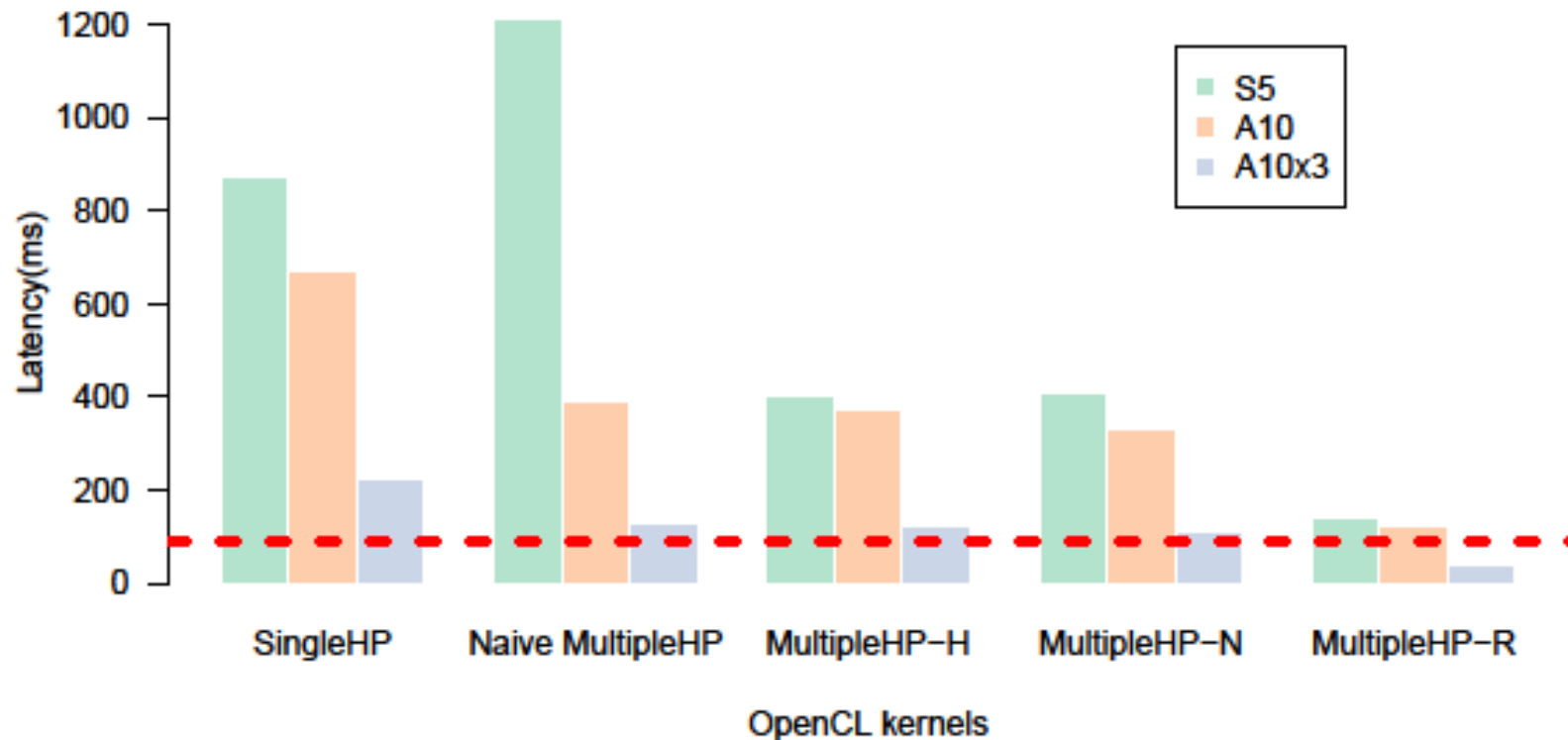
Processing a single harmonic plane at a time

- **MultipleHP method**

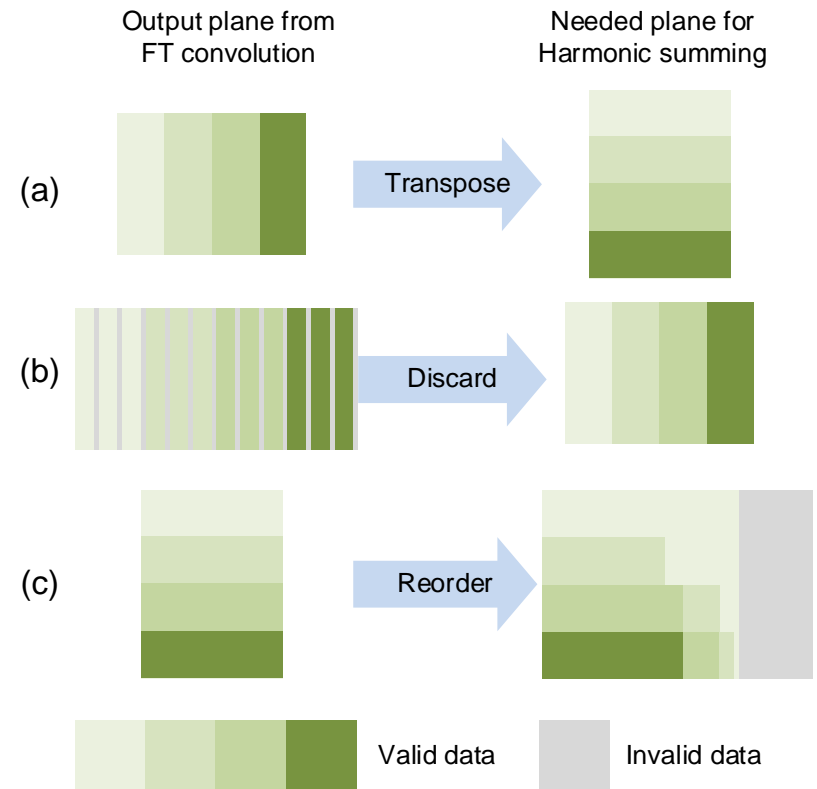
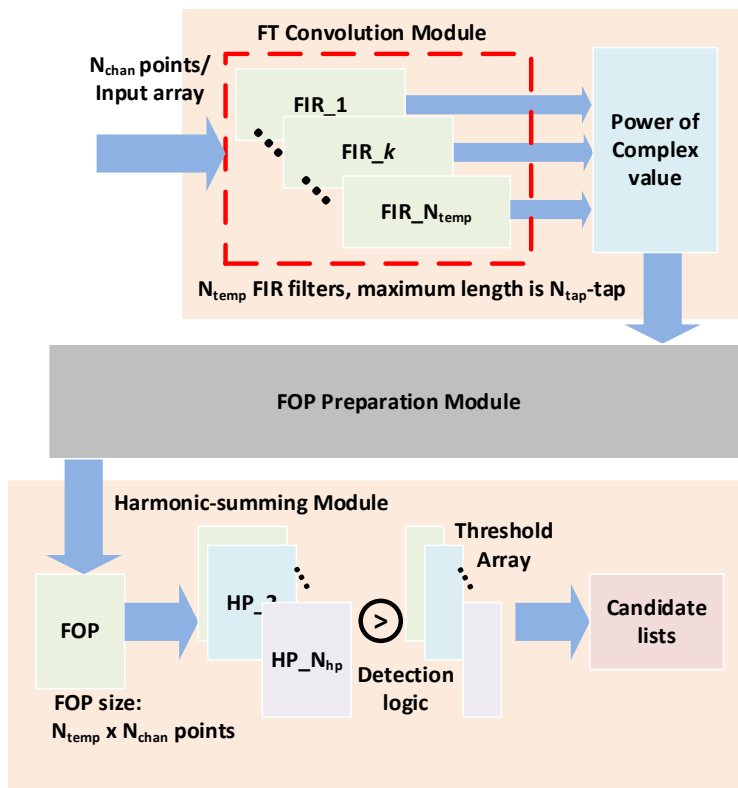
Processing multiple harmonic planes in parallel

- Naïve MultipleHP
- **MultipleHP-H**: preloading data with high toughing frequency
- **MultipleHP-N**: loading necessary data to calculate a block of points
- **MultipleHP-R**: based on MultipleHP-N and the FOP is reordered to achieve stream mode

MultipleHP-R method performs best among the proposed methods and it meets the required time limitation

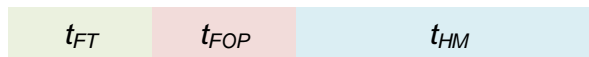


An FOP preparation module is added to connect the FT convolution module and harmonic summing module

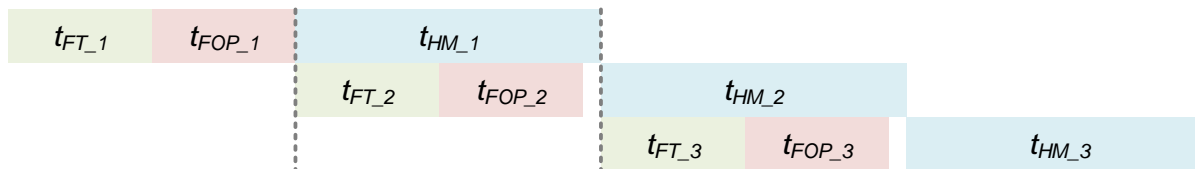


$$\text{Max}(t_{FT}, t_{FOP}, t_{HM}) \geq t_{FDAS} / 2$$

Single



Multiple

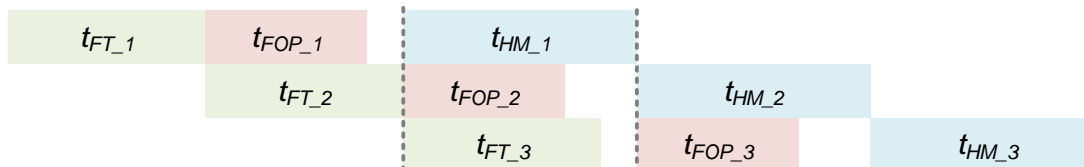


$$\text{Max}(t_{FT}, t_{FOP}, t_{HM}) < t_{FDAS} / 2$$

Single



Multiple



Form of stencil computation:

- Recomputing an array element based on neighbour values
- In median filter, target value is replaced by median of all values in a window around target value
- => sorting all values, taking middle value

Example:

- 2D array, window (stencil) is 3 x 3

3	7	2
1	2	4
6	4	0



0 1 2 2 3 4 4 6 7

Large window size.

- 63 frequency channels, 5MHz, n .
- 1023 time steps, 65ms, m .

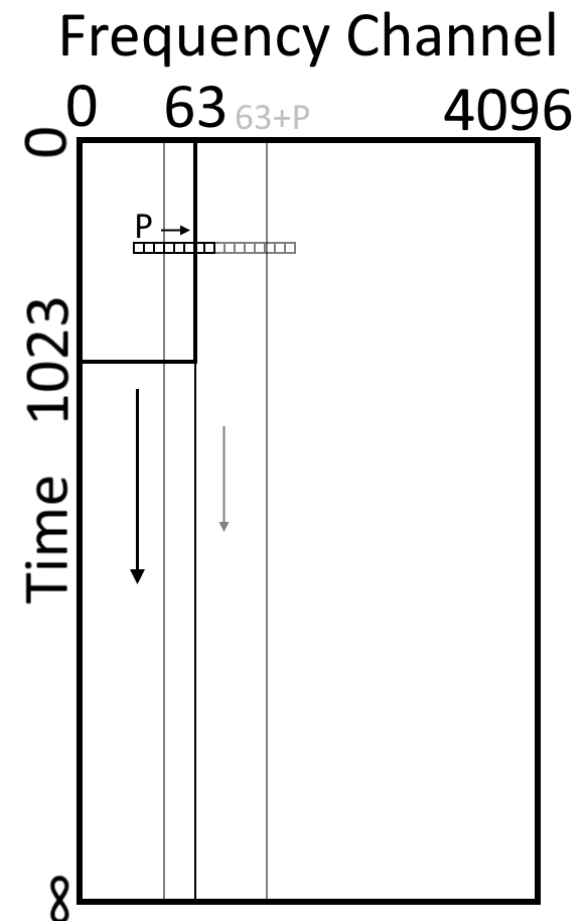
Array Size.

- Frequency coverage 4096 channels, 300MHz.
- Time coverage potentially infinite, steps of $64\mu s$.

Throughput requirements.

- 8bit input data.
- 4096 channels every $64\mu s$.
- 64 Million values per second.

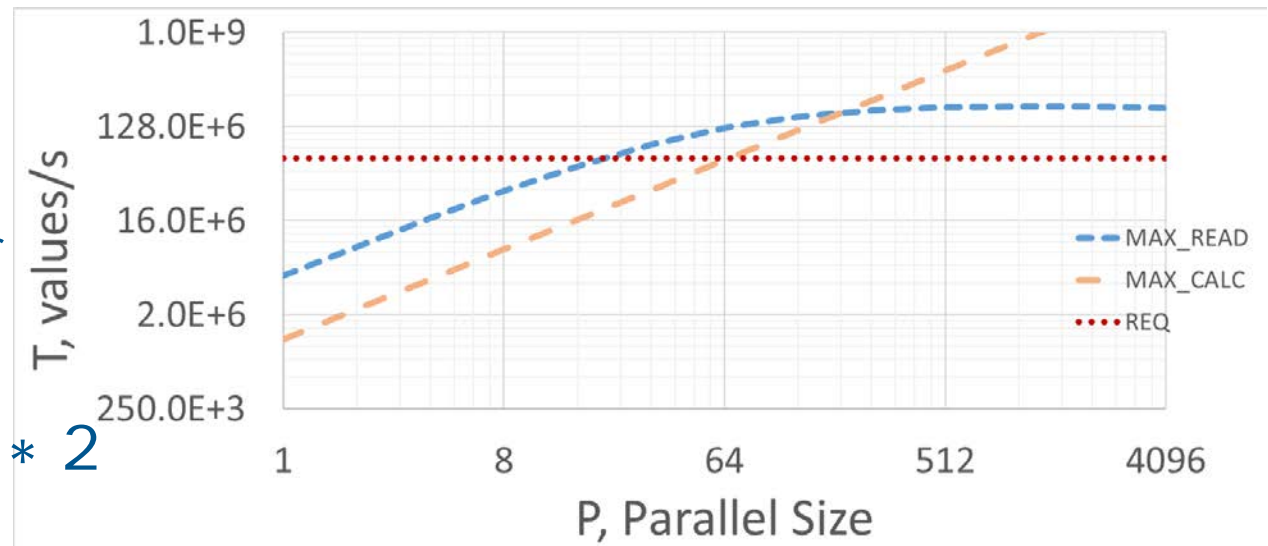
- Choice of algorithmic approach dominates performance
- Using best of various state-of-the-art methods
- Histograms
 - 256 bins for 8bit data
 - Each bin counts that value in current window.
- Sliding
 - Adjacent windows overlap in data
 - Reuse previous window, replacing oldest row with next row



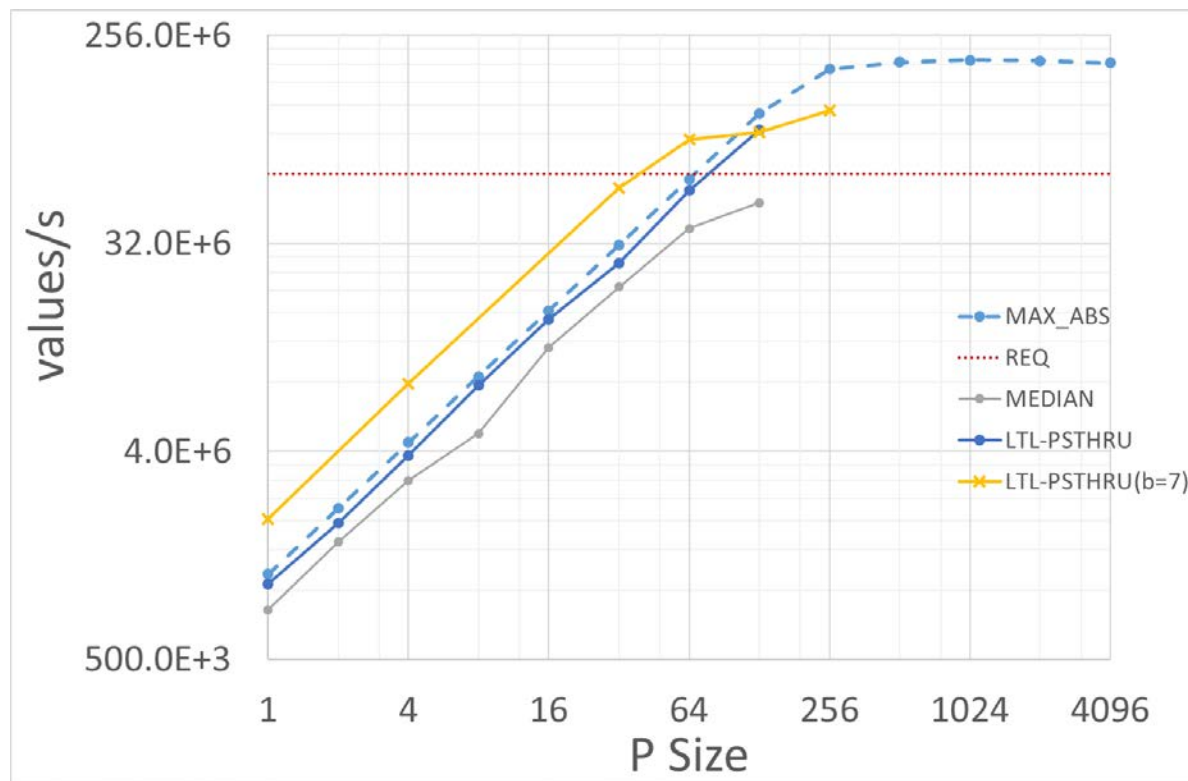
- Developed to determine theoretical performance.
- Limitations placed on loads per cycle and calculations per cycle.
- $n=63$ requires $P > 64$.

T is throughput, f is frequency, L is loads per row,
 I is loads per cycle

$$T = (V/c) * f$$
$$V/c = I/L * P$$
$$L = (n + P - 1) * 2$$

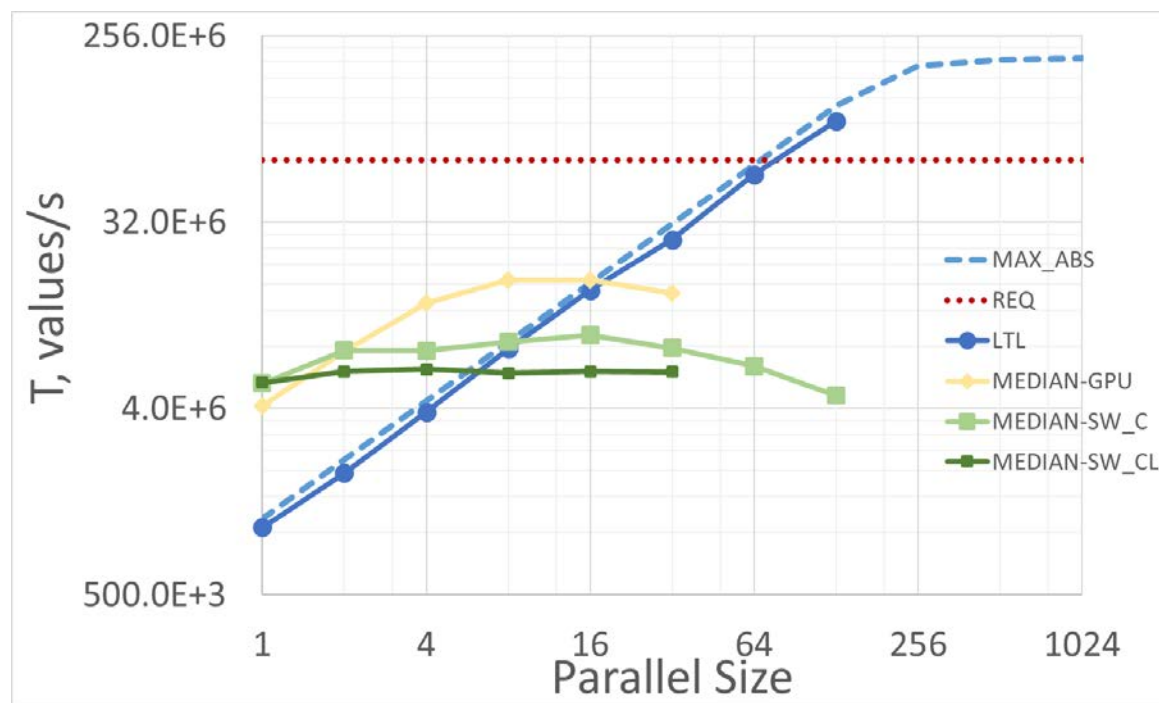


- Pipeline balancing
- Able to achieve required throughput
- LTL-PSTHRU($P=128$) reaching 99.3Mvalue/s.
 - Theoretical performance of 117.7M, limited by calculation rate.
 - 84.4% → Room for improvement.
- LTL-PSTHRU($P=128$, $b=7$) reduces bit depth for higher performance



Comparison with GPU&CPU

- CPU & GPU versions not able to perform as well – despite OpenCL adjustments
- CPU appears not to operate in parallel
- GPU quickly drops off
 - Appears similar in shape to roofline models
 - Suggests possible resource constraints
- FPGA starts slower, much better scaling with P



- Investigated the optimization of FT convolution and harmonic-summing and median filtering modules
- Pipelining modules on single device (no host-FPGA communication)
- Using multiple FPGAs
- High-level approach is employed to implement the optimized methods
 - Covered large design space
 - Easy porting and sharing with partners

- Designing FPGA implementation for FLDO –
Candidate Folding and Optimisation

Long term:

- Creating OpenCL library (target FPGA) for all
(computationally intensive) Pulsar Search
Pipeline stages

Thanks

Q&A