

Reduced-Precision Acceleration of Radio-Astronomical Imaging on Reconfigurable Hardware

Stefano Corda

stefano.corda@epfl.ch

07-07-2022

Scientific Computing Accelerated on FPGAs, Maison de la Simulation (Saclay)



SKAO

SKAO Mission and Vision

Our mission

Our mission is at the core of what we are here to deliver as an organisation. It outlines our deliverables, while setting a bold ambition for us to achieve with regards to our impact in the world.

“The SKAO’s mission is to build and operate cutting-edge radio telescopes to transform our understanding of the Universe, and deliver benefits to society through global collaboration and innovation.”

The SKA Observatory Convention defines the purpose of the SKAO as to facilitate and promote a global collaboration in radio astronomy with a view to the delivery of transformational science.

Our vision

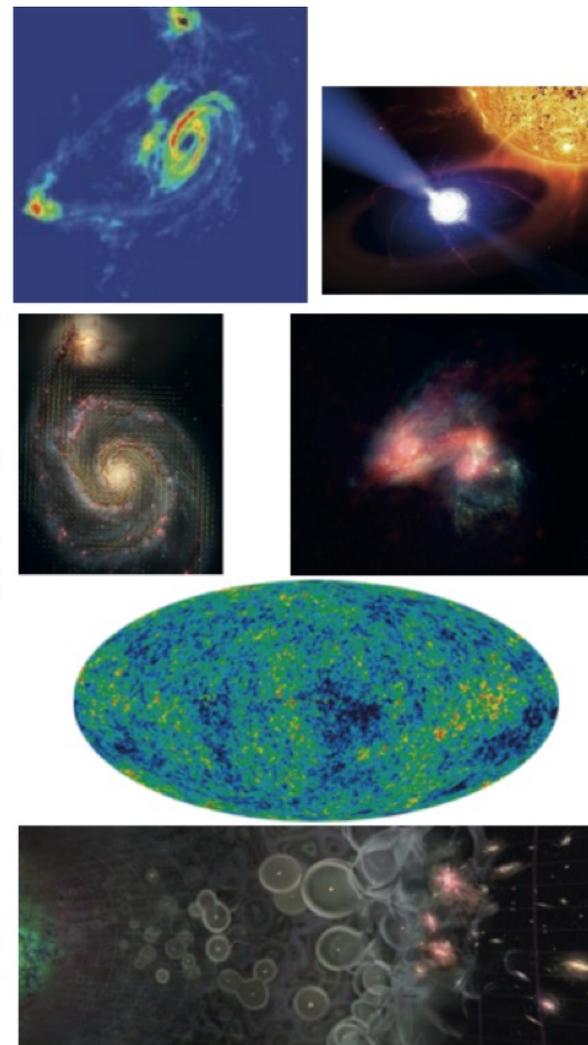
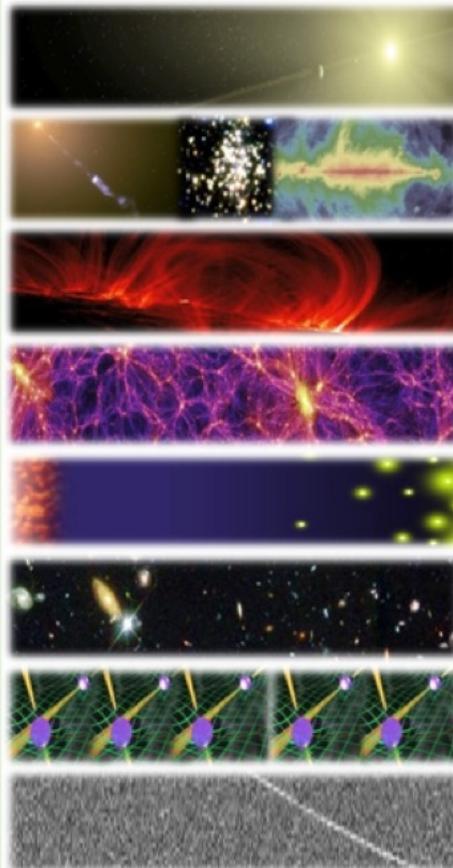
“The SKAO is one observatory, with two telescopes, on three continents; a 21st century observatory and an inter-governmental organisation with sustainability and respect to all our communities at its heart, driven by a commitment to fundamental science and technology.”



SKAO

SKAO Key Science Drivers

- **The Cradle of Life & Astrobiology**
 - *How do planets form? Are we alone?*
- **Strong-field Tests of Gravity with Pulsars and Black Holes**
 - *Was Einstein right with General Relativity?*
- **The Origin and Evolution of Cosmic Magnetism**
 - *What is the role of magnetism in galaxy evolution and the structure of the cosmic web?*
- **Galaxy Evolution probed by Neutral Hydrogen**
 - *How do normal galaxies form and grow?*
- **The Transient Radio Sky**
 - *What are Fast Radio Bursts? What haven't we discovered?*
- **Galaxy Evolution probed in the Radio Continuum**
 - *What is the star-formation history of normal galaxies?*
- **Cosmology & Dark Energy**
 - *What is dark matter? What is the large-scale structure of the Universe?*
- **Cosmic Dawn and the Epoch of Reionization**
 - *How and when did the first stars and galaxies form?*



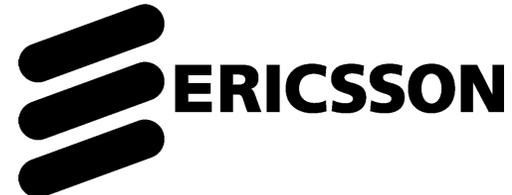
Reduced-Precision Acceleration of Radio-Astronomical Imaging on Reconfigurable Hardware

Stefano Corda, Bram Veenboer, Ahsan Javed Awan, John Romein,
Roel Jordans, Akash Kumar, Albert-Jan Boonstra, Henk Corporaal



TECHNISCHE
UNIVERSITÄT
DRESDEN

ASTRON



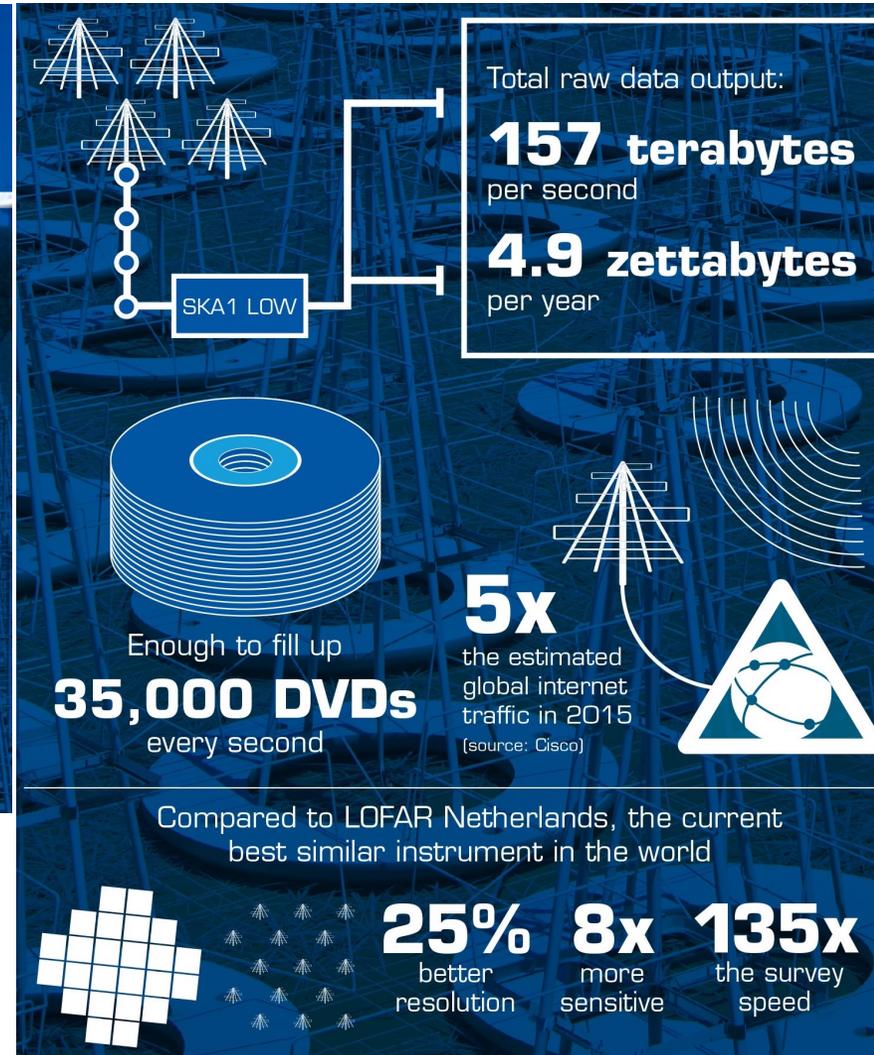
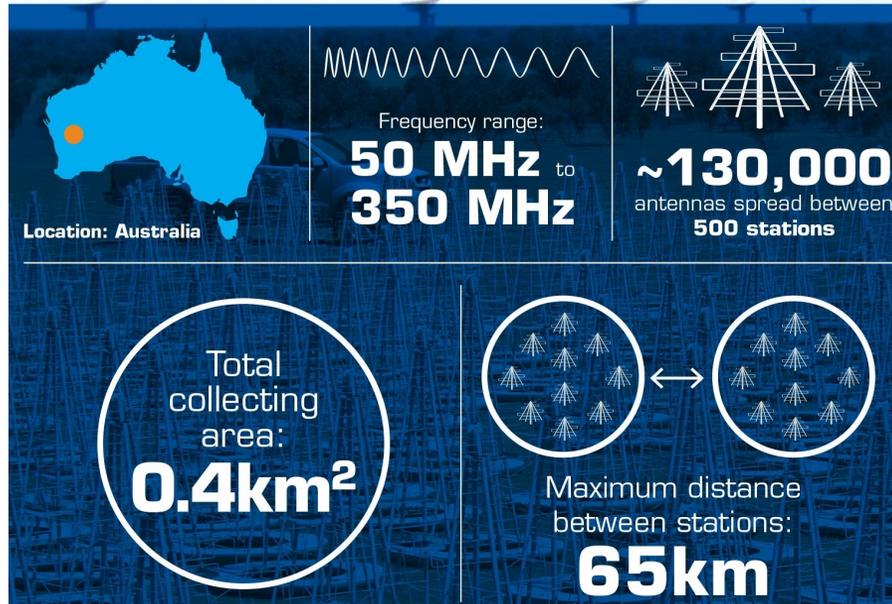
Square Kilometre Array Challenge

SKA (Square Kilometre Array) requirements (per Science Data Processor (SDP) site):

- ~ 260 PFlop/s
- ~ 157 TB/s
- ~ 5 MW

SKA1 LOW - the SKA's low-frequency instrument

The Square Kilometre Array (SKA) will be the world's largest radio telescope, revolutionising our understanding of the Universe. The SKA will be built in two phases - SKA1 and SKA2 - starting in 2018, with SKA1 representing a fraction of the full SKA. SKA1 will include two instruments - SKA1 MID and SKA1 LOW - observing the Universe at different frequencies.



HPC and FPGAs

Rank	System	Rpeak (PFLOP/s)	Power (MW)	HPCG %
1	Frontier	1685.65	21.10	N. A.
2	Fugaku	537.21	29.90	3.62
3	LUMI	214.35	2.94	1.27
4	Summit	200.79	10.10	1.97
5	Sierra	125.71	7.44	1.90



“When a large-scale HPC system wastes only 1% to 10% of its computing cycles, it wastes energy that could support a small city.”

Why FPGAs?:

- Floating-point support.
- Custom HW for domain-specific applications.
- High-level synthesis tools (reduced programming effort).

Outline

- Problem statement and contributions
- Background
- Methodology
- Application analysis
- Accelerator architecture
- Architecture evaluation and discussion
- Related work
- Conclusions and future work

Problem Statement

- High-Performance requirements (SDP → almost ExaFlop/s).
- Image-Domain Gridding (IDG) is highly efficient in single-precision on GPUs.
- FPGA technology and toolchain improvement.
- Reduced precision for noise-tolerant applications.

Challenges:

- Is reduced precision applicable to the radio-astronomical imaging domain?
- Can we profit from low precision using FPGAs for radio-astronomical imaging?

R. V. van Nieuwpoort et al., "Correlating radio astronomy signals with many-core hardware", IJPP 2010

R. Jongerius et al., "An end-to-end computing model for the square kilometre array", Computer 2014

B. Veenboer et al., "Radio-astronomical imaging on graphics processor", ASCOM 2020

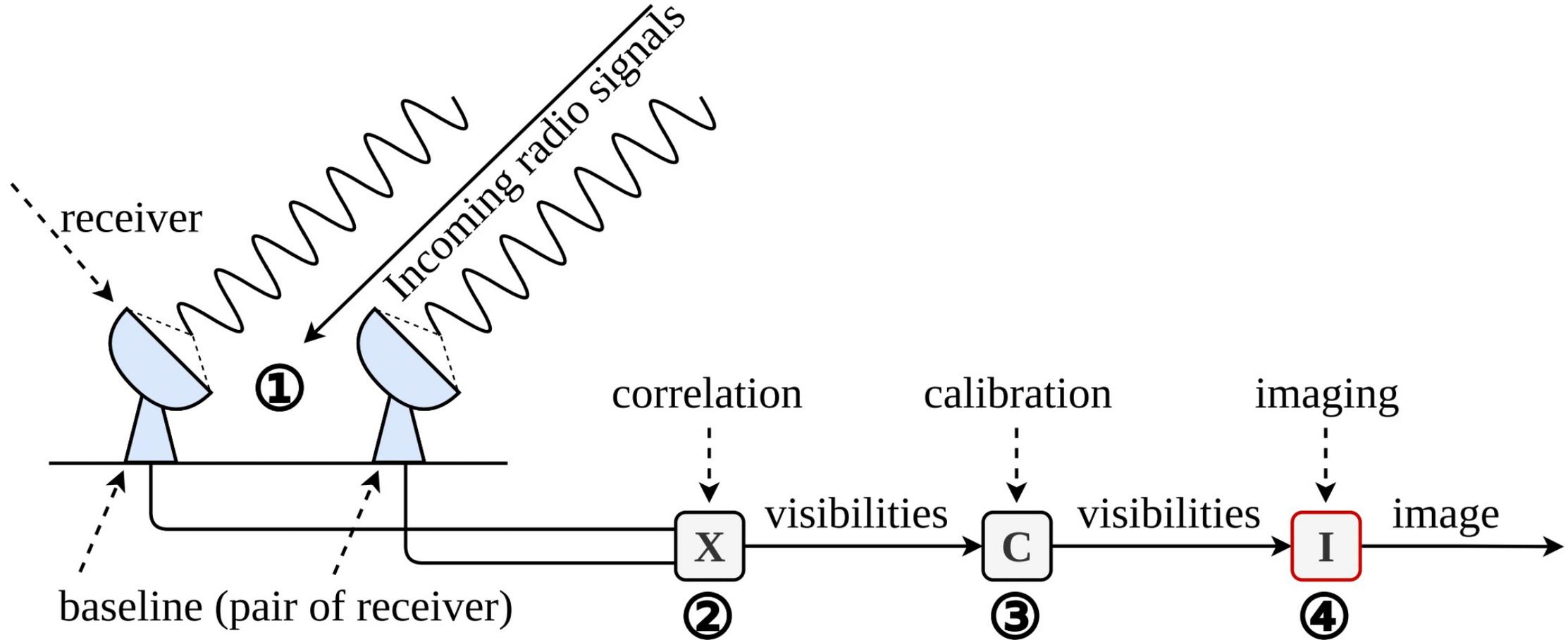
Intel, "Enabling High-Performance Floating-Point Designs," <https://www.intel.com/content/dam/www/programmable/us/en/pdfs/literature/wp/wp-01267-fpgas-enable-high-performance-floating-point.pdf>

S. Cherubin et al., "Tools for Reduced Precision Computation: A Survey", ACM Comput. Surv. 2020

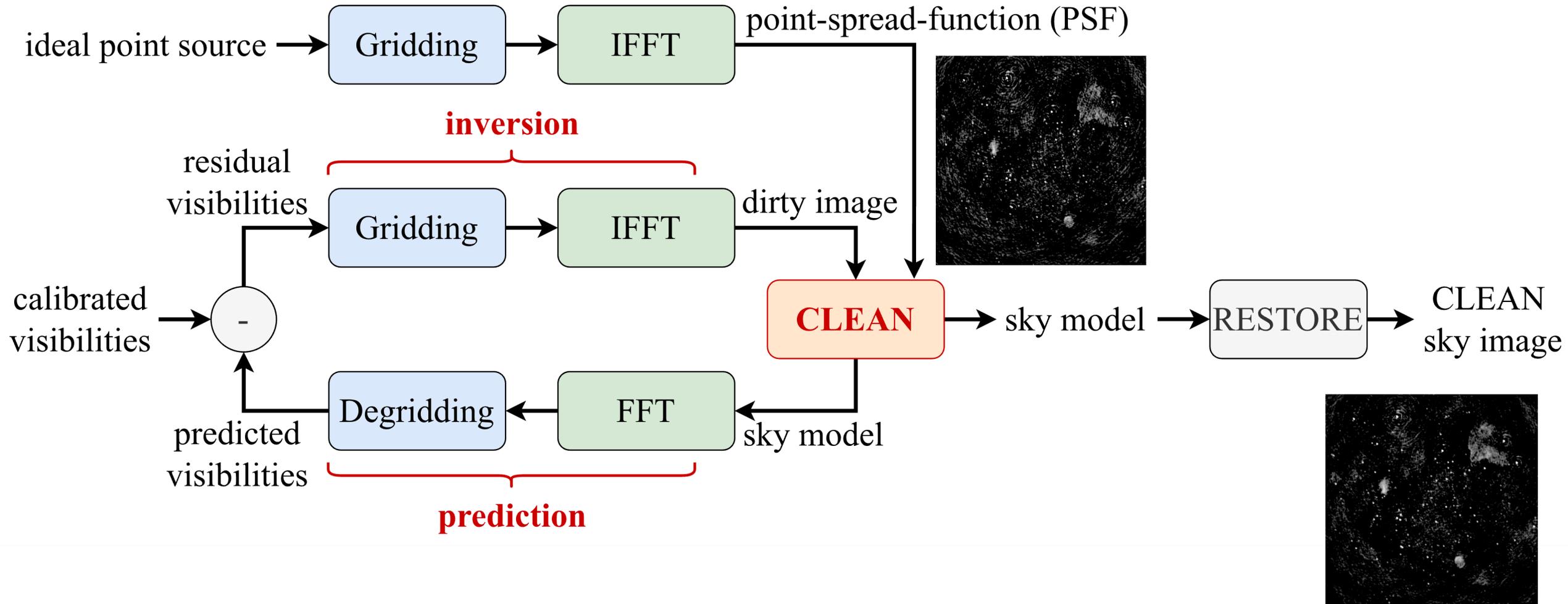
Contributions

- An in-depth analysis to determine the precision requirements for Image-Domain Gridding, included in the state-of-the-art imager WSClean.
- The first custom floating-point Gridding accelerator on reconfigurable hardware.
- An in-depth performance evaluation of our accelerator prototypes and state-of-the-art architectures with similar features (peak performance, thermal design power (TDP) and lithography technology).

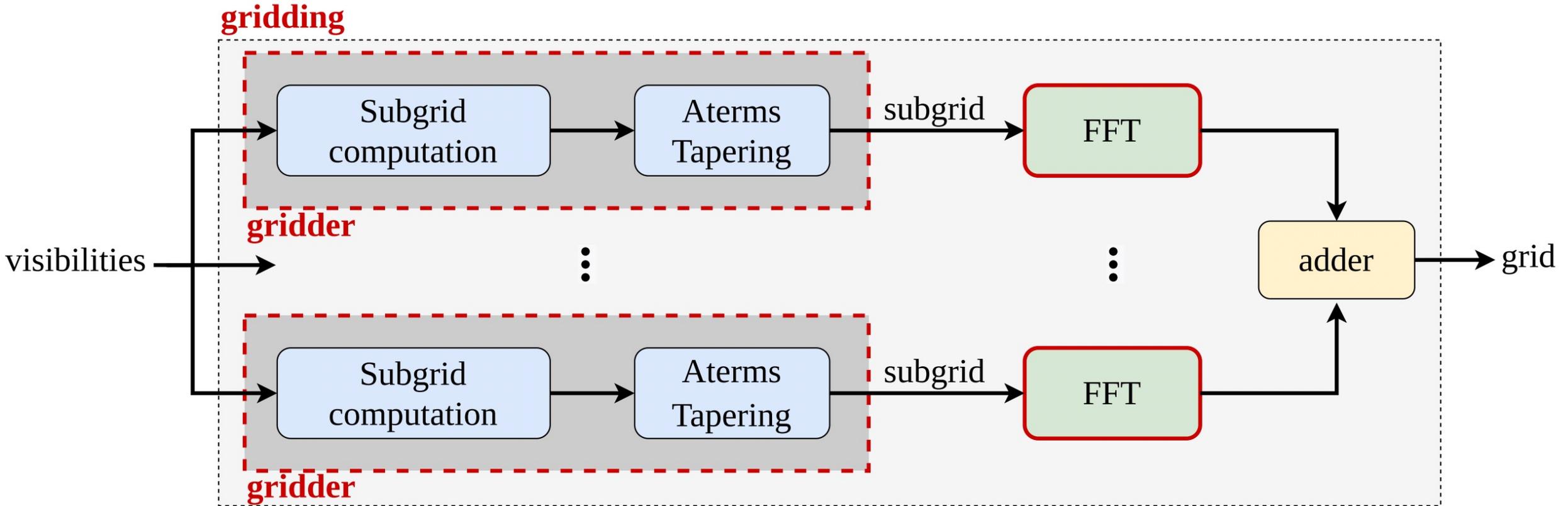
Interferometry



Radio-astronomical imaging



Gridding algorithm



S Van der Tol et al., "Image Domain Gridding: a fast method for convolutional resampling of visibilities", A&A 2018

B. Veenboer et al., "Radio-astronomical imaging on graphics processor", ASCOM 2020

B. Veenboer et al., "Radio-Astronomical Imaging: FPGAs vs GPUs", EuroPar19

A. R. Offringa et al., "Precision requirements for interferometric gridding in the analysis of a 21 cm power spectrum", A&A 2019

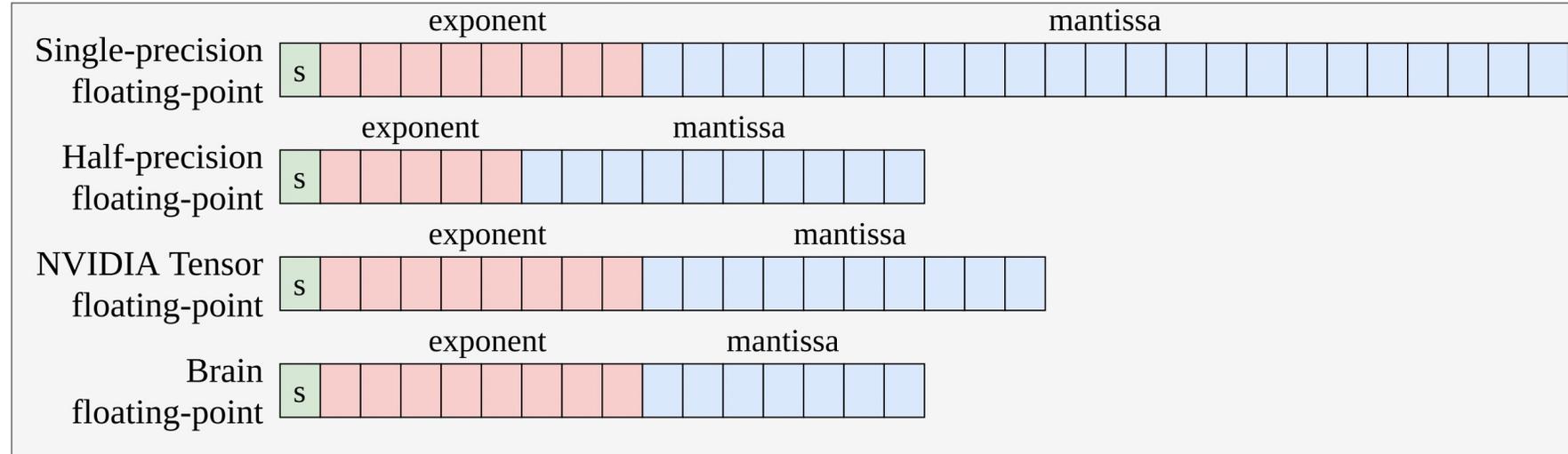
Reduced-precision

Software and hardware technique consisting in employing smaller data types to improve performance.

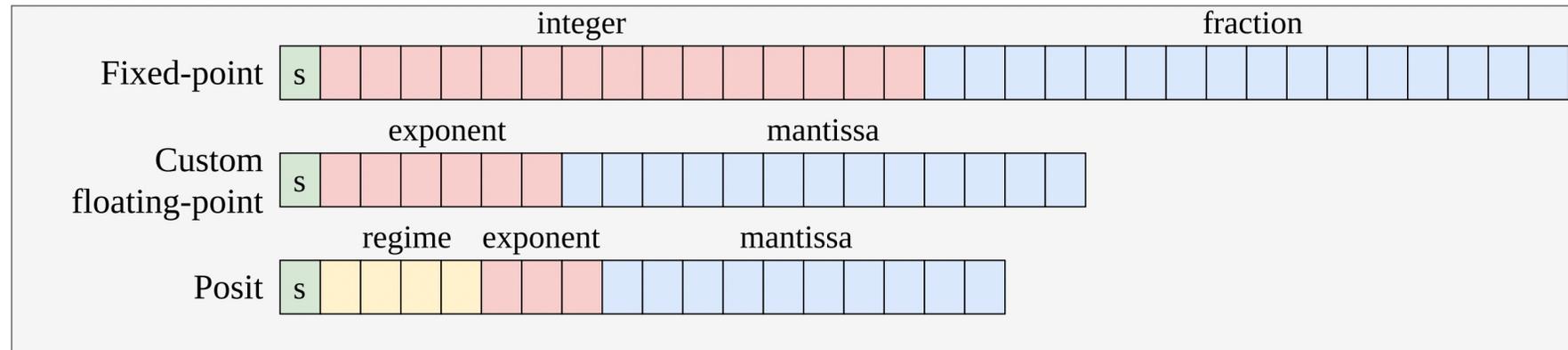
State-of-the-art:

- Automated/assisted precision tuning tools.
- Reduced-precision emulation libraries.
- Reduced-precision HLS libraries.

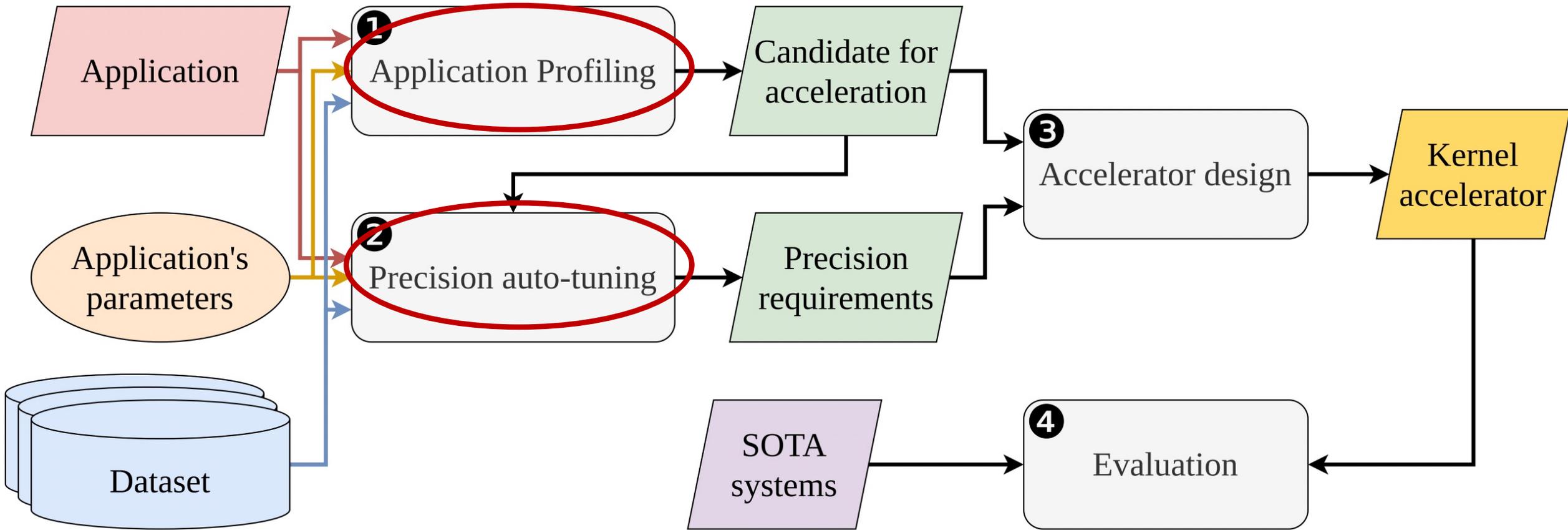
Standard arithmetic number formats



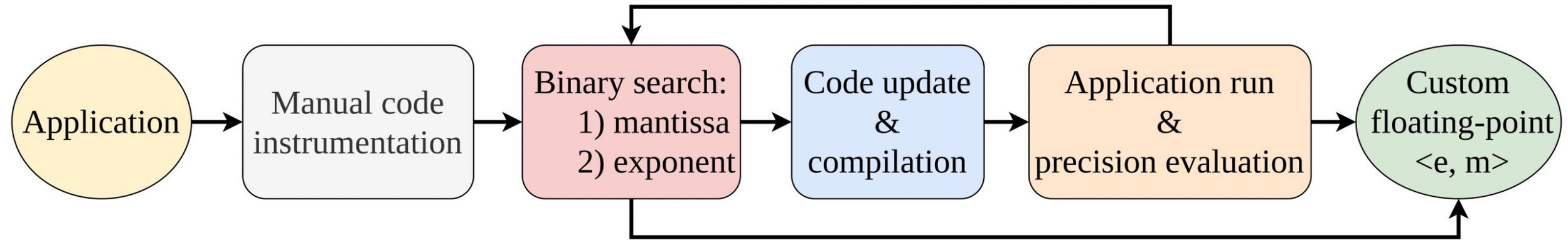
Custom arithmetic number formats



Methodology



Application profiling & precision auto-tuning



Structural Similarity Index Measure (SSIM) to assess image quality

Name	Description
Central Frequency	120.1172-125.7812 MHz
Channels per subband	4
Channel width	48 828.125 Hz
Declination	50.9410-54.8590
Duration	7199 s
Integration interval	2.002 78 s
Right Ascension	311.2500-318.7500

Datasets (LOFARSCHOOL) parameters

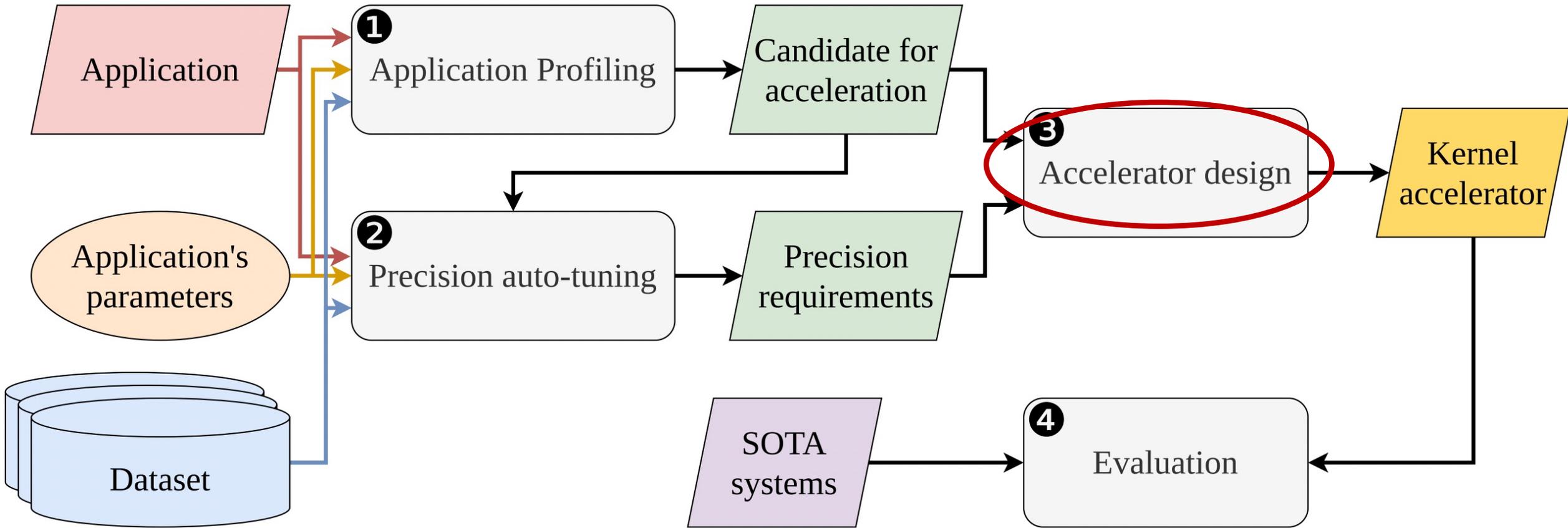
Software	Version
boost	1.68
OpenBLAS	3.9
python	3.8
wcslib	6.3
cfitsio	3.450
casacore [62]	3.3.0
dysco [63]	1.2
IDG [64]	master 011dfb18
WSClean [65]	master 2680c6a

SW versions

Parameter	Value	Description (unit)
size	6000 6000	output x and y dimensions (pixels)
scale	5 asec	scale of a pixel (degrees)
use-idg	active	-
auto-threshold	3	CLEAN stop condition (sigma)
niter	50000	number of minor CLEAN iterations
mgain	0.85	gain per major CLEAN iteration
weight	briggs 0	weighting mode and robustness
taper	gaussian 2amin	

Imager parameters

Methodology



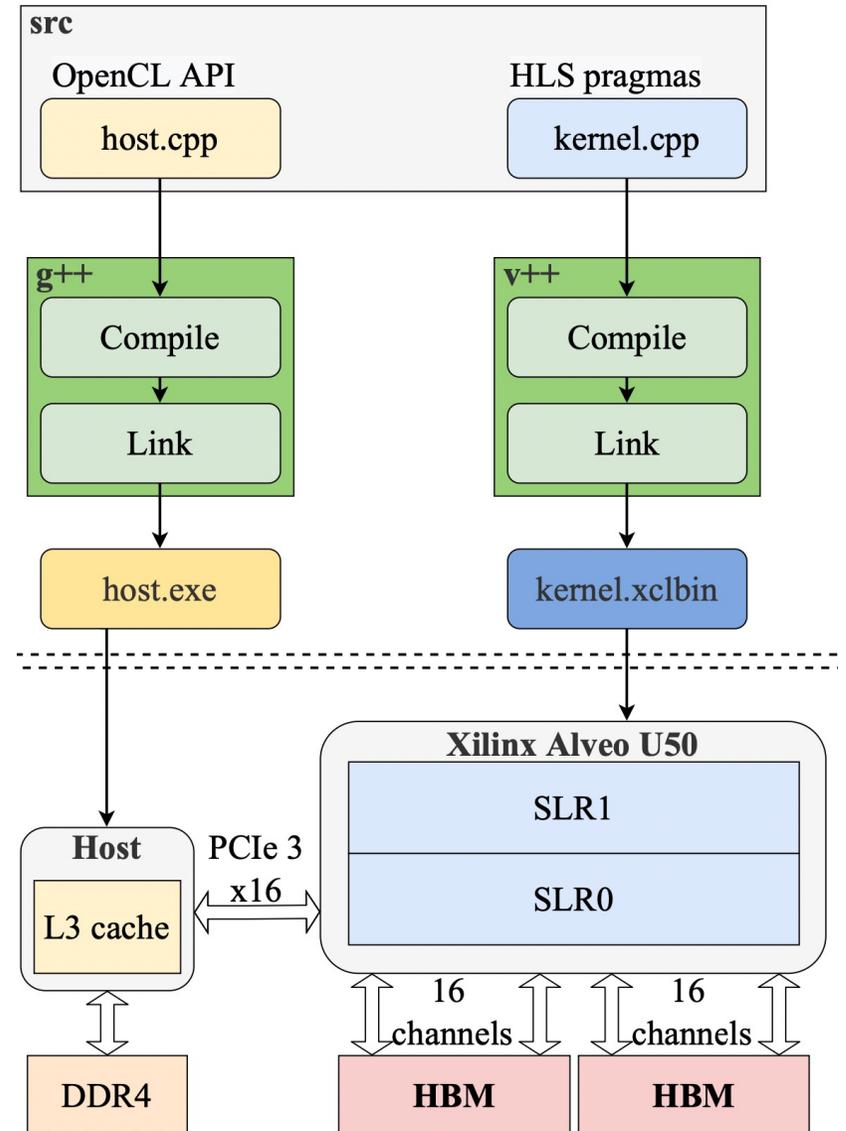
Accelerator design

Xilinx Vitis 2020.2:

- Host code with OpenCL API.
- Accelerator (Kernel) code with HLS pragmas.

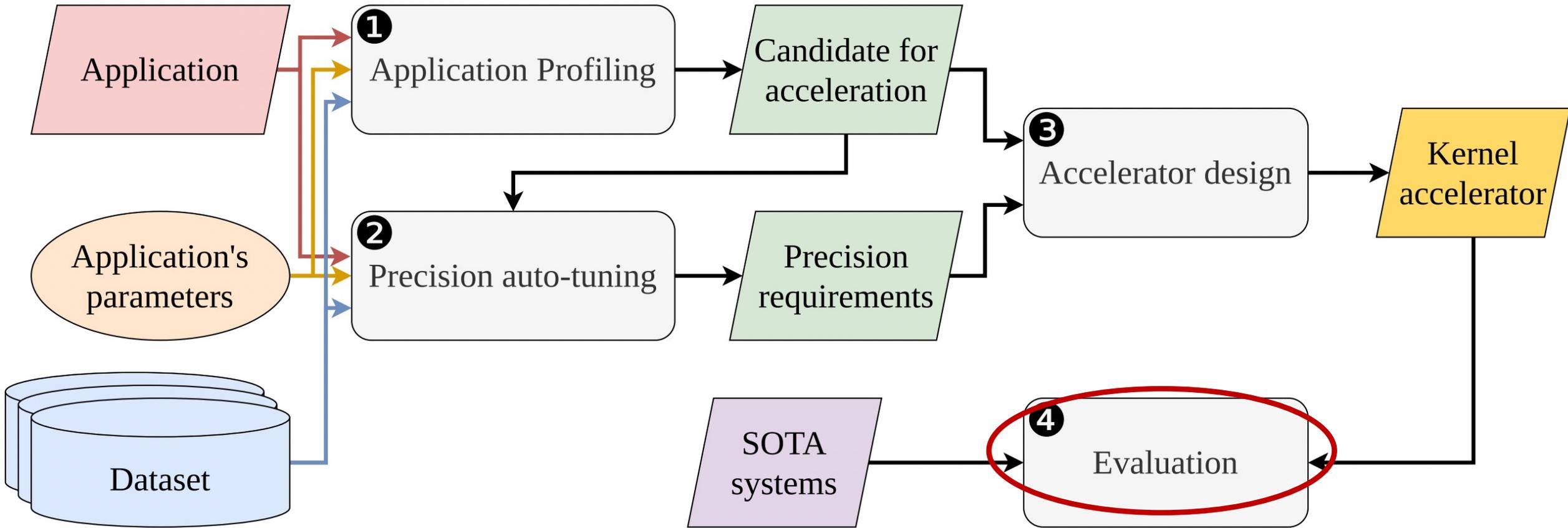
THLS (templatised soft floating-point for high-level synthesis) to map custom floating-point on FPGA.

Xilinx Vitis tool flow



Deployment system

Methodology



Evaluation

Intel i9 9900k	8 cores, 2 threads per core, 4.0 GHz all cores, 16 MB L3 Cache, 64 GB DDR4 3600 MHz
NVIDIA GTX 750	512 CUDA cores, 1.14 GHz, 2 MB L2 Cache, 2 GB GDDR5
AMD RX 550	8 compute units, 1.09 GHz, 512 KB L2 Cache, 4 GB GDDR5
Xilinx Alveo U50	872 K LUTs, 1743 K Registers, 5952 DSPs, 8 GB HBM2

Out-of-the-box TDP is lower than advertised

Performance evaluation:

- Libpowersensor: GPUs and FPGAs power
- Perf: CPU power, FLOP and DRAM traffic
- NVIDIA nvprof: NVIDIA GPU FLOP and DRAM traffic
- AMD CodeXL: AMD GPU FLOP and DRAM traffic

Architecture	Peak Performance	Bandwidth	TDP	Energy efficiency	Process
Intel i9 9900k	1.024 TFLOP/s	57.60 GB/s	95 W	10.79 GFLOP/W	14 nm Intel
NVIDIA GTX 1050 Ti	2.138 TFLOP/s	112.1 GB/s	75 W	28.50 GFLOP/W	14 nm Samsung [80]
NVIDIA GTX 750	1.164 TFLOP/s	80.19 GB/s	38 W	30.63 GFLOP/W	28 nm TSMC [81]
AMD RX 550	1.097 TFLOP/s	96.00 GB/s	35 W	31.34 GFLOP/W	14 nm GlobalFoundries [82]
Xilinx Alveo U50	Peak Performance	Bandwidth	TDP	Energy efficiency	Process
Theoretical (724 MHz)	1.547 TFLOP/s	316 GB/s	75 W	19.77 GFLOP/W	16 nm TSMC
Theoretical (300 MHz)	0.641 TFLOP/s	316 GB/s	75 W	8.55 GFLOP/W	16 nm TSMC
Empirical (292 MHz)	0.535 TFLOP/s	316 GB/s	75 W	6.84 GFLOP/W	16 nm TSMC

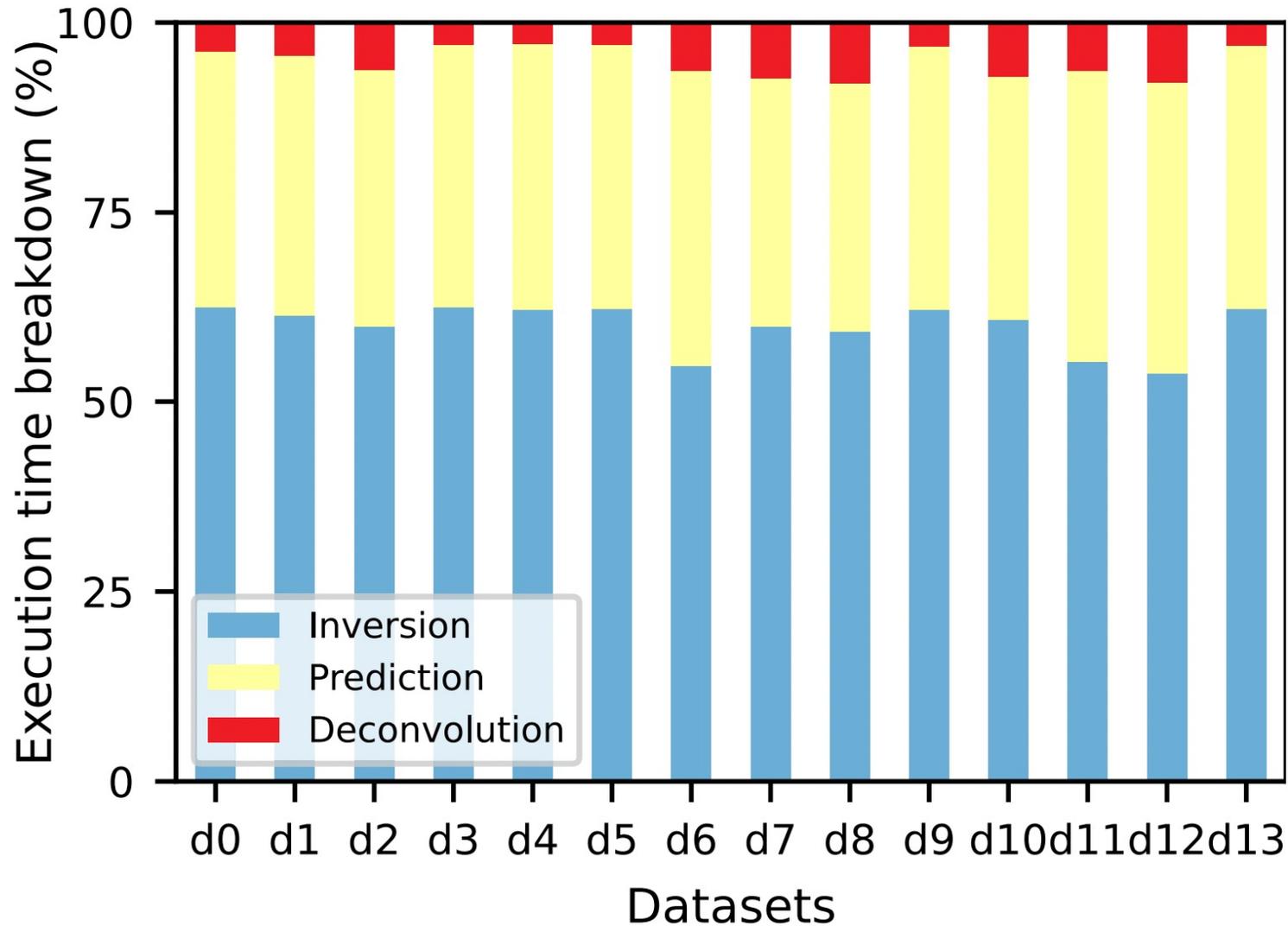
Similar peak performance

GPUs should be more energy-efficient

Similar lithography technology

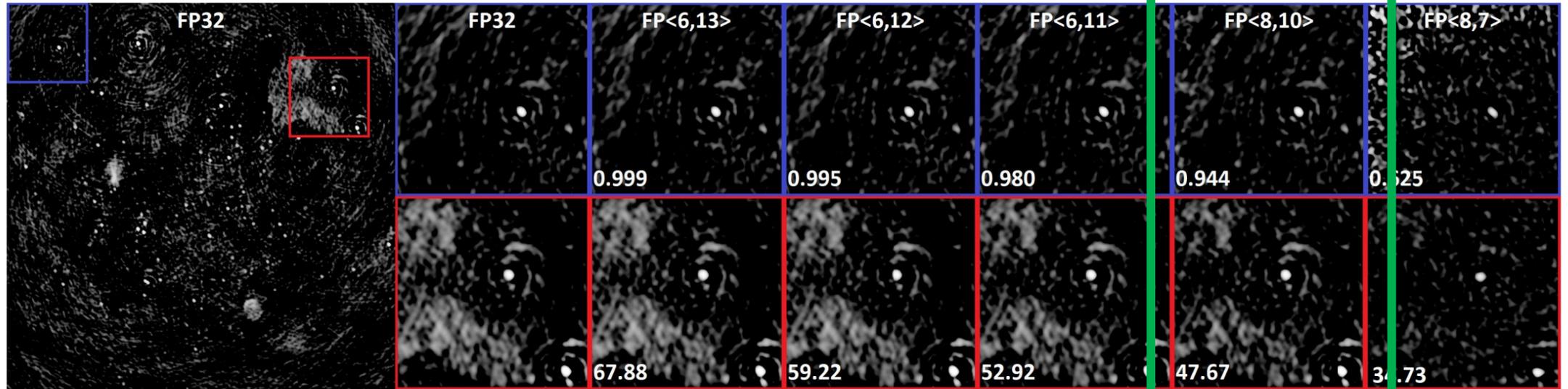
NVIDIA GTX 1050/1050ti has defective power measurement counters

Bottleneck analysis



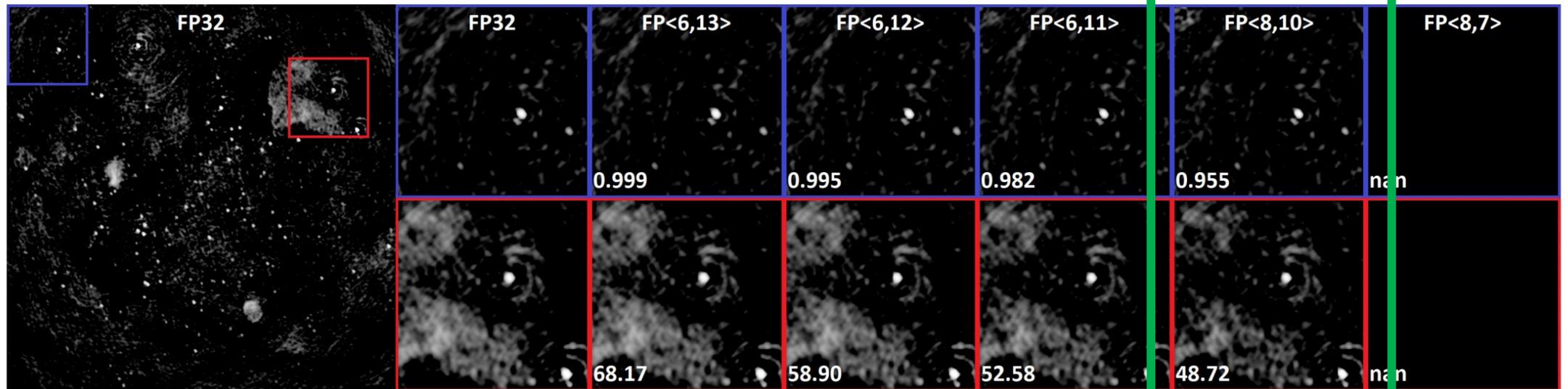
Inversion is the main bottleneck and Gridding is the largest kernel

Precision auto-tuning evaluation



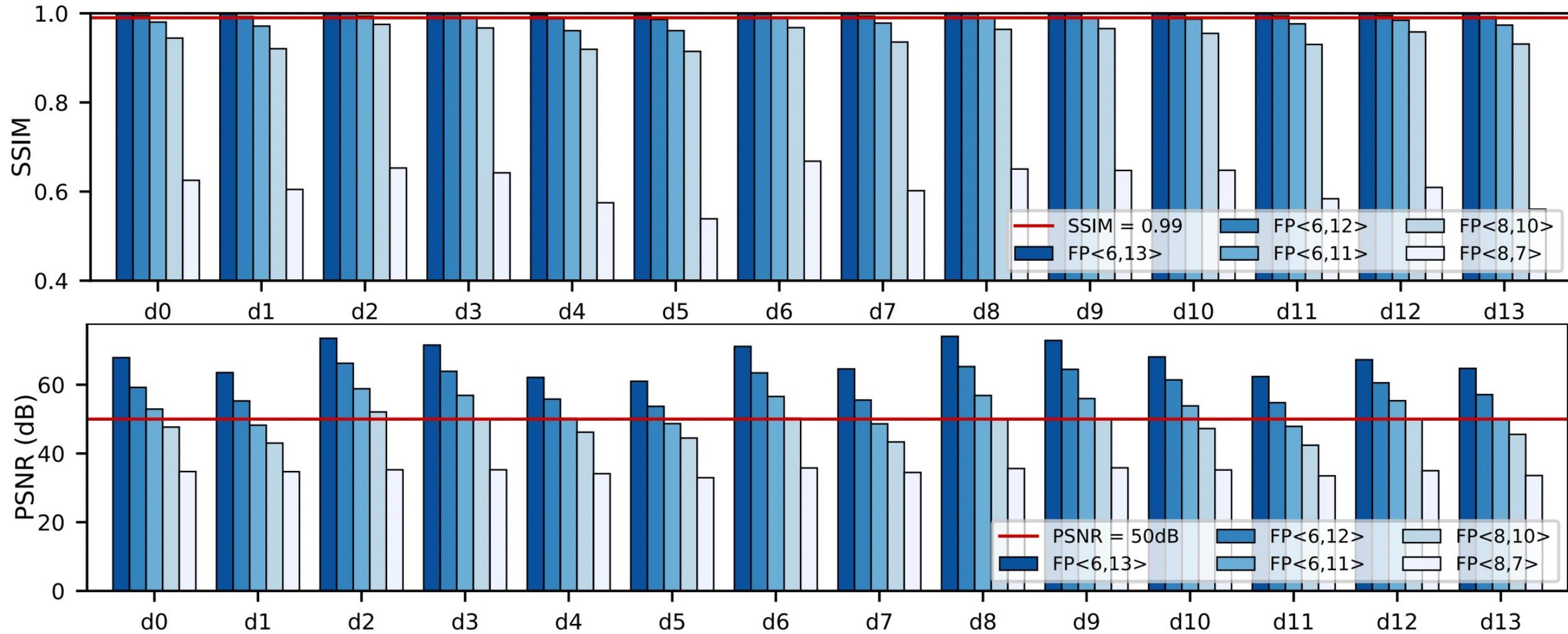
(a) Dirty images.

Artefacts

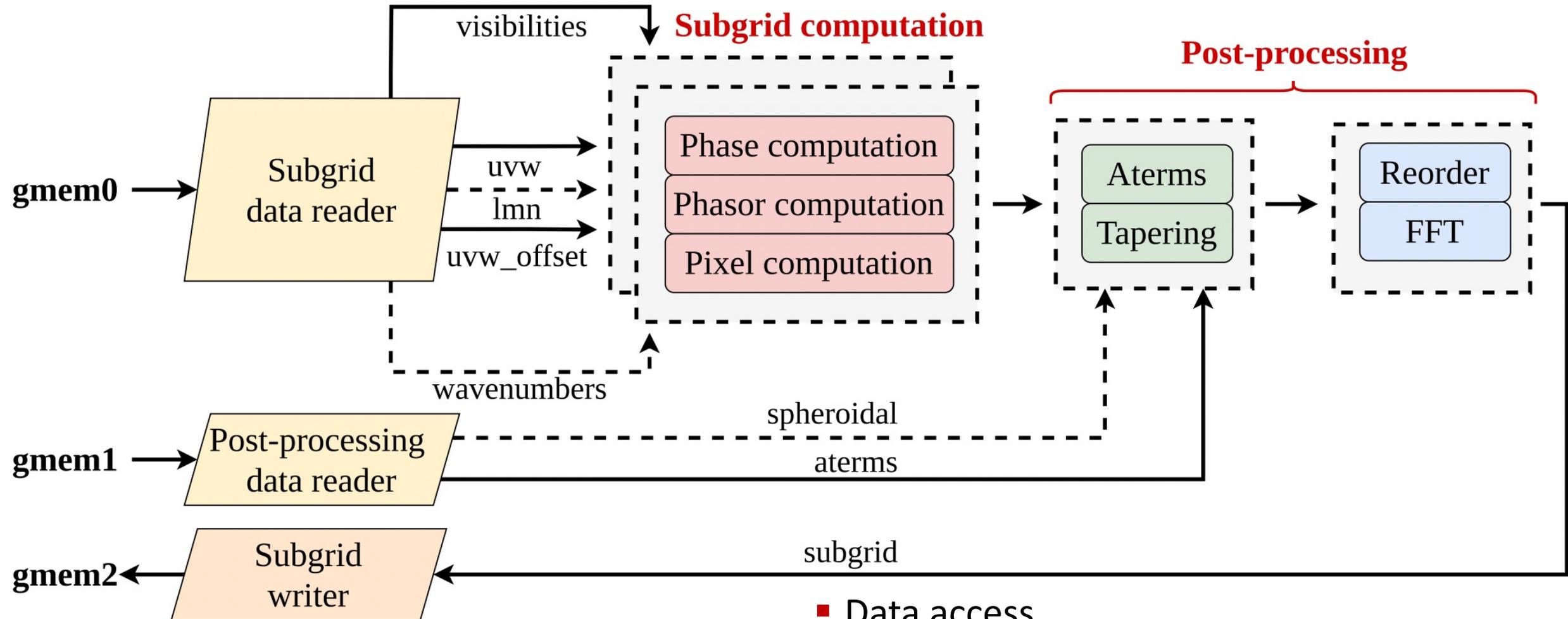


(b) Clean images.

Precision auto-tuning evaluation



Accelerator architecture



- Data access
- Initiation interval
- Parallelism

Accelerator architecture

Algorithm 1 Subgrid Computation HLS Pseudocode

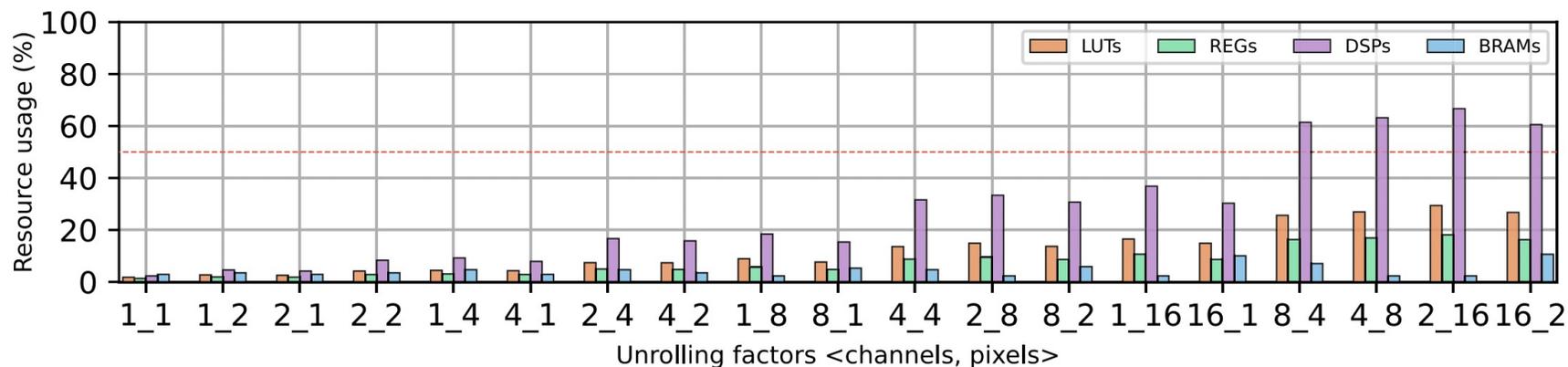
Input: visibilities, wavenumbers, uvw, uvw_offset, lmn

Result: subgrids

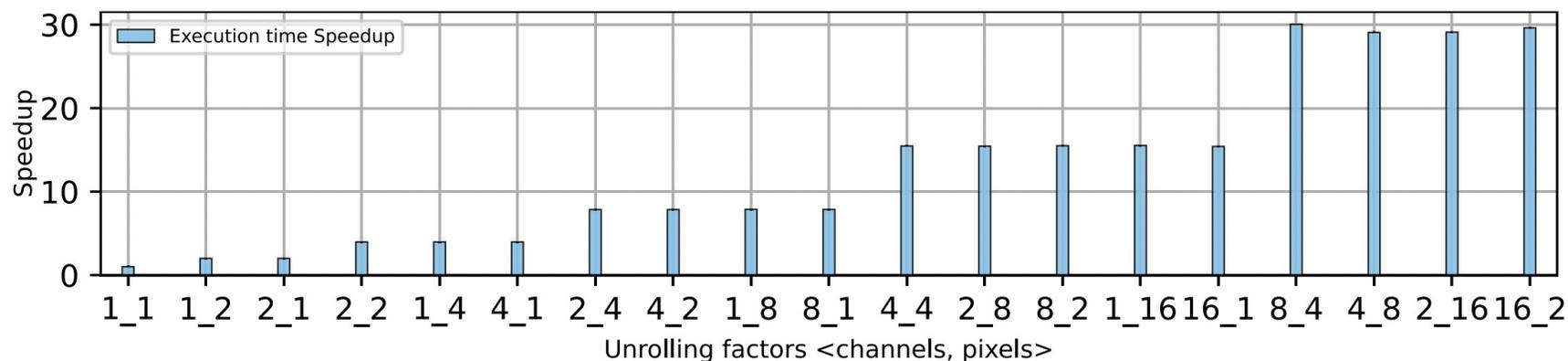
```
1 subgrids ← 0;
2 for s in subgrids_per_cu do
3   for t in timesteps do
4     for c in channels do
5       #pragma unroll factor = UNROLL_CHANNELS
6       for p in pixels do
7         #pragma unroll factor = UNROLL_PIXELS
8         complex<float> pixel[pol]
9         float lmn [3] ← lmn[p]
10        float phase_offset ← compute_phase_offset(uvw_offsets, lmn)
11        float phase_index ← compute_phase_index(uvw, lmn)
12        float phase ← compute_phase(phase_index, phase_offset, wavenumbers)
13        float phasor [2] ← cosisin(phase)
14        for pol in polarizations do
15          #pragma unroll
16          complex<float> pixel[pol] += visibilities[t][c][pol] * phasor
17        end
18      end
19    end
20  end
21 end
```

- Data access
- Initiation interval
- Parallelism

Subgrid computation: unrolling factors



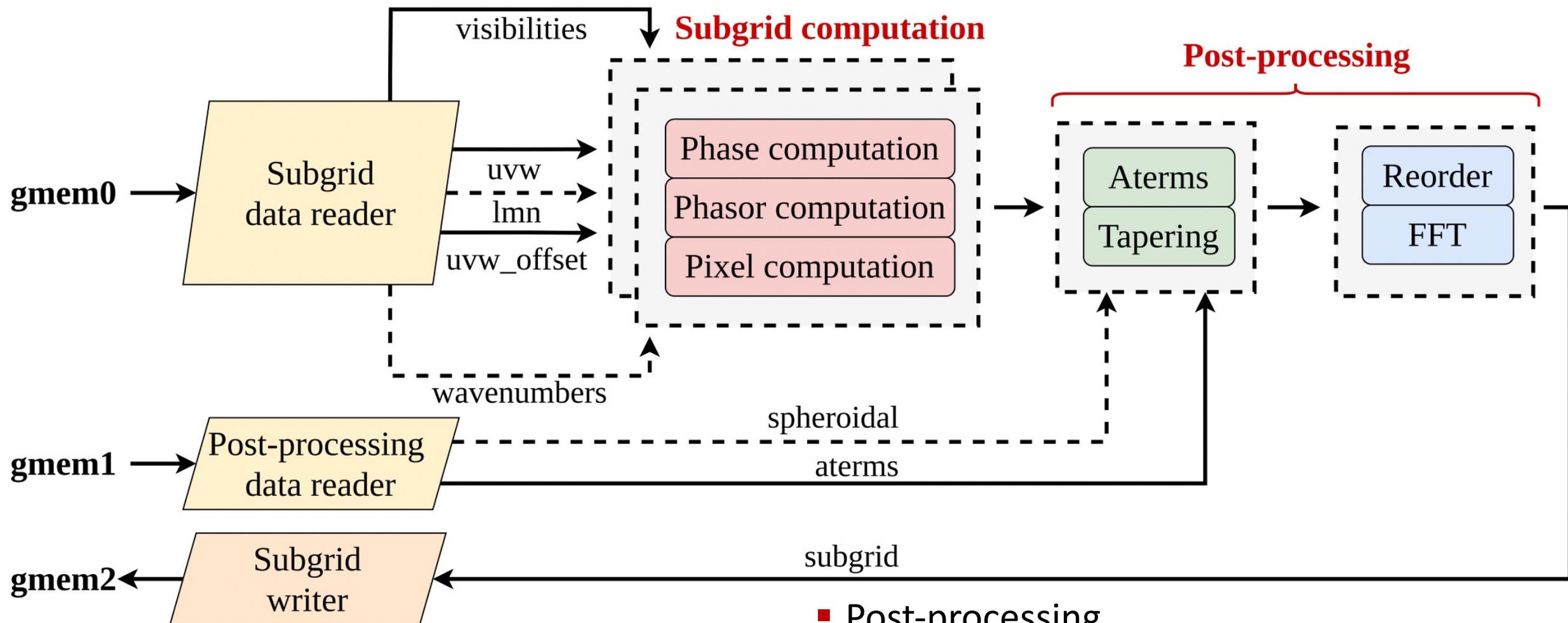
(a) Resource usage.



(b) Speedup.

- Similar performance with different unrolling factor combinations
- Unroll over channels → more BRAMs
- Unroll over pixels → more DSPs (more cosine/sine computations)
- 4_4 is a good trade-off for a balanced use of BRAMs and DSPs

Accelerator architecture



- Post-processing
- Cosine/sine
- Reduced precision

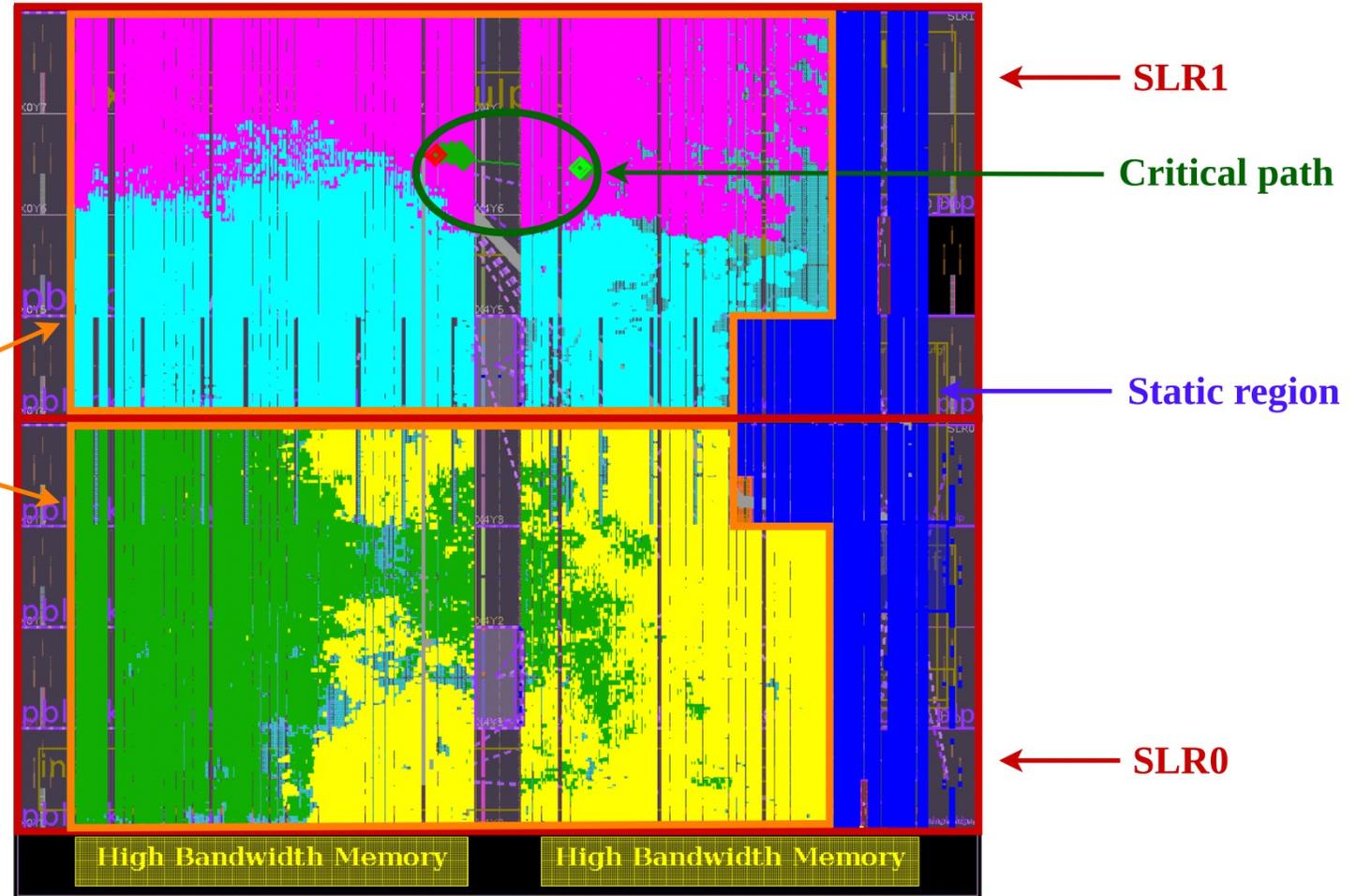
Device specific considerations

Resources type	Total	Dynamic region	Available (%)
LUTs	872 K	731 K	83.83%
REGs	1743 K	1462 K	83.88%
DSPs	5952	5340	89.72%
BRAMs	1344	1128	83.93%
URAMs	640	608	95.00%

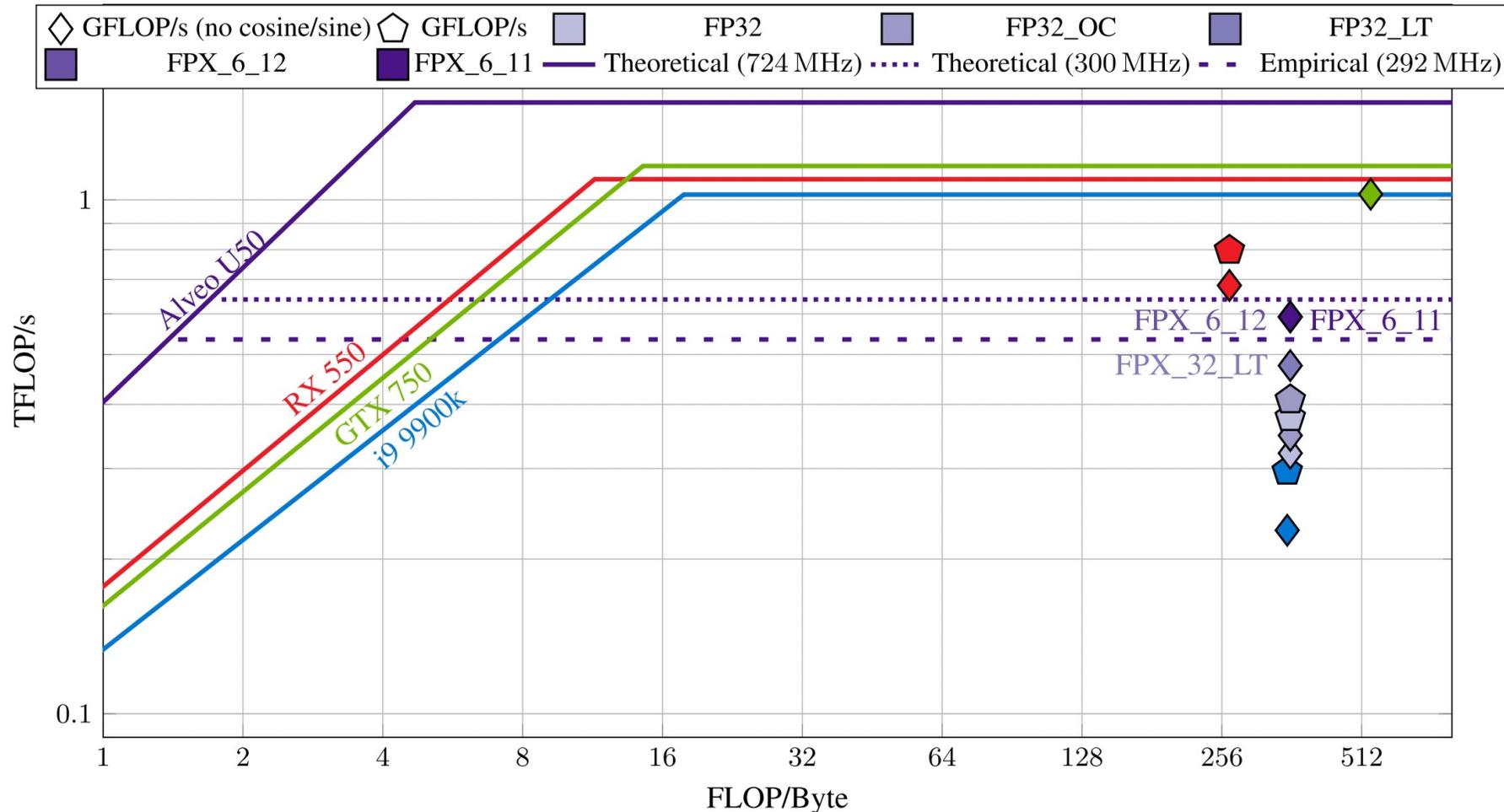
Mapping

- Static region
- Super Logic Regions (SLRs) and intra-gap
- Xilinx Vitis/Vivado strategy and pblock placement
- HBM channels
- Frequency

Dynamic Region



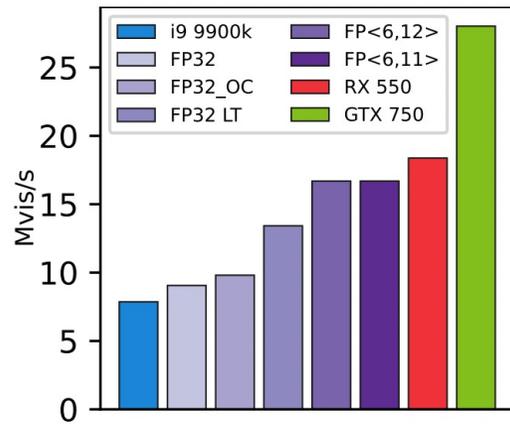
Performance evaluation



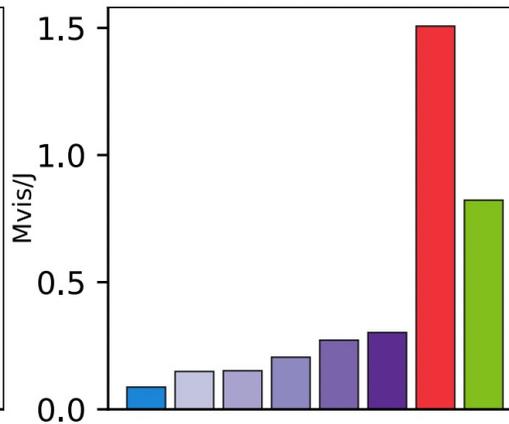
- NVIDIA has special units (SFU) for sine/cosine.
- AMD sine/cosine $\frac{1}{4}$ time compared to e.g. multiplication.
- Throughput vs HW utilization.

Area, throughput and energy efficiency

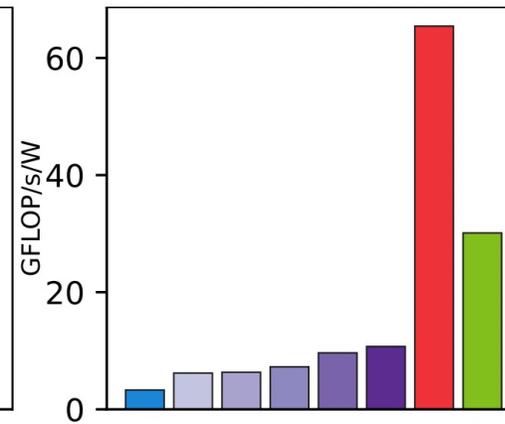
Version	LUTs	FFs	DSPs	BRAMs	Frequency
FP32	434 k (49.88%)	604 k (34.72%)	4114 (69.12%)	454 (33.74%)	300 MHz
FP32_OC	435 k (49.99%)	640 k (36.78%)	4114 (69.12%)	454 (33.74%)	346 MHz
FP32_LT	649 k (74.58%)	614 k (35.29%)	3142 (52.71%)	1045 (77.75%)	296 MHz
FPX_6_11	642 k (74.81%)	754 k (43.33%)	1956 (32.81%)	818 (60.86%)	300 MHz
FPX_6_12	656 k (75.39%)	767 k (44.10%)	1956 (32.81%)	818 (60.86%)	300 MHz



(a) Performance in terms of Mvis per second.



(b) Performance in terms of Mvis per Joule



(c) Energy efficiency in terms of GFLOP/s per Watt

- Power limit FP32 & timing closure issues.
- Lookup-table (+50% computation) & reduced precision (+100% computation).
- Xilinx Alveo U50 outperforms Intel i9 9900k (up to 2.12x in throughput and 3.46x in energy efficiency) and has a comparable throughput wrt. To AMD RX 550 (~11% slower).
- GPUs are more efficient (GTX 750 highest energy efficiency → scaling technology with DeepScaleTool).

Related work

Work	Date	Application	Platform	Optimization
Offringa [20, 35]	2014-2017	W-Stacking, CLEAN (Högbom, Cotton-Schwab, Multiscale)	CPU	Optimized full imager (WSClean)
Veenboer [6, 98]	2017-2020	Image-Domain Gridding	CPU/GPU	code optimization (added to WSClean)
Grel [99]	2018	Högbom CLEAN	FPGA	Custom accelerator of Högbom CLEAN formulated as a Compressive Sensing problem
Veenboer [11]	2019	Image-Domain Gridding	FPGA	custom accelerator
Seznec [15]	2019	Generic deconvolution	GPU	half-precision deconvolution
Hou [100]	2020	W-Projection	FPGA	custom accelerator
Cordeiro [4]	2020	Large 2D FFT	FPGA	NMC acceleration

What is new?

- Reduced-precision evaluation for radio-astronomical imaging.
- Porting of the Image-Domain Gridding Algorithm on a Xilinx FPGA.
- Reduced-precision Image-Domain Gridding algorithm.

Conclusions

Summary:

- An in-depth analysis to determine the precision requirements for Image-Domain Gridding (IDG), included in the state-of-the-art imager WSClean.
- The first custom floating-point Gridding accelerator on reconfigurable hardware.
- An in-depth performance evaluation of our accelerator prototypes and state-of-the-art architectures with similar features: peak performance, thermal design power (TDP), and lithography technology.

Lessons learned:

- Reduced precision suitability for radio-astronomical imaging.
- Benefits of reduced precision in the radio-astronomical imaging application domain.
- FPGAs vs CPUs vs GPUs.
- Xilinx Alveo U50.

Current and Future work

- CLEAN study and acceleration on FPGA/GPU.
- AI HW, e.g. NVIDIA tensor cores for imaging pipeline.
- HW/SW co-design for SKA:
 - HPC benchmarking: SKA-SDP Benchmark Suite: <https://gitlab.com/ska-telescope/sdp/ska-sdp-benchmark-tests>
 - Porting Image-Domain Gridding with HIP: <https://gitlab.com/ska-telescope/sdp/ska-sdp-idg-bench>
 - Evaluation and SW optimization of upcoming HW (Intel Ponte Vecchio, NVIDIA Grace – Hopper, ARM CPUs, AMD CDNA2/3, FPGAs...)

Reduced-Precision Acceleration of Radio-Astronomical Imaging on Reconfigurable Hardware

Thanks!

Stefano Corda

stefano.corda@epfl.ch

07-07-2022

Scientific Computing Accelerated on FPGAs, Maison de la Simulation (Saclay)

