



Energy-Efficient FPGA Solutions for Large-Scale FFTs and Non-Uniform FFTs

A Software-Hardware Co-Design Approach for Radio Interferometry

**Rubén Rodríguez Álvarez, Denisa Constantinescu, Miguel Peón Quirós, Adrien
Devresse, Hamza Chouh, Shreyam Krishna, Etienne Orliac, David Atienza**

EPFL - Embedded Systems Laboratory, EcoCloud, SCITAS

ruben.rodriiguezalvarez@epfl.ch, denisa.constantinescu@epfl.ch

Scope

Energy-efficient computing with domain-specific accelerators
Multi-scale hardware-software co-design approach

Group Members

France

- INSA Rennes: Jean-F. Nezan, Mickaël Dardaillon, Hugo Miomandre, Jacques Morin
- OCA: Shan Mignot, Alain Miniussi, Chiara Ferrari, André Ferrari
- OP: Damien Gratadour

Switzerland

- EcoCloud: Miguel Peon Quiros, David Atienza
- ESL: Denisa Constantinescu, Rubén R. Álvarez, Basile Darne, David Atienza
- SCITAS: Adrien Devresse, Hamza Chouh, Etienne Orliac, Gilles Fourestey

Partners: EPFL, Laboratory of Astrophysics, MeerKAT, SKACH, SKAO

Pipelines Profiling & Specification

Single-Node SW/HW co-design

Multi-Node Scale-up

Integration and Testing

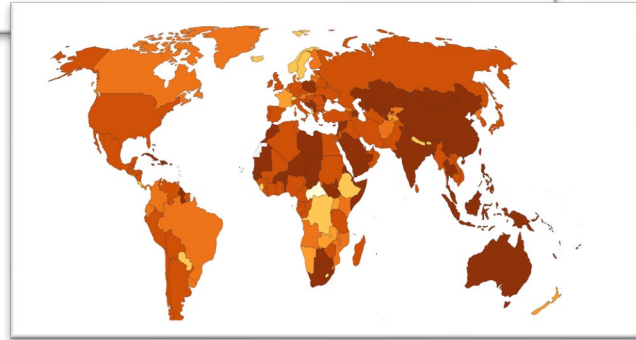
On-Field Demonstrator



seams-project.com

Electricity: Berlin's shock plan to adapt to the weather

Companies may soon have to adapt their production to the strength of the wind and the duration of sunshine, in order to relieve the electricity networks, put to the test by the intermittency of renewable energies. This is the option proposed by the Ministry of Economy and Climate in a note published in July. Enough to trigger the ire of the business world.



Global concern for energy consumption and lower carbon emission factor

Use domain-specific accelerators (GPUs, FPGAs, ASICs)

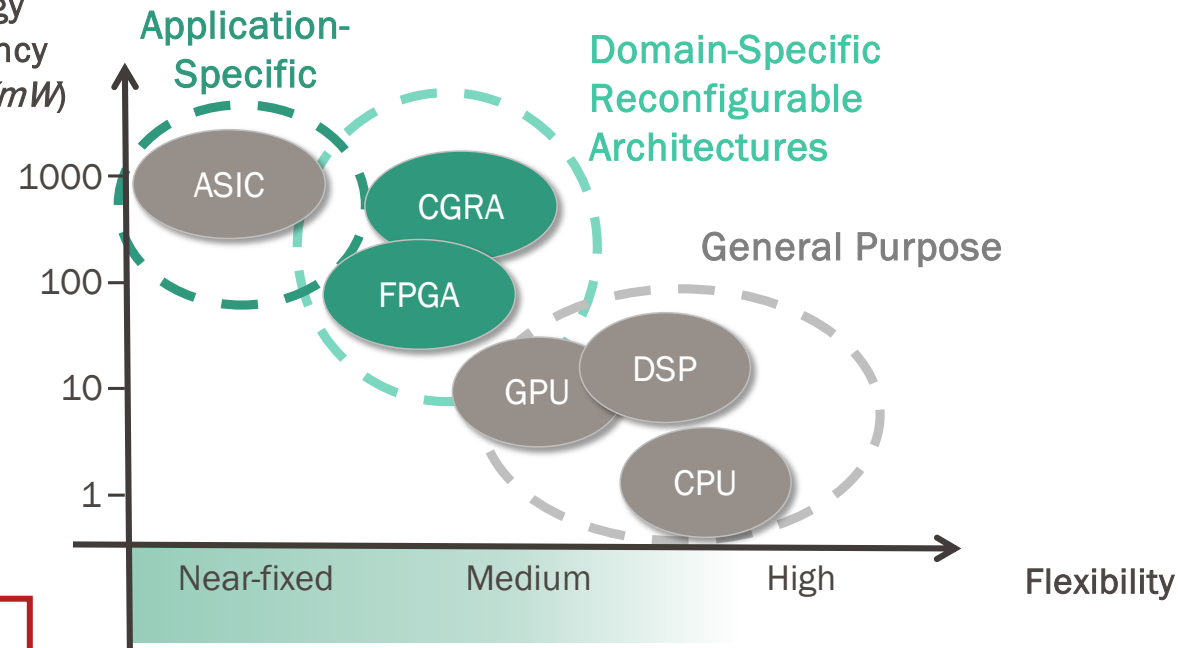


Shift the focus to Energy-to-Completion

SKA

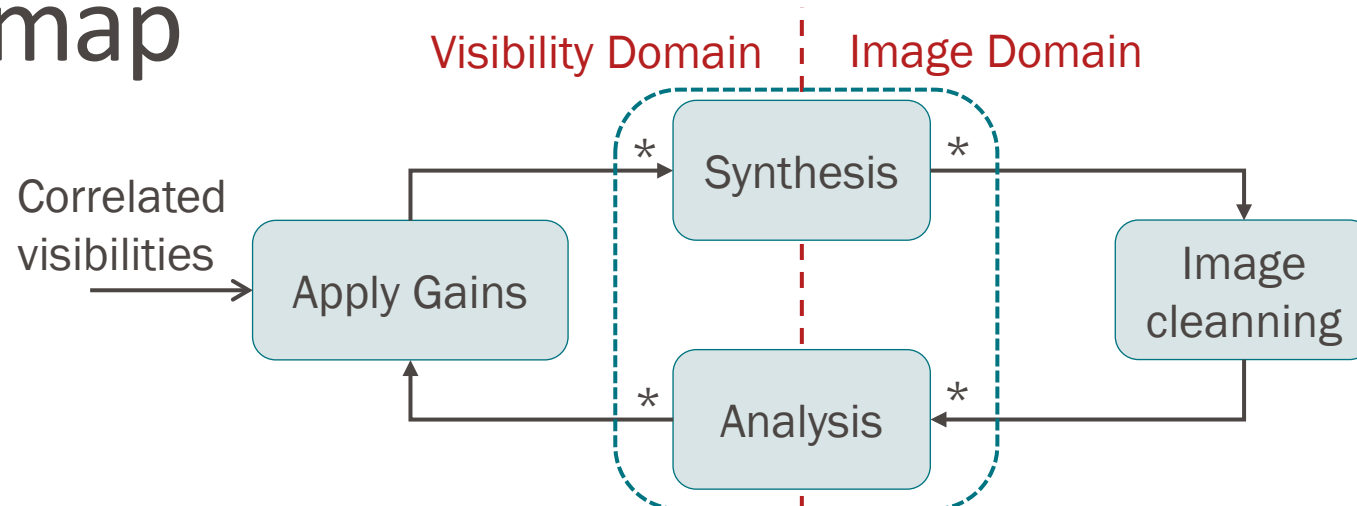
- Exa-scale amounts of data
- Large computation -> Scalability in Multi-node
- Multi-kernel -> Domain-specific computation

Energy Efficiency (MOPS/mW)



Adapted from Kevin J M Martin. IPDPSW 2022

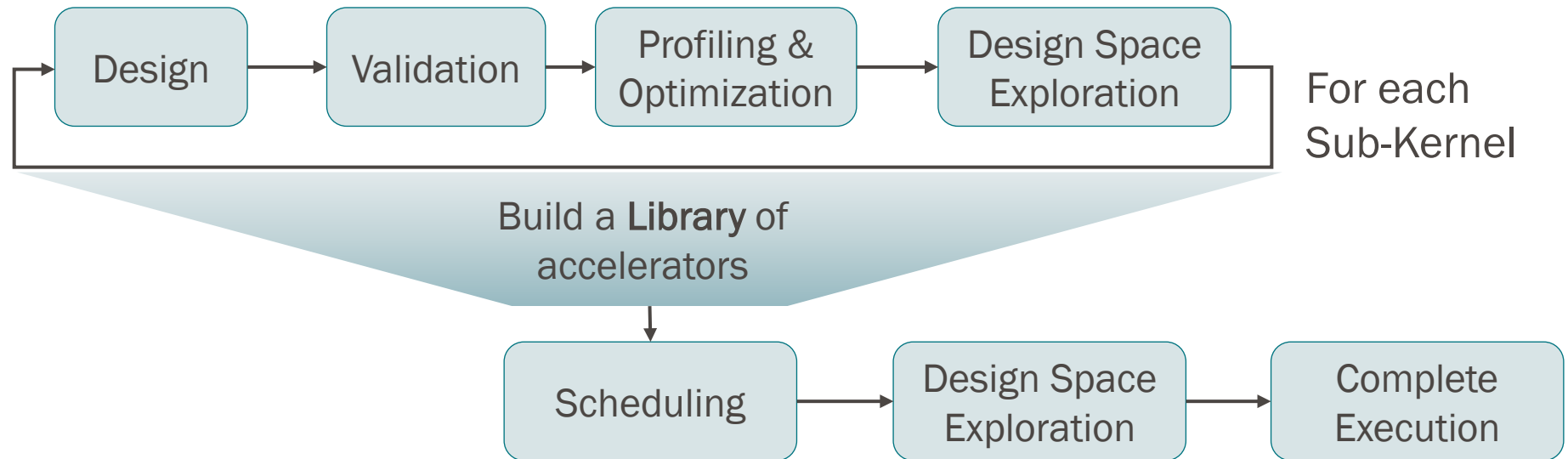
Roadmap



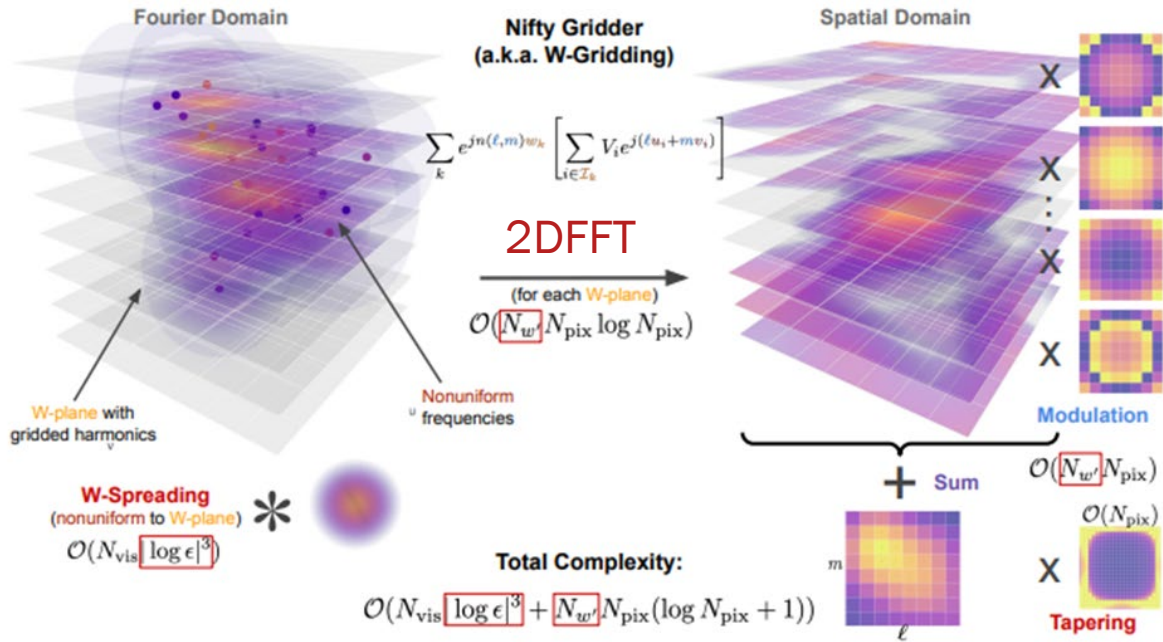
Divide the operation in **Sub-Kernels**

*Non-uniform points

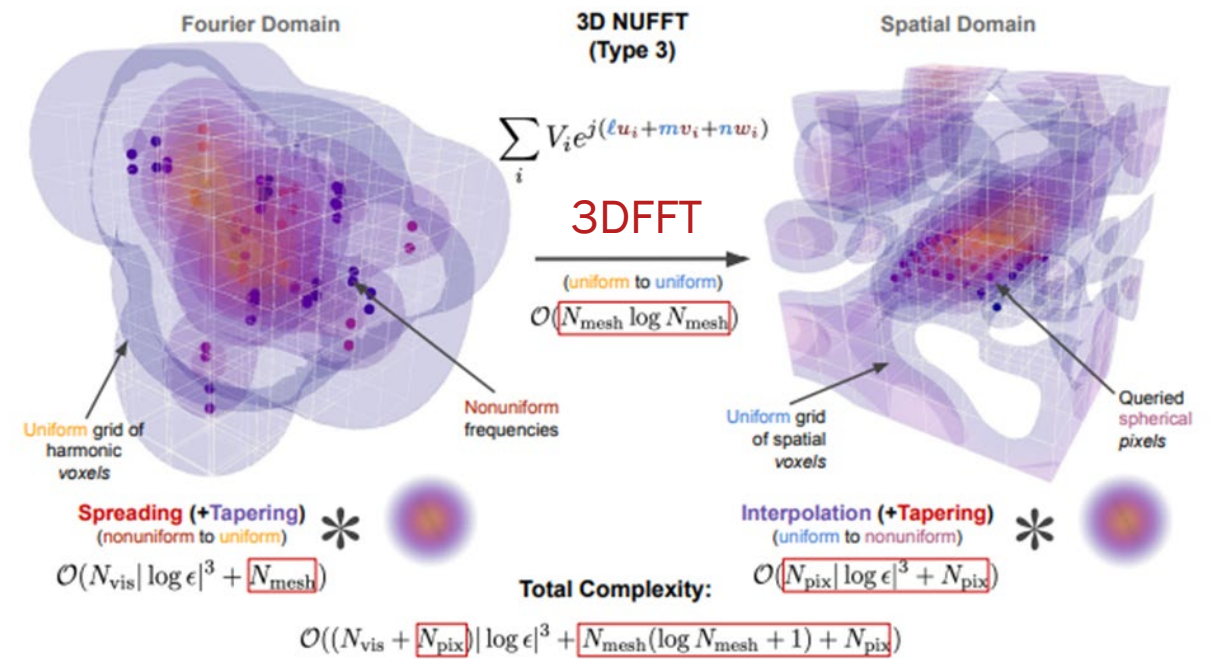
FPGA Flow:



W-gridder



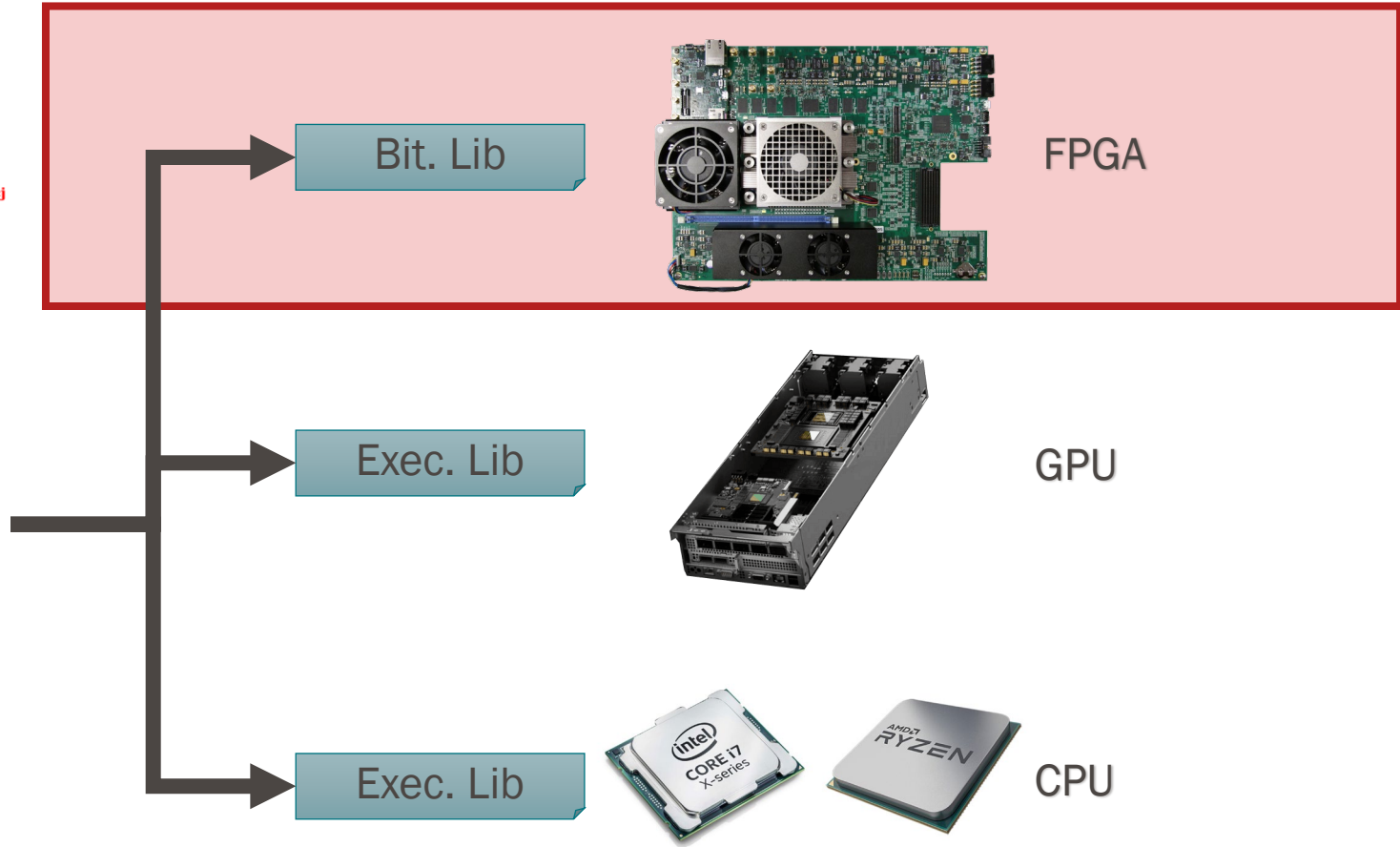
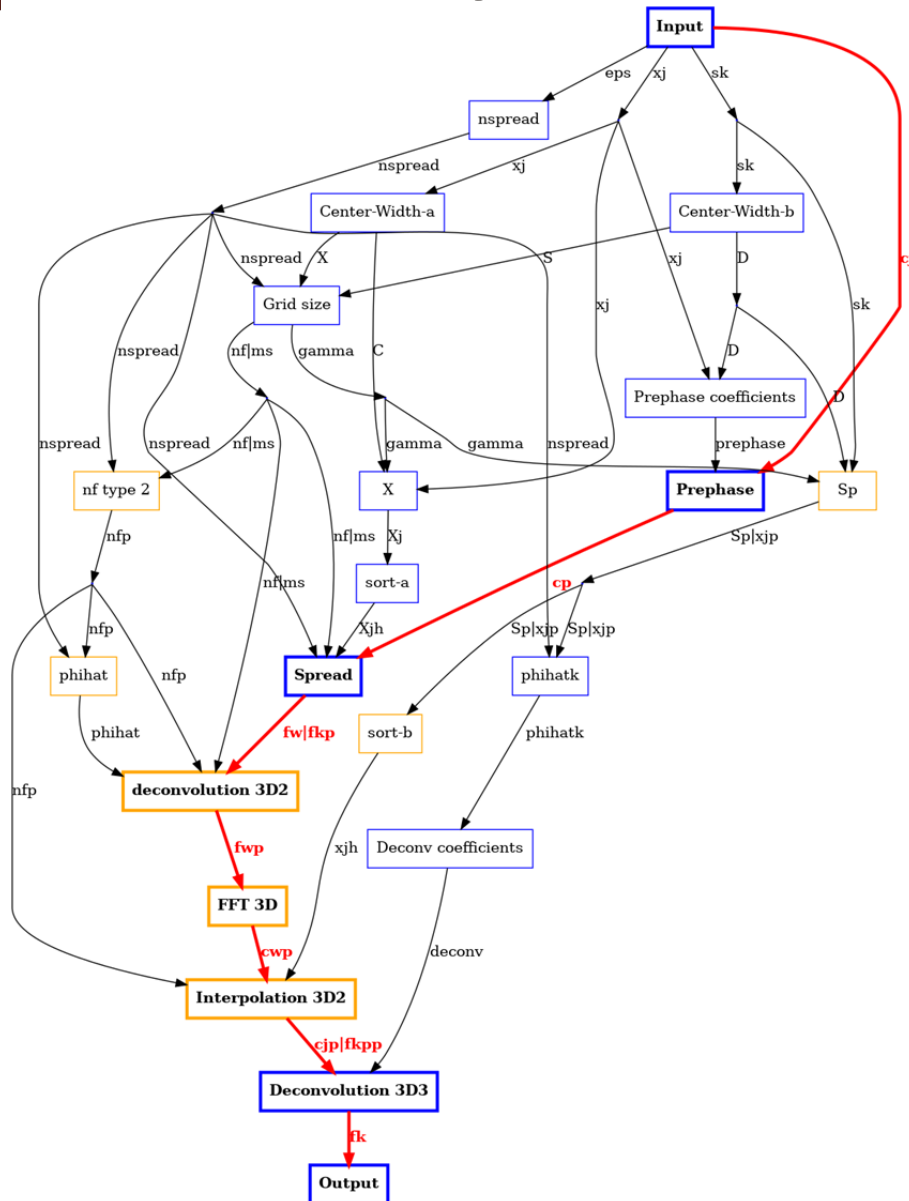
Finufft



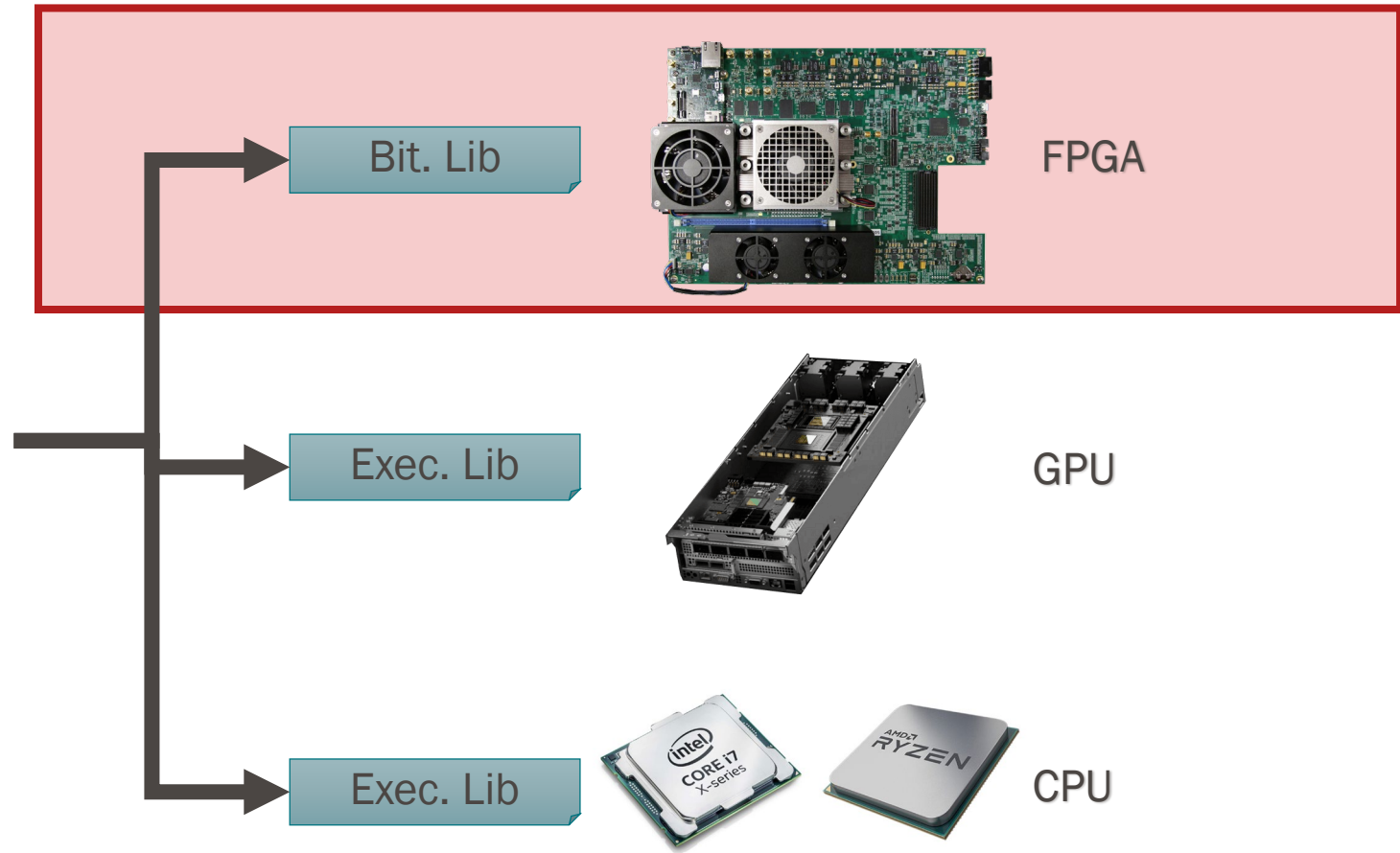
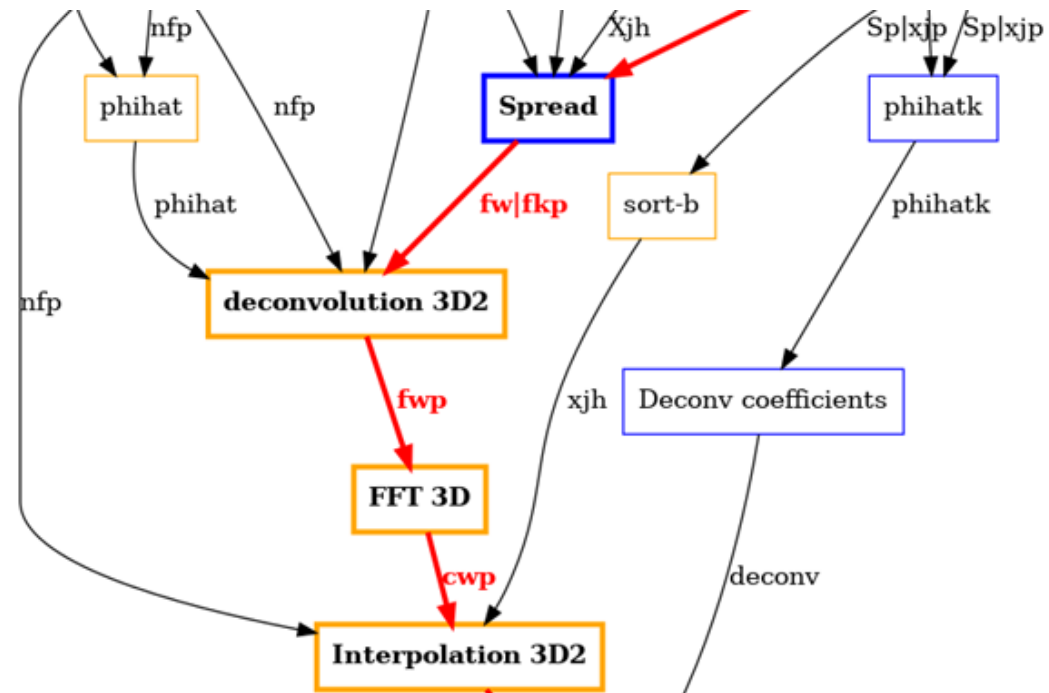
- Kashani et al. "HVOX: Scalable Interferometric Synthesis and Analysis of Spherical Sky Maps." (2023).
- Tolley et al. "BIPP: An efficient HPC implementation of the Bluebild algorithm for radio astronomy." (2023).
- Corda et al. "Near memory acceleration on high resolution radio astronomy imaging." MECO. IEEE, (2020).

How do these algorithms map into an FPGA?

Finufft Synthesis

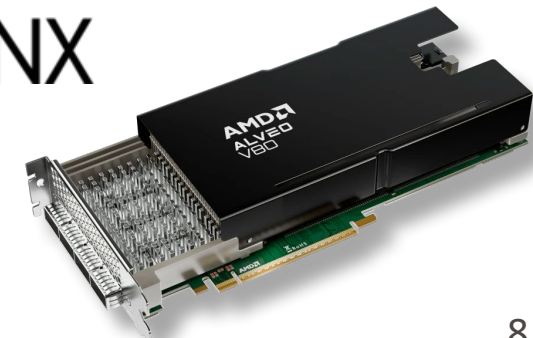


Finufft Synthesis

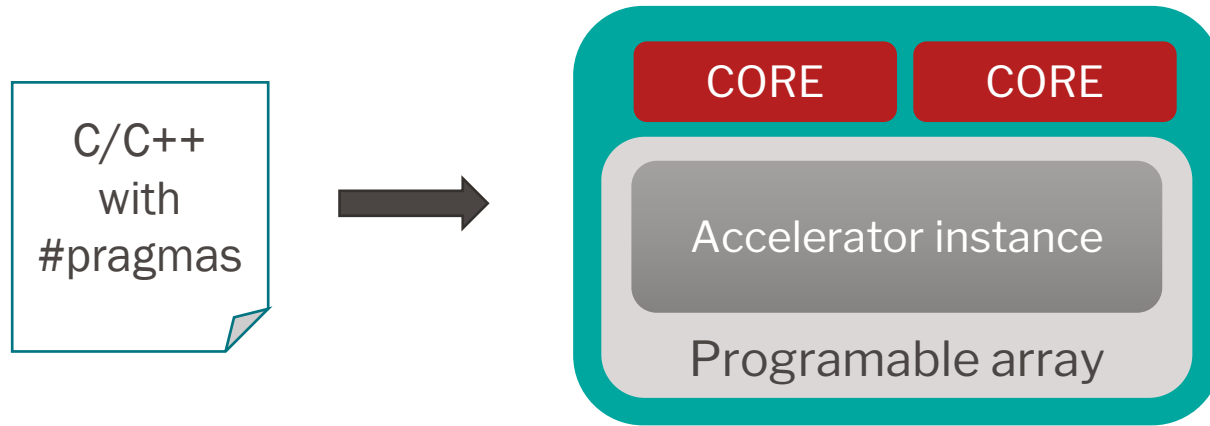


3D FFT takes 40%-90% of the computation

Characteristics	Agilex 7 M-Series Dev Kit	Alveo V80 Card
Internal memory	370Mb BRAM	132Mb BRAM + 541Mb URAM
High Bandwidth Memory (HBM2e)	32GB @ 1TB/s	32GB @ 810GB/s
Compute Elements	3.9M LEs + 12.3K DSPs + 1.3M ALMs	2.6M LUTs + 10.8K DSPs
Max Power (TDP)	(2x) 240 Watts	190 Watts
Global Memory (DDR4/5)	64 GB	32 GB
Comms	16x PCIe 5, CXL, GbE 116Gbps, fiber optic	2x PCIe 5
Technology	7nm Intel	7nm TSMC
Max Clock Freq	500MHz-1GHz	600MHz-1GHz



High-Level Synthesis (HLS) for FPGAs



Characteristics:

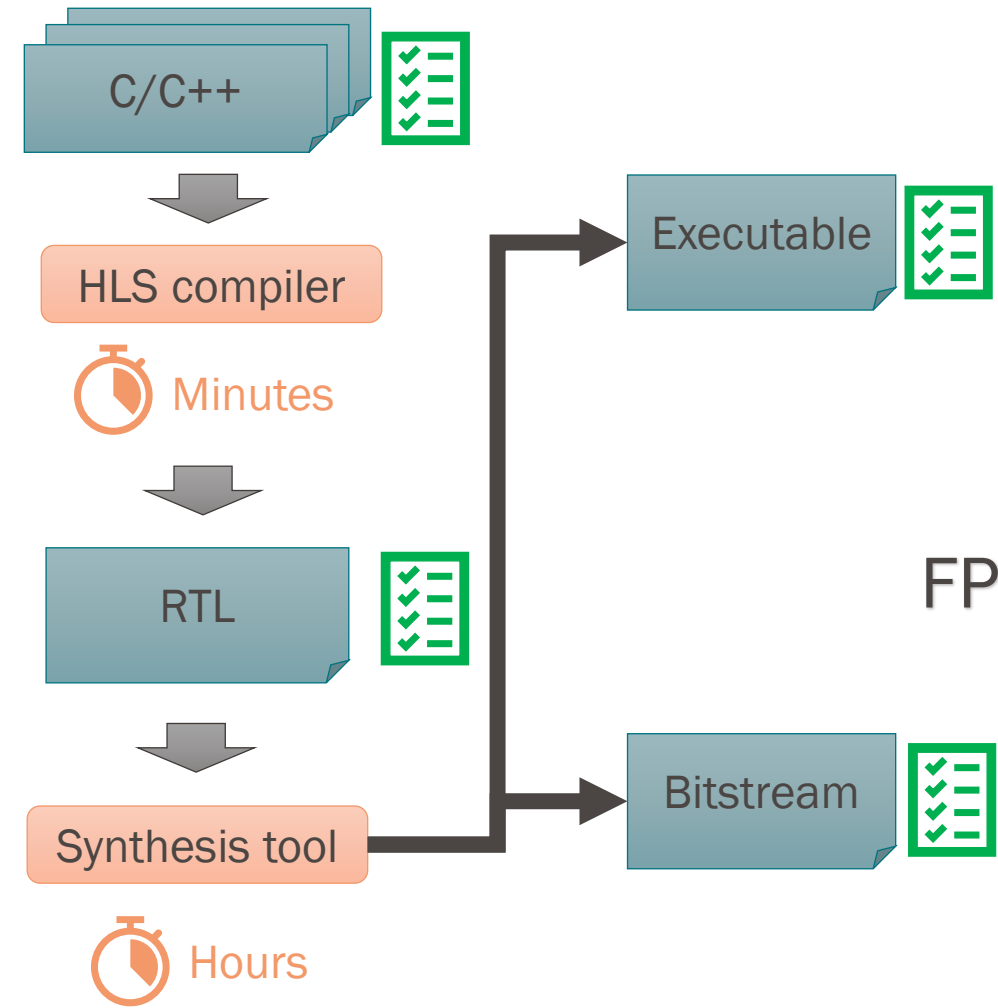
- Mixed precision data types
- **Parallel, pipeline and serial**
- Resources constraints
- Code breakdown
- Highly parametrizable

Characteristics	HLS FPGA	CUDA GPU
Programming support	High	High
Productivity (design time)	Medium	High
Energy Efficiency	High	Low-Medium
Latency	Medium	Low
Scalability	High	High
Flexibility	High	Limited

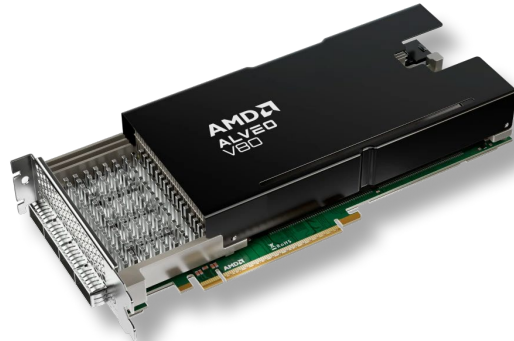
We teach HLS and Co-Design; used it to accelerate **CNNs** and **genome alignment** applications

HLS is a good fit for changing SW, portable HW, and design explorations

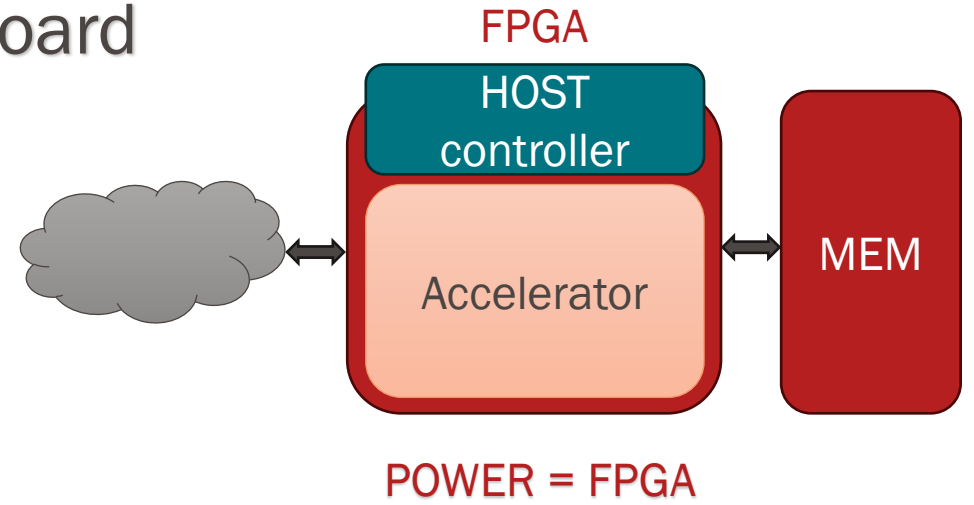
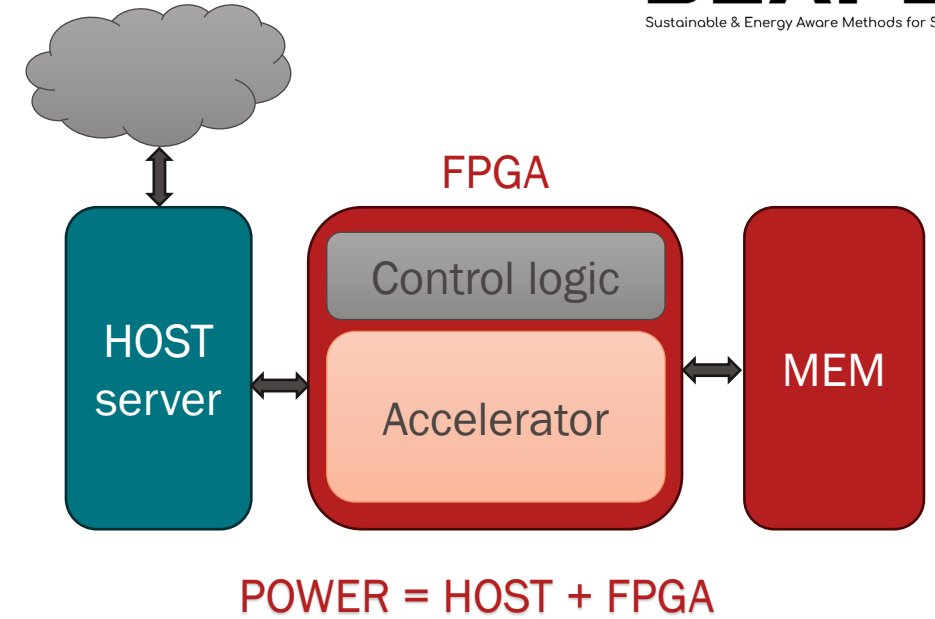
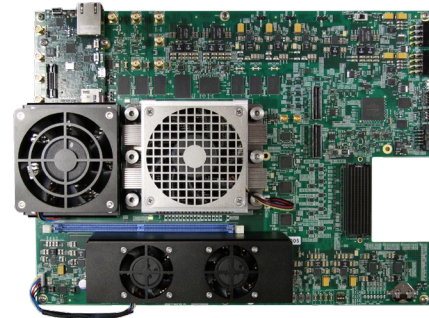
HLS Design Flow



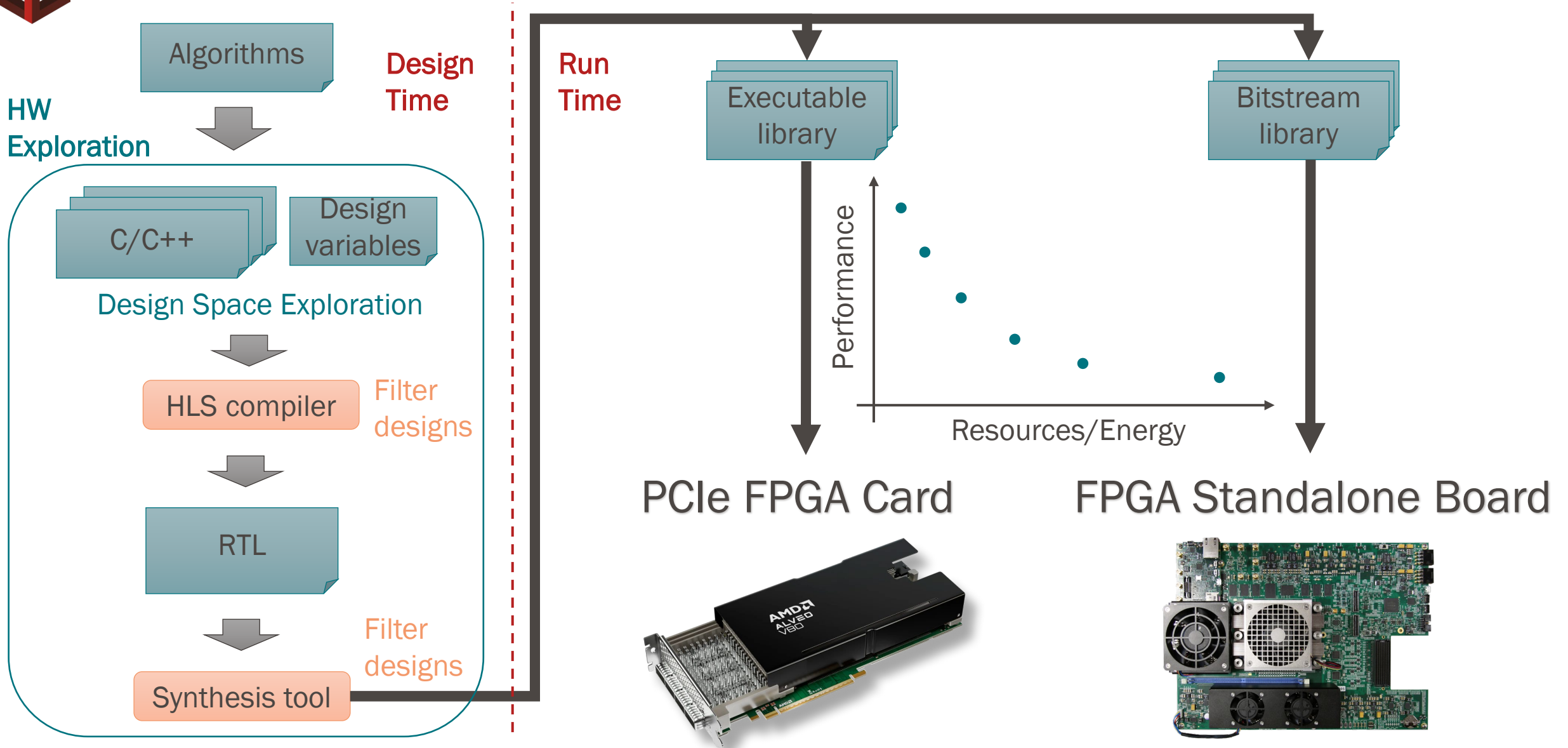
PCIe FPGA Card



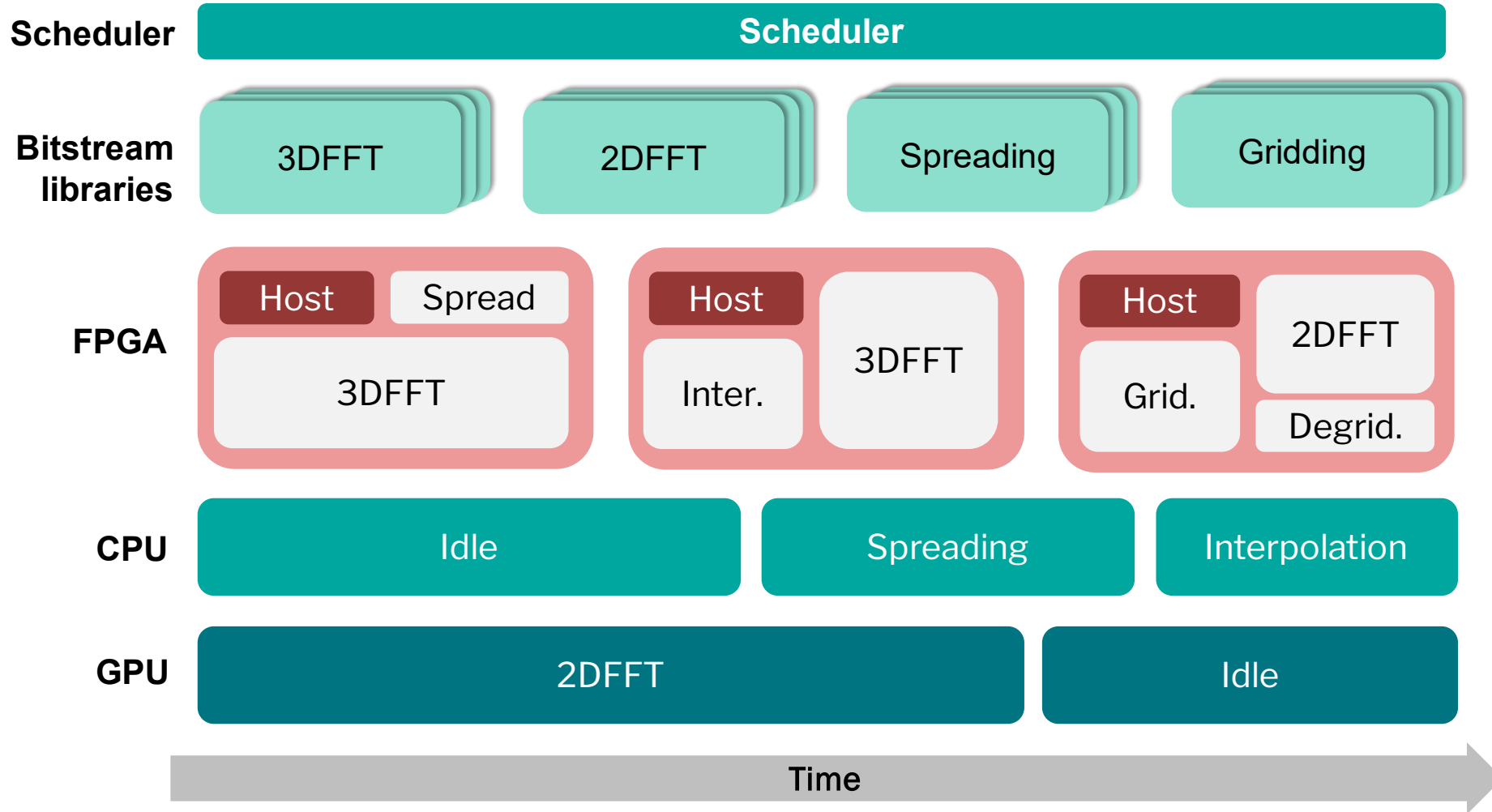
FPGA Standalone Board



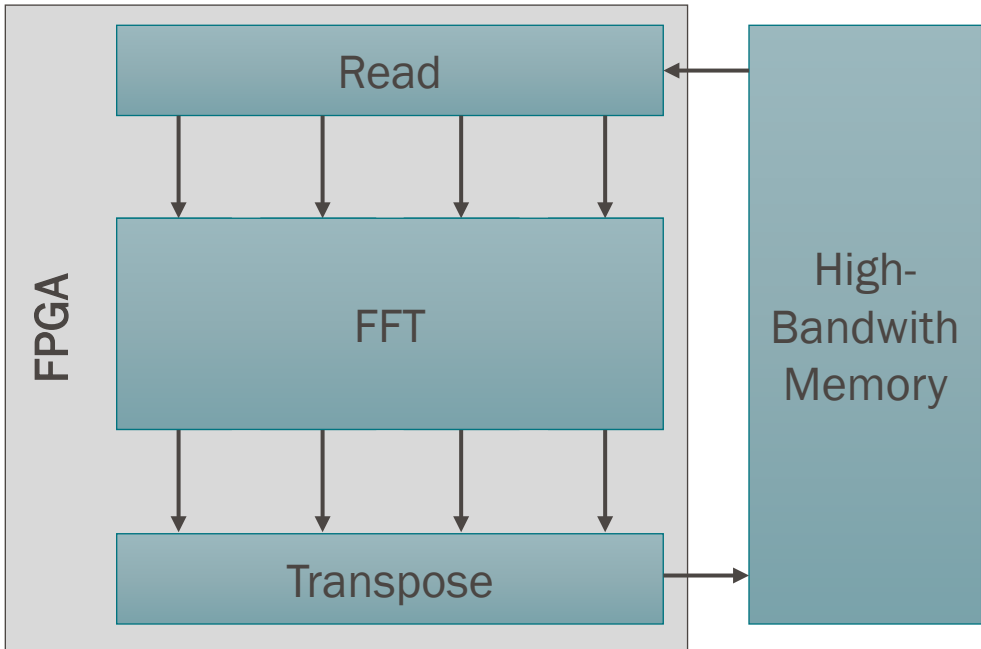
Design Space Exploration Design Flow



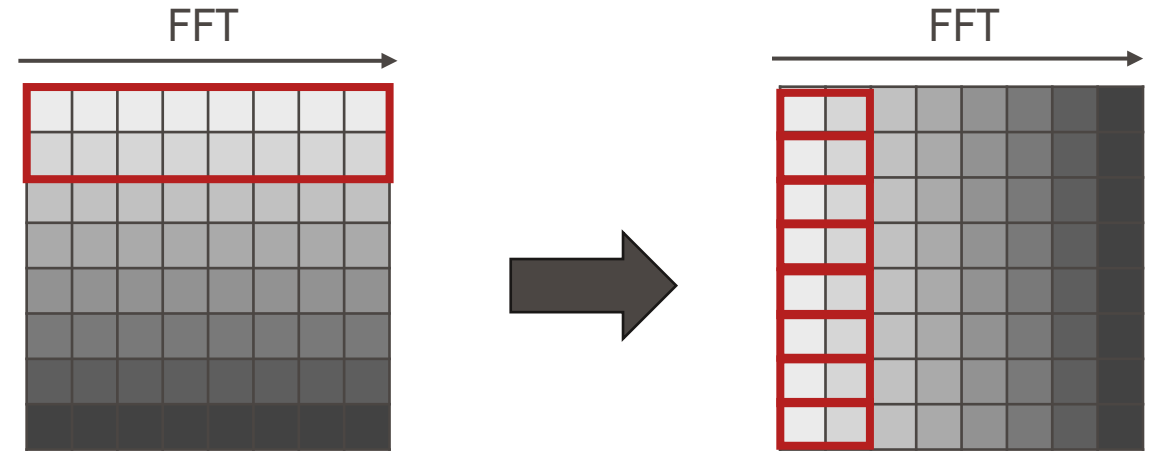
Run-time



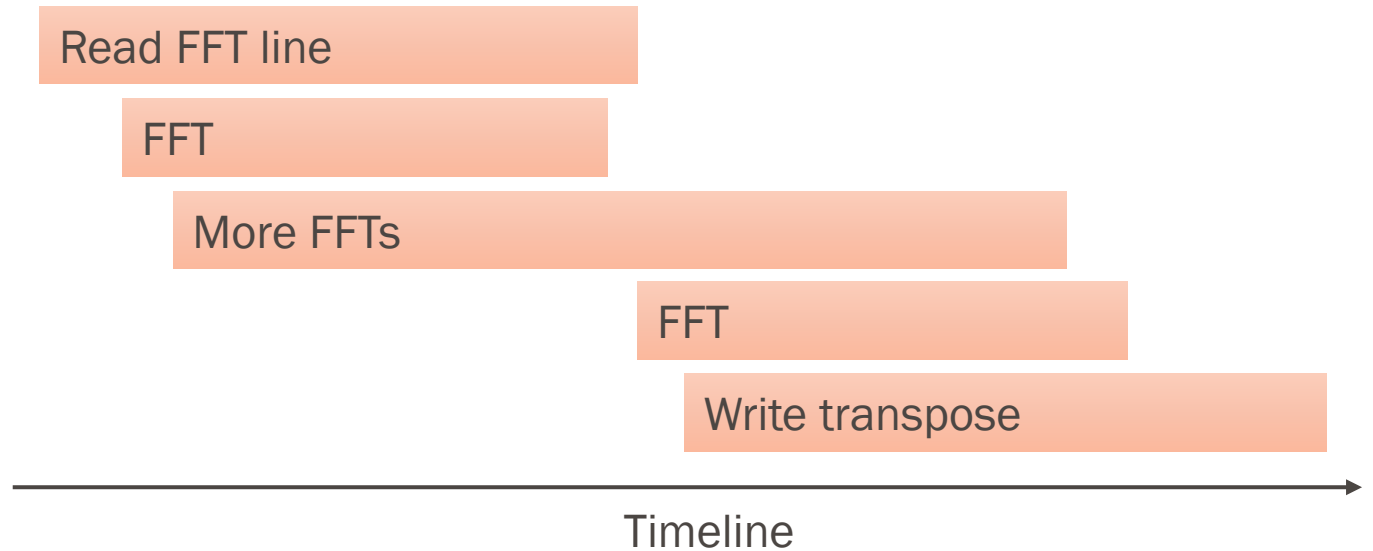
FFT HW Design and Exploration



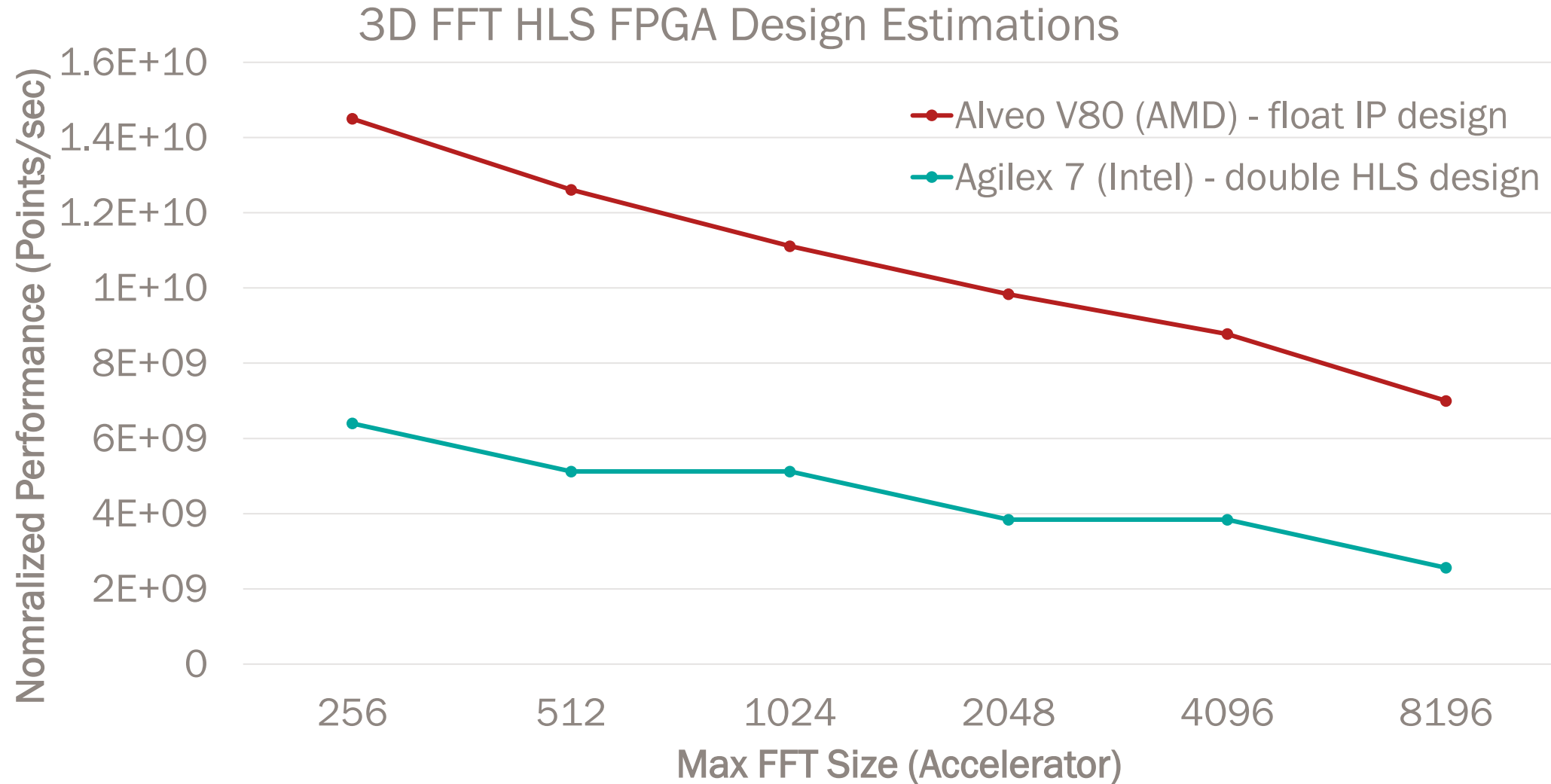
- Design Variables**
- # parallel FFTs
 - Transpose buffer
 - FFT stages
 - Data format
 - FFT Max FFT size



Consecutive transfers to memory takes less time and energy



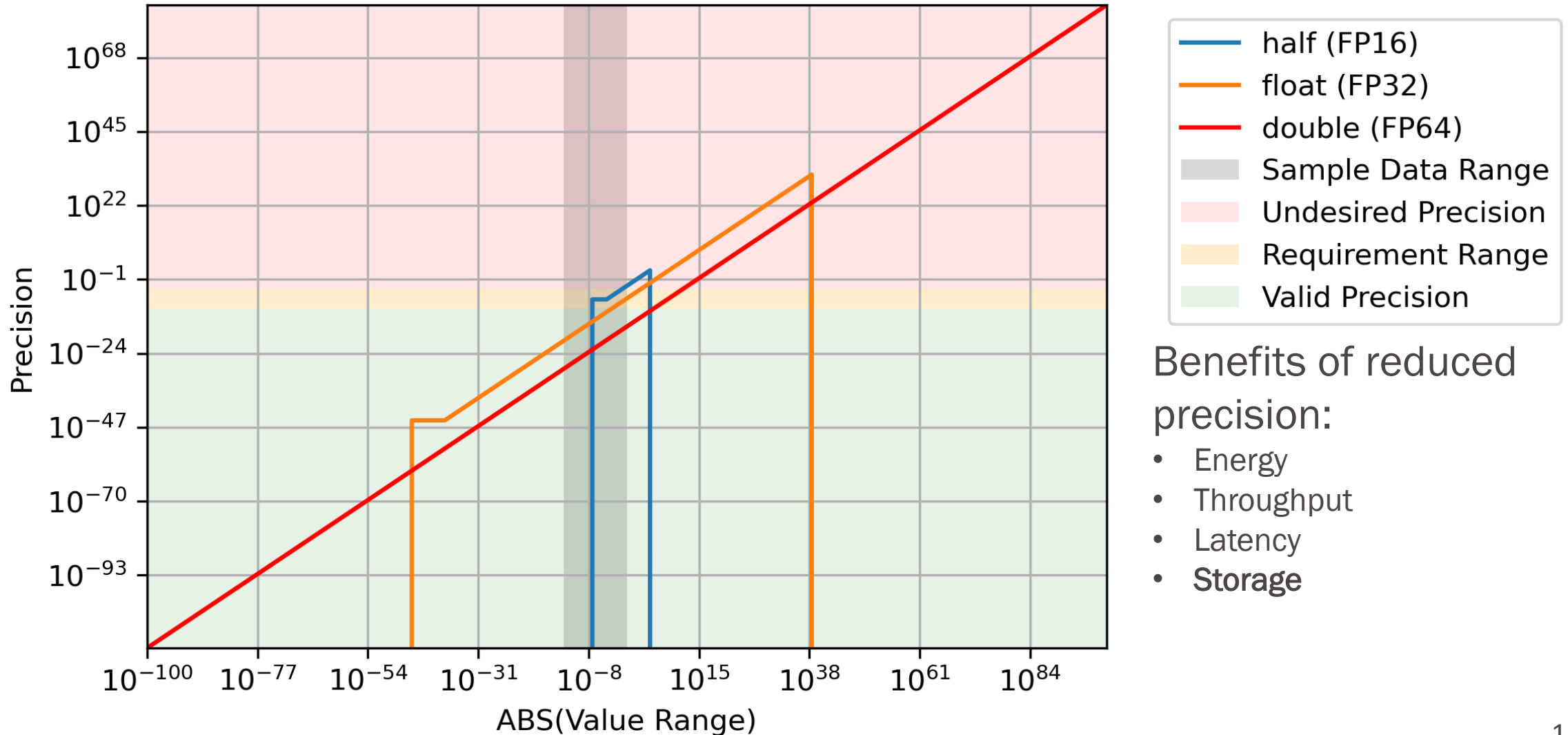
3D FFT HW Preliminary Results



Precision in for FINUFFT Synthesis (BIPP)

Sample data extracted from bipp execution, simulated with OSKAR for SKA-Low configuration

Precision ranges for different data representations



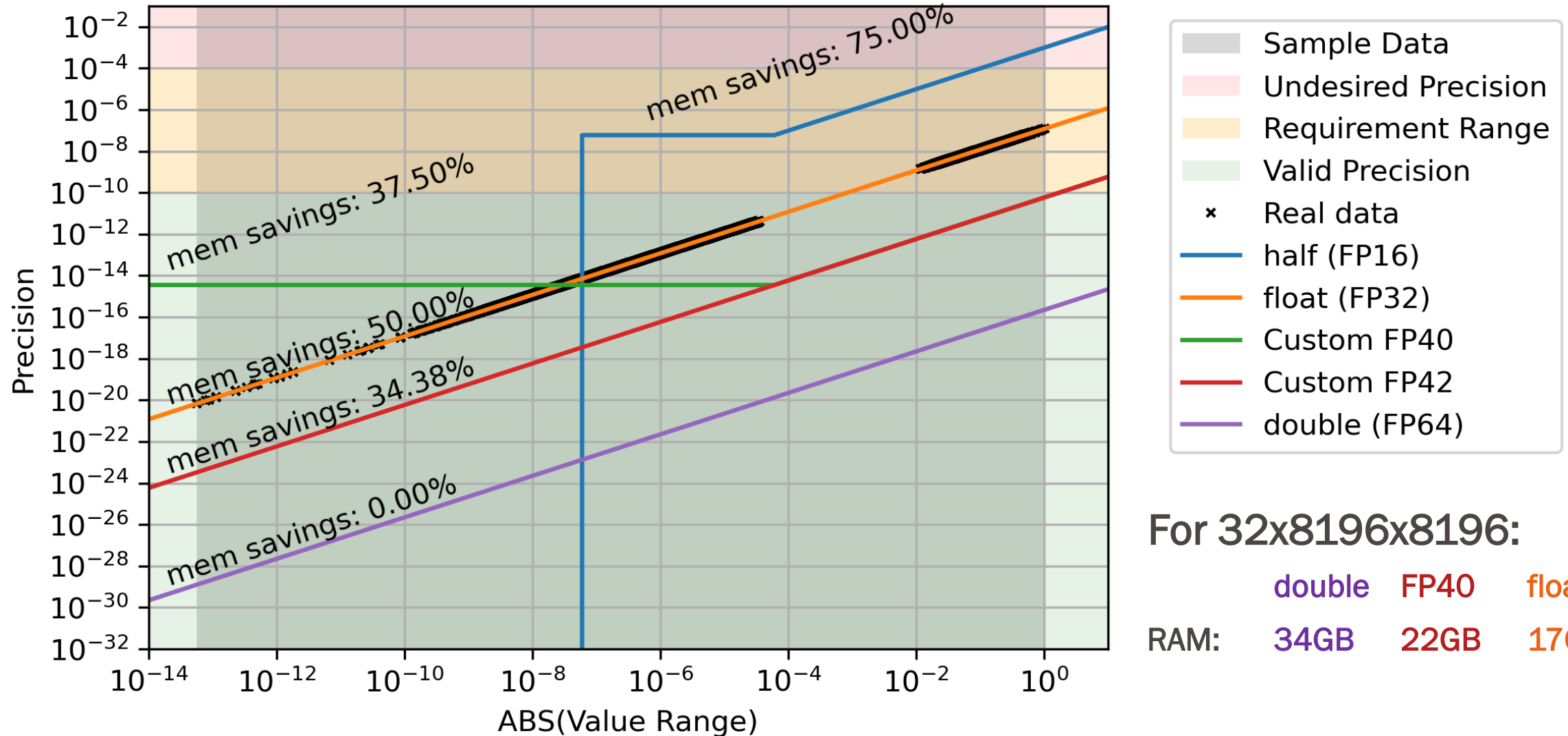
Benefits of reduced precision:

- Energy
- Throughput
- Latency
- **Storage**

Precision in for FINUFFT Synthesis (BIPP)

Sample data extracted from bipp execution, simulated with OSKAR for SKA-Low configuration

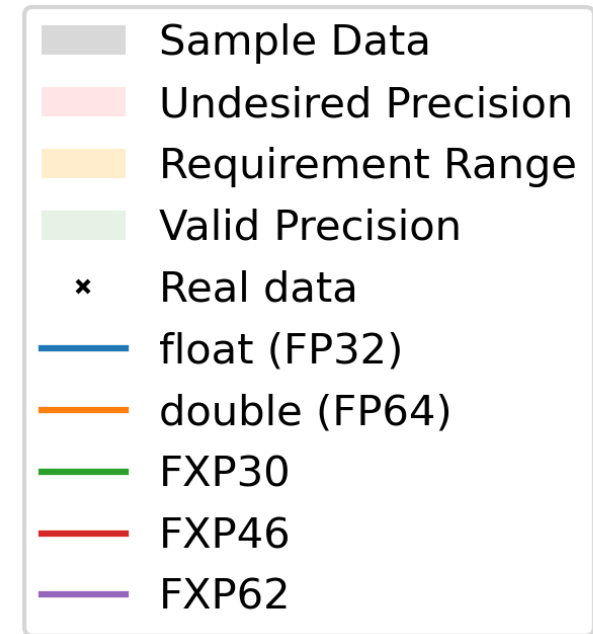
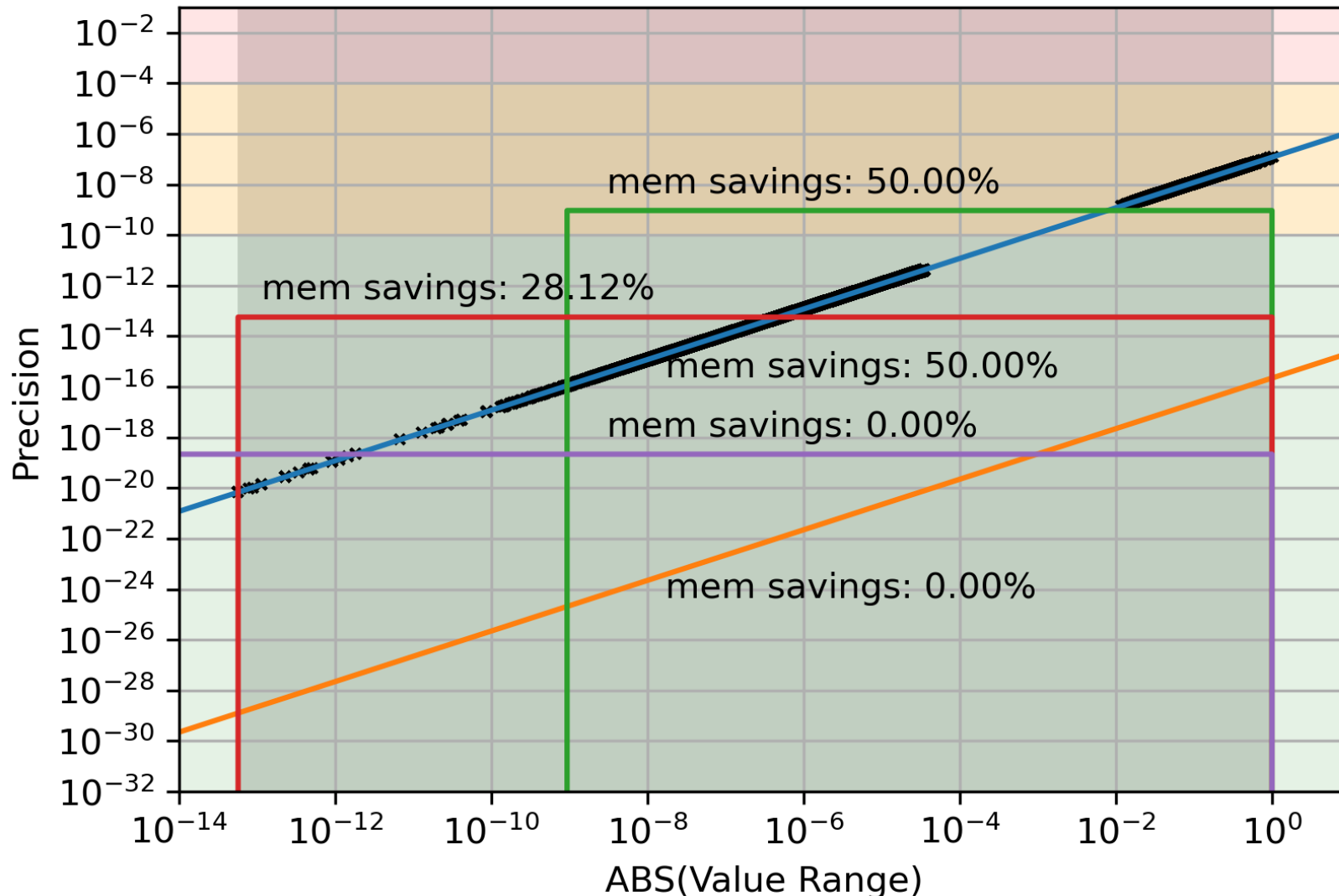
Precision of Floating-Point Formats



Precision in for FINUFFT Synthesis (BIPP)

Sample data extracted from bipp execution, simulated with OSKAR for SKA-Low configuration

Precision of Fixed-Point Formats



For 32x8196x8196:

	double	FP40	float
RAM:	34GB	24GB	17GB

Done

- Deploy flexible algorithms using FPGAs
- Accelerate kernels with an FPGA
- Explore the Design Space
- Share resources among different kernels

Ongoing Exploration

- FPGAs improve the energy consumption
- FPGAs match (even increase) the performance of GPUs
- Custom precision data formats are beneficial
- Solve memory-bounded workloads in FPGAs
- Reconfigure the FPGA at run-time for dynamic workloads

Inputs Needed

- Other algorithms to accelerate (i.e. ML)
- Dynamic range of real data (at different stages)
- Precision & latency requirements for different use cases
- Precision metrics (i.e. SNR)
- Scalability of algorithms



Thank you!

Ruben

EPFL - Embedded Systems Laboratory
ruben.rodriquezalvarez@epfl.ch
denisa.constantinescu@epfl.ch



Backup Slides



Types of NUFFTs algorithms

Method	Spread	FFT	Interpolation
NUFFT ₁	$N_{\text{vis}} \log \epsilon ^2$	$N_{\text{pix}} \log N_{\text{pix}}$	N_{pix}
W-gridding	$N_{\text{vis}} \log \epsilon ^3$	$N_{w'} N_{\text{pix}} \log N_{\text{pix}}$	$N_{w'} N_{\text{pix}}$
NUFFT ₃	$N_{\text{vis}} \log \epsilon ^3 + N_{\text{mesh}}$	$N_{\text{mesh}} \log N_{\text{mesh}}$	$N_{\text{pix}} \log \epsilon ^3 + N_{\text{pix}}$



Kashani, Sepand, et al. "HVOX: Scalable Interferometric Synthesis and Analysis of Spherical Sky Maps." *arXiv preprint arXiv:2306.06007* (2023).